

DISSERTATION

# **DISTRIBUTED VIDEO TRACKING**

by

SHARRANTH RAGU

16002514

Dissertation submitted in partial fulfilment of  
the requirements for the  
Bachelor of Engineering (Hons)  
(Computer)

JANUARY 2022

Universiti Teknologi PETRONAS  
Bandar Seri Iskandar  
31750 Tronoh  
Perak Darul Ridzuan



## CERTIFICATION OF APPROVAL

### **Distributed Video Tracking**

by

SHARRANTH RAGU

16002514

A project dissertation submitted to the  
Computer Engineering Programme  
Universiti Teknologi PETRONAS  
in partial fulfilment of the requirements for the  
BACHELOR OF ENGINEERING (Hons)  
(COMPUTER)

Approved by,



---

(Dr. Patrick Sebastian)

UNIVERSITI TEKNOLOGI PETRONAS

TRONOH, PERAK

January 2022

## CERTIFICATION OF ORIGINALITY

This is to certify that I am responsible for the work submitted in this project, that the original work is my own except as specified in the references and acknowledgements, and that the original work contained herein have not been undertaken or done by unspecified sources or persons.

A handwritten signature in black ink, appearing to read "Sharranth Ragu". It is written in a cursive style with a horizontal line underneath it.

SHARRANTH RAGU

## **ACKNOWLEDGEMENT**

For the past six months, I have been working on my final year project and dissertation, titled Distributed Video Tracking. Throughout the entire journey in completing this relatively challenging project, I have received immense support and assistance. With that being said, I would first take this opportunity to thank and express my extreme gratitude towards Dr. Patrick Sebastian, my final year project supervisor. My supervisor largely assisted in facilitating the overall project, providing ideas and solutions along the way.

Besides, I would like to appreciate Universiti Teknologi PETRONAS in providing such pragmatic and well organised syllabus, which made the progress of the whole project a smooth sail from day one. I believe that the final year project allows students to work on various skills, whilst honing technical understanding that may be taught in the years of pursuing degree, or knowledge that requires completely new self-learning.

I would want to specially thank my friends, that may have directly or indirectly helped me in this project. My parents and family were also my pillar of support throughout the project, both financially and emotionally.

## ABSTRACT

Distributed video tracking, or widely known as *Person Re-identification* is a process or the ability of a system to associate images or frames from videos of the same person, which could be sourced from both pre-recorded and real-time footages. The developed system should be able to identify and retrieve the person of interest by combining images from various non-overlapping cameras. The purpose of Re-ID is to specifically determine whether the person of interest has appeared in another location at a different time or by the same camera at a different instant. Hence, the approach in solving the issue is to first identify distinct characteristics that define a particular person, perhaps by feeding the *training* model with a dataset of cropped images of the non-stationary subject as there is movement across the camera [2]. The project can be divided into two main algorithms; the first algorithm is mainly catered for Pedestrian detection and tracking whilst the second focuses in reidentifying the pedestrian. The project would utilize OpenCV for capturing devices, TensorFlow object detection API (*Faster-RCNN*) for person detection (*forming bounding boxes*) and *SORT* (*Simple Online and Realtime Tracking*) for tracking the person during a video input . With prebuilt OpenCV functions, cropped images would be cropped from identifying the coordinates of bounding boxes, and saved in a local database. The stored images would then be the input material to be fed into the *Siamese Neural Network* developed through TensorFlow Keras. The Siamese Model was trained individually on three different camera inputs, where two of them were from pre-recorded videos taken from a phone camera, whilst one was taken from an internet protocol camera. With an input from the first angle, an individual would linger around for an unfixed amount of time. This is where the trained images would be collected, before being fed into the Siamese Neural Network for training. The model was then tested on another field of view, that was not overlapping, by allowing the same individual or a random person to walk across the frame. This was repeated similarly for all the different types of camera inputs as well. The overall performance of the neural networks were relatively decent, in fact excellent in some cases, based on the extensive testing done and challenges done. That being said, for all the three Siamese Model developed, the area under the curve in the receiver operating characteristic curve were very well in the accepted range, varying from the lowest of 0.64 and highest of 0.80. The average precision recorded in the precision recall curve ranged anywhere from 0.61 to 0.84, which comes to show that all the developed models have done excellently in predicting. Meanwhile, the time taken to train for the 10 epochs were approximately 190 to 310 seconds.

## TABLE OF CONTENTS

<b>CHAPTER 1: INTRODUCTION.....</b>	<b>1</b>
<b>1.1 Background of Study .....</b>	<b>1</b>
<b>1.2 Problem Statement.....</b>	<b>4</b>
<b>1.3 Objectives.....</b>	<b>5</b>
<b>1.4 Scope of Study .....</b>	<b>5</b>
<b>CHAPTER 2: LITERATURE REVIEW.....</b>	<b>6</b>
<b>2.2 Object Detection.....</b>	<b>6</b>
<b>2.2.2. <i>Fast R-CNN</i>.....</b>	<b>6</b>
<b>2.2.3. <i>Faster R-CNN</i> .....</b>	<b>7</b>
<b>2.2.4. <i>SSD Mobilenet</i>.....</b>	<b>7</b>
<b>2.2.5. <i>Mask R-CNN</i> .....</b>	<b>8</b>
<b>2.3 Object tracking.....</b>	<b>9</b>
<b>2.3.1. <i>Centroid Tracker</i> .....</b>	<b>9</b>
<b>2.3.2. <i>SORT (Kalman filters)</i> .....</b>	<b>10</b>
<b>2.3.3. <i>DCFnet</i> .....</b>	<b>11</b>
<b>2.4 Siamese Neural Network (One-shot Image Recognition).....</b>	<b>12</b>
<b>2.5 Non-Maximum Suppression (NMS) .....</b>	<b>14</b>
<b>CHAPTER 3: METHODOLOGY.....</b>	<b>15</b>
<b>3.1 Project Proposed Method.....</b>	<b>15</b>
<b>3.2 Project Methodology.....</b>	<b>16</b>
<i>i. Algorithm 1: Pedestrian Detection and Tracking .....</i>	<b>16</b>
<i>ii. Algorithm 2: Pedestrian Identification .....</i>	<b>17</b>
<b>3.2.1 Pedestrian Detection and Tracking .....</b>	<b>18</b>
<b>3.2.2 Pedestrian Identification .....</b>	<b>19</b>
<b>3.3 Software, Hardware and Tools.....</b>	<b>20</b>
<b>3.4 Gantt Chart .....</b>	<b>22</b>
<b>CHAPTER 4: RESULT &amp; DISCUSSION.....</b>	<b>24</b>
<b>4.1 Initial Phase .....</b>	<b>24</b>
<b>4.2 Person Detection and Tracking .....</b>	<b>26</b>
<b>CHAPTER 5: CONCLUSION AND RECOMMENDATIONS .....</b>	<b>48</b>
<b>5.1 Conclusion .....</b>	<b>48</b>
<b>5.2 Recommendation.....</b>	<b>49</b>
<b>CHAPTER 6: REFERENCES.....</b>	<b>50</b>

## LIST OF FIGURES

<b>FIGURE 1.1:</b> Projectile growth of world sales of Physical security products (2014-2024).....	1
<b>FIGURE 1.2:</b> Overview of object recognition tasks.....	2
<b>FIGURE 1.3:</b> Example of person re-id principle.....	3
<b>FIGURE 2.11:</b> Working principle of Fast R-CNN.....	6
<b>FIGURE 2.12:</b> Working principle of Faster R-CNN.....	7
<b>FIGURE 2.13:</b> Time Speed tests.....	7
<b>FIGURE 2.14:</b> Working principle of SSD MobileNet Framework.....	8
<b>FIGURE 2.15:</b> An example of the Mask-RCNN Framework.....	8
<b>FIGURE 2.16:</b> Assignment of IDs based on detected individuals.....	9
<b>FIGURE 2.17:</b> Example of the implementation of SORT .....	10
<b>FIGURE 2.18:</b> DCFnet architecture.....	11
<b>FIGURE 2.19:</b> Typical flow of person re-identification.....	12
<b>FIGURE 2.20:</b> Verification of tasks during training .....	13
<b>FIGURE 2.21:</b> One-shot Image Recognition .....	13
<b>FIGURE 2.22:</b> Working principle of the NMS .....	14
<b>FIGURE 2.23:</b> Removal of noise (additional bounding boxes) .....	14
<b>FIGURE 3.1:</b> Project Block Diagram .....	15
<b>FIGURE 3.2:</b> Pedestrian Detection and Tracking .....	16
<b>FIGURE 3.3:</b> Pedestrian Identification .....	17
<b>FIGURE 4.11:</b> Comparison between the different Network architectures.....	24
<b>FIGURE 4.12:</b> TensorFlow versions and their compatibilities .....	24
<b>FIGURE 4.13:</b> DLINK-DCS5030L configuration page .....	25
<b>FIGURE 4.14:</b> TP-LINK R1043ND router configuration page .....	25
<b>FIGURE 4.15:</b> Demonstrates the working principle of the algorithm .....	26
<b>FIGURE 4.16:</b> Video input from a pre-recorded video (1) .....	27
<b>FIGURE 4.17:</b> Video input from a pre-recorded video (2) .....	27
<b>FIGURE 4.18:</b> To test the Siamese Model on the video used for the training itself .....	28
<b>FIGURE 4.19:</b> Random video input (1) .....	29
<b>FIGURE 4.20:</b> Random video input (2) .....	29
<b>FIGURE 4.21:</b> A pre-recorded video, from a different angle (1) .....	30
<b>FIGURE 4.22:</b> A pre-recorded video, from a different angle (2) .....	30
<b>FIGURE 4.23:</b> Average fps on pre-recorded video input from two distinct angles .....	31

<b>FIGURE 4.24: ROC Curve and Precision Recall curves from a pre-recorded video (1) .....</b>	<b>33</b>
<b>FIGURE 4.25: ROC Curve and Precision Recall curves from a pre-recorded video (2) .....</b>	<b>33</b>
<b>FIGURE 4.26: Video input from live stream HTTP (1) .....</b>	<b>34</b>
<b>FIGURE 4.27: Video input from live stream HTTP (2) .....</b>	<b>34</b>
<b>FIGURE 4.28: Video input from live stream HTTP (3) .....</b>	<b>35</b>
<b>FIGURE 4.29: Video input from live stream HTTP (4) .....</b>	<b>35</b>
<b>FIGURE 4.30: Video input from live stream, for testing (1) .....</b>	<b>36</b>
<b>FIGURE 4.31: Video input from live stream, for testing (2) .....</b>	<b>36</b>
<b>FIGURE 4.32: Video input from live stream, from a different angle (1) .....</b>	<b>37</b>
<b>FIGURE 4.33: Video input from live stream, from a different angle (2) .....</b>	<b>37</b>
<b>FIGURE 4.34: Average fps on stream input from two distinct angles .....</b>	<b>38</b>
<b>FIGURE 4.35: ROC Curve and Precision Recall curves from a stream (1) .....</b>	<b>39</b>
<b>FIGURE 4.36: ROC Curve and Precision Recall curves from a stream (2) .....</b>	<b>39</b>
<b>FIGURE 4.37: Input from recorded video (1) .....</b>	<b>40</b>
<b>FIGURE 4.38: Input from recorded video (2) .....</b>	<b>40</b>
<b>FIGURE 4.39: Input from recorded video, from a different angle (1) .....</b>	<b>41</b>
<b>FIGURE 4.40: Input from recorded video, from a different angle (2) .....</b>	<b>41</b>
<b>FIGURE 4.41: Input from recorded video, from a different angle (3) .....</b>	<b>42</b>
<b>FIGURE 4.42: Input from recorded video, with a different person (1) .....</b>	<b>42</b>
<b>FIGURE 4.43: Input from recorded video, with a different person (2) .....</b>	<b>43</b>
<b>FIGURE 4.44: Input from recorded video, with a different lighting (1) .....</b>	<b>43</b>
<b>FIGURE 4.45: Input from recorded video, with a different lighting (2) .....</b>	<b>44</b>
<b>FIGURE 4.46: Input from recorded video, with a different lighting (3) .....</b>	<b>44</b>
<b>FIGURE 4.47: Average fps on pre-recorded video inputs from two distinct angles .....</b>	<b>45</b>
<b>FIGURE 4.48: ROC Curve and Precision Recall curves from the pre-recorded videos (1)....</b>	<b>46</b>
<b>FIGURE 4.49: ROC Curve and Precision Recall curves from the pre-recorded videos (2)....</b>	<b>46</b>
<b>FIGURE 4.50: Average time taken to run 10 epochs .....</b>	<b>47</b>

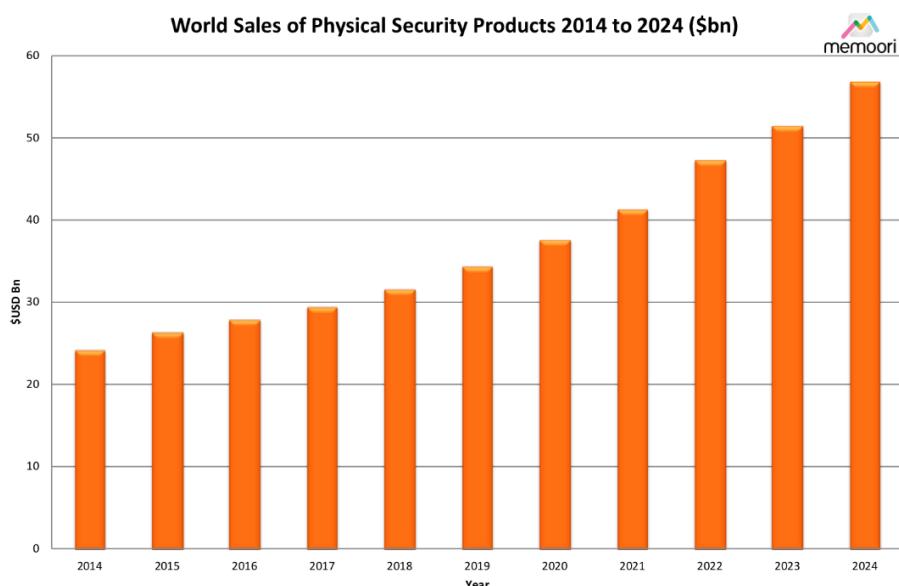
## LIST OF TABLE

<b>TABLE 1: Software, Hardware and Tools.....</b>	<b>20</b>
---	-----------

# CHAPTER 1: INTRODUCTION

## 1.1 Background of Study

Distributed video tracking or Person Re-identification has seen rising interest in this rapidly growing global market, especially in terms of surveillance monitoring. With the increased use of machine learning techniques in wide area of technological fields, it has also seen some heightened potential benefits within the safety and security branches. Monitoring can be defined as a systematic, incessant, active, or passive observation of places, objects, processes, or individuals [3]. Contrarily, surveillance emphasises on individuals, buildings or even vehicles that is targeted or in cases where the subject detected may be deemed as suspicious, based on reliable data.



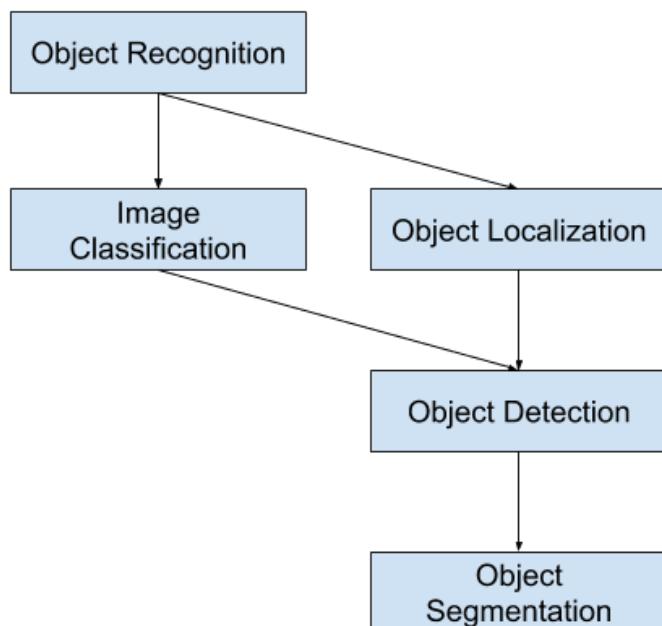
**FIGURE 1.1 – Projectile growth of world sales of Physical security products (2014-2024)**

Physical security has been expanding exceptionally, particularly in the last 10 years as the market has surged to an annual compound growth rate of 6.27 percent. The figure 1.1 above projects that the market would hover around a value of \$56.76 billion by 2024 with a compound annual growth rate of an excellent 10.72 percent [4]. With the increasing demand of surveillance monitoring, it is significant to develop more advanced and complex solution which should be more accurate, complex, and cost-effective, thus reducing the requirement of hiring actual supervisors or guards for the jobs that AI could simply handle.

The relatively new implementation of this AI monitoring technique is a widely used technological advancement in liberal countries such as Switzerland, Norway, and Sweden. According to a research conducted [5], artificial intelligence surveillance systems are used by around 51% of advanced democracies. From safe city platforms to facial recognition cameras, democratic governments or countries are employing a myriad of monitoring technology, which does not automatically imply that democracies abuse their systems' rights.

Fundamentally, computer vision tasks can be divided into three branches; Image Classification, Object Localization and Object Detection [6]. For instance, in image classification, the type of classes of objects can be determined and identified with single object inputs such as photographs which, in return, would provide an output of class labels.

The existence of objects can be detected in an image, and the bounding box would be used to determine the position in object localization. The inputs for such cases include photographs with more than a single object whilst bounding boxes with defined heights and widths are the delivered outputs. Similar to object localization, object detection instead would be able to detect each of the bounding boxes' class labels [6].

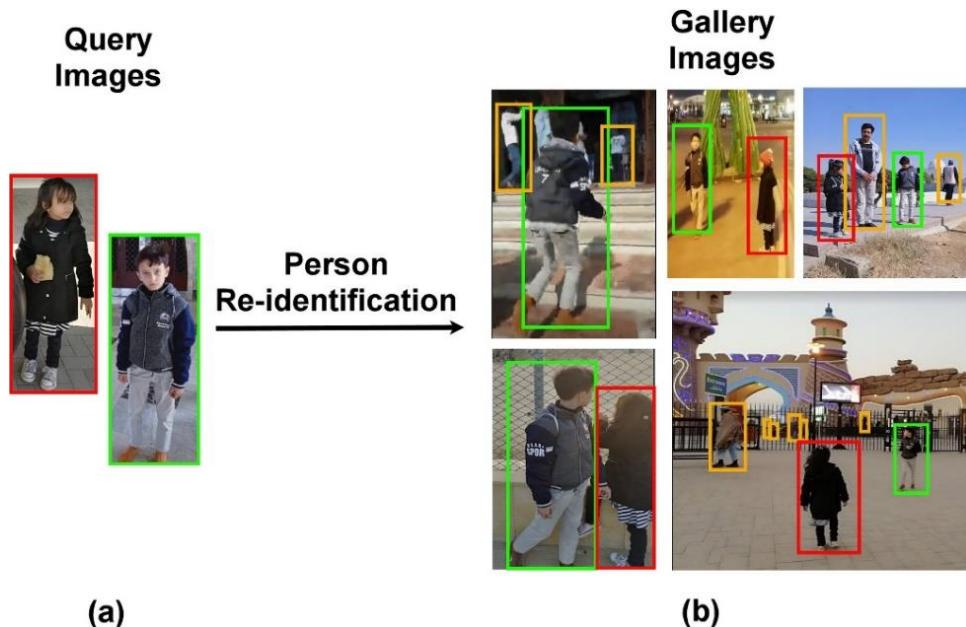


**FIGURE 1.2 – Overview of object recognition tasks**

The figure 1.2 represents the overview of object recognition tasks which can further be broken down to several more subsections. In image classification, algorithms would generate a list of all the object categories identified in the image or frame, while in single object localization provides an additional information of the bounding boxes denoting position and scale of instances in each categories of objects.

The Person Re-identification framework provides a variety of advantages, for instance the most common scenario in which it could be utilised is for video surveillance in security assurance like areas around a shopping mall, parking lot, university campus, or scenarios that involve several non-overlapping cameras. Undeniably, there has been extensive research that has been conducted, that primely focuses in studying the numerous methods of re-identification techniques. Re-ID, this relatively new technology has the potential to observe high growth demand in the upcoming years, due to its extensive beneficial use cases.

However, some of the contradicting issues or problems in distributed video tracking include occlusions, poor lighting conditions, understanding the implementation of colour spaces between different cameras used and appearance variations [7],[8]. Besides, the several methods that exists all provide various sets of percentage performance. Therefore, practicing a suitable re-ID technique is vital in maintaining the accuracy, efficiency, and future deployment of such technology, in the real world.



**FIGURE 1.3 – Example of person re-id principle**

## **1.2 Problem Statement**

The Person Re-identification and detection framework in video surveillance allows the tracking of individuals from security cameras; internet protocol cameras, closed circuit television (CCTV). Distributed video tracking, involves challenges that require the retrieval of a specific person from several non-overlapping, discontinuous cameras; attempts to establish whether the person being tracked has appeared in another location or not. This technology, with the aid of artificial intelligence and deep learning mechanism would provide a wide range of advantages that can be utilised in various fields.

Currently, the majority of surveillance monitoring that exists, especially the conventional ones that are out there in the market are limited, due to the dependency of actual human being needed to operate such systems. This caveat that persists is an important factor, because having hired individuals to do such jobs are not necessarily cost-effective, and may lead to unexpected human error. However, with the assistance of distributed video tracking, higher product efficiency may be achieved.

For example, security surveillance monitoring need guards to actually observe whether or not a particular person has appeared multiple times on different, non-overlapping cameras. Hence, guards must closely observe for any suspicious behaviour continuously or view long hours of security footages in case something goes wrong in the future. This would lead to risks or issues such as false identification, delayed recovery, and extended periods of working hours. Some of these are critical scenarios, as it may require a higher budget or time which is not favourable at all times.

Instead, with distributed video tracking and the framework for pedestrian detection would provide smoother yet safer security surveillance options. Since all the information gathered from the cameras are labelled and classified, it would ultimately be saved to a database; can be used to study future research's that may be conducted from the collected data. A well-made and optimised framework, could even eliminate the need of constant physical monitoring, as everything has been done automatically and simultaneously.

### **1.3 Objectives**

The primary objectives for the project are:

- To understand the various methods for capturing parameters for the person of interest.
- To develop algorithms that would be able to detect, track and re-identify the person of interest from non-overlapping cameras.
- To verify and evaluate an existing network in a real-world scenario and on a new data.

### **1.4 Scope of Study**

The project would have an emphasis in understanding the different methods or procedures for capturing parameters on the person of interest. Some of the existing framework for person re-identification that exists are *Comparative Attention Networks (CAN)*, *Part-Aligned Representation (PAR)*, *Deep Contents Aware Features (DCAF)*, *PersonNet*, *Learning Discriminative Neural Space (LDNS)* and *Gated Siamese*. The differences and benefits of picking these different methods would be explained further in this report. Besides, this project would cover in understanding the types of different trackers; CNN-based and Correlation filters. This project, however, utilizes the *Faster R-CNN* architecture for person detection, which is simply based on the initial *R-CNN* with major improvements on efficiency and accuracy.

For image processing and the capturing of reference images, OpenCV would be used. It is a programming functions package that focuses mostly on real-time computer vision tasks. TensorFlow would be used for the building the proposed model for the project, as it can be used in a variety of applications, although its focus is primarily on deep neural network training and inference. Both of the machine learning frameworks would be coded completely in Python (*Visual Studio Code as the opted IDE*). Keras, a python-based deep learning API that runs on top of TensorFlow, would allow for the development of more complex deep learning models, and in the case of this project, the Siamese Neural Network.

## CHAPTER 2: LITERATURE REVIEW

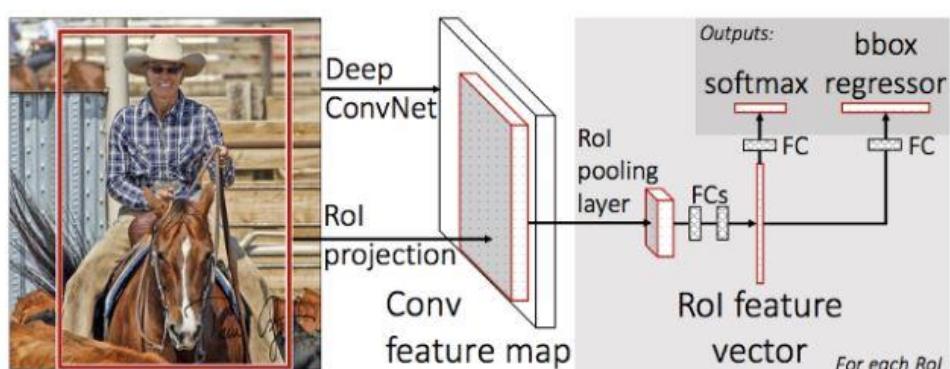
### 2.2 Object Detection

With the aid of Convolutional Neural Networks architectures, detecting various objects and segregation according to their classes becomes possible. For example, in a surveillance monitoring system, the developed algorithm should be able to detect individuals and assign respective IDs to them, or a set a count to visualise the number of people are there at any particular instance. Fast R-CNN, Faster R-CNN [19], and Mask R-CNN are more recent architectures that are based on their predecessor, CNN.

#### 2.2.2. *Fast R-CNN*

Instead of a pipeline, Fast R-CNN is presented as a model that learns and outputs regions and classes instantaneously [21]. Typically, the model would take the image as inputs and generates a collection of area suggestions; processed by a deep CNN network. A distinct layer, Region of Interest Pooling Layer, is added at the conclusion of the deep CNN to extract characteristics relevant to their respective inputs [21]. Contrary to its predecessor, Fast R-CNN inputs the entire image into the Deep ConvNet to obtain the output, Feature MAP [22].

Similarly, *Selective Search* will obtain the bounding boxes for the input images. However, the bounding box will be projected proportionally to the Feature MAP [22], also known as *region of interest (Roi) projection*. Since the fully connected network layer does not work in different sizes, the Roi pooling layer would resize the Roi projection, consecutively classify and predict the bounding boxes.

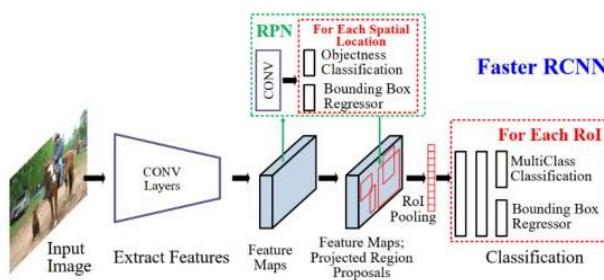


**FIGURE 2.11 – Working principle of Fast R-CNN**

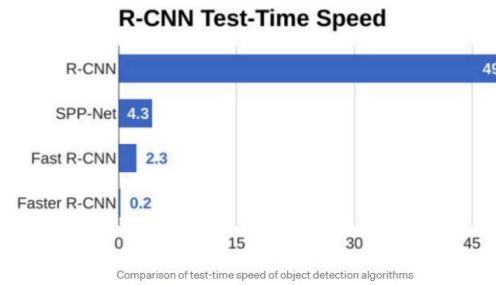
### 2.2.3. Faster R-CNN

The Faster R-CNN this technique uses the Region Proposal Network (RPN), which is way more cost-effective and computationally cheaper than both its predecessors, R-CNN and Fast R-CNN [21]. The input image would be parsed into the convolutional layer at once and sent to the Feature MAP, similar to Fast R-CNN. Instead of *Selective Search*, RPN would be used for each spatial location [22], [23] for object classification and regressor for bounding boxes. Rather than discovering the region proposals and then passing thru the CNN network for features extraction [13], it merely requires a single time to extract features for the entire region.

Hence, the Feature MAP would then be used to predict the region from the input image. For the purpose of this project, this algorithm would be utilized due to its extremely high efficiency and response time. Despite producing way more accurate results compared to its competitors, Faster R-CNN tend to be computationally intensive and requires high CPU and Graphics processing power usage [29],[30].



**FIGURE 2.12 – Working principle of Faster R-CNN**

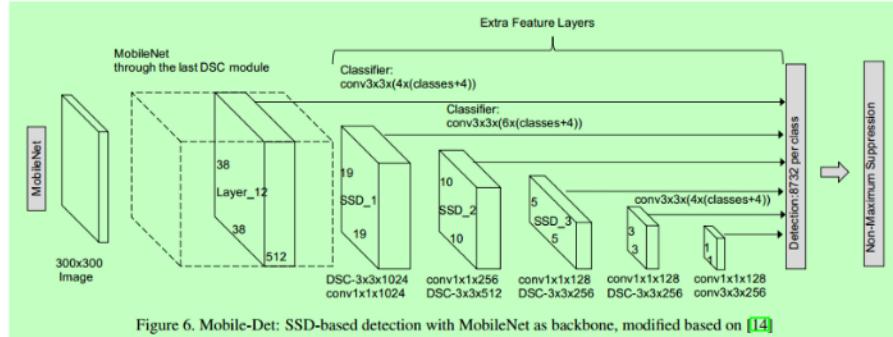


**FIGURE 2.13 – Time Speed tests**

### 2.2.4. SSD Mobilenet

The Caffe framework was used to implement the SSD MobileNet, which was pre-trained with the PASCAL VOC 2007 dataset [25]. The SSD stands for Single Shot Detector architecture. Since it is pretty accurate, quick, compact, lightweight and has the simplicity to adjust or tune, this framework was excellent for real-time use [26]. However, SSD MobileNet is accompanied with some very obvious limitations; has the ability to only identify and detect 20 different classes. Besides, object categories with relatively smaller size would be harder to detect, because the performance of SSD MobileNet is solely dependent on object dimensions [25],[26].

The SSD framework utilizes a system known as the feed-forward convolutional network to generate a fixed size set of bounding boxes keep tracks of the existence of object classes [30]. Non maximum suppression would then occur which removes the overlapping bounding boxes and reduces the overall noise, before producing the final results [30].



**FIGURE 2.14 – Working principle of SSD MobileNet Framework**

## 2.2.5. Mask R-CNN

The Mask R-CNN is an extension of the Faster R-CNN framework, but differs by having the addition of a branch that would be used in the prediction of masking a particular object [28]. Now since the architecture is able to generate masks in pixel-level, it would be an easy task to perform other more complicated jobs, like predicting various human postures [29]. On a myriad of tasks, Mask R-CNN can outperform all single-model competitors that exists as of today [29]. However, being able to execute such intricate challenges comes with a price; tend to be computationally intensive and requires longer processing time [29].



**FIGURE 2.15 – An example of the Mask-RCNN Framework**

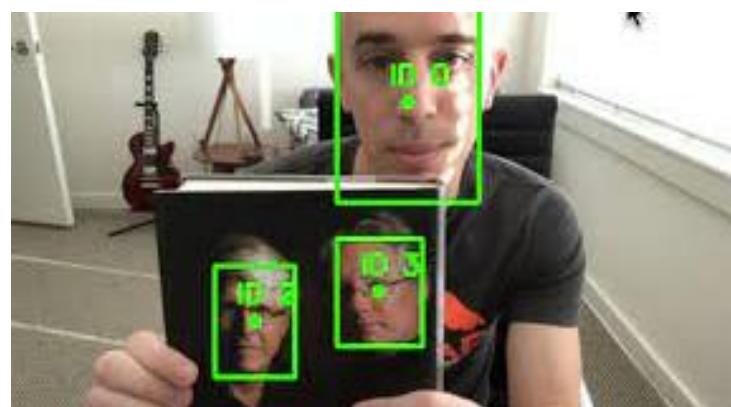
## 2.3 Object tracking

This deep learning application utilizes a set of tools to collect the initial objects' observation which allows the identification and creation of different classes of objects as they move across frames of a particular video. It would then assign respective IDs to each of the identified items. In a more define manner, object tracking accurately detects objects in a video and understand them as a collection of trajectories [16]. Object tracking in machine learning can be divided into two different categories [13]; CNN-based and Correlation filters. Despite being only slightly better than the Correlation filters, CNN-based approaches tend to take longer time and requires prior training. Contrarily, DCF-CSR (Discriminative Correlation Filter with Channel and Spatial Reliability) [18], provides much accurate results whilst being a durable tracker that is efficient enough to run in real-time.

### 2.3.1. Centroid Tracker

This tracking algorithm is highly effective; the way it works is, it depends on the Euclidean distance among existing object centroids and new object nodes [31]. It does this by identifying the centroid of the bounding boxes to which new IDs are assigned [31],[32]. Prior to these steps, object detection should have been performed on each of the corresponding frames. Between each of these frames, the centroids must be near together for it to work efficiently. Due to its extreme simplicity, centroid tracker would not be able to perform well while tracking a number of objects in real-time.

Centroid tracking merely consists of a few steps; head-starting with a series of object detections, and for each of the respective first detections, unique IDs would be created. Hereafter, the respective IDs should be maintained as the objects move across frames, without assigning a new ID for the particular object, based on the limitations set.



**FIGURE 2.16 – Assignment of IDs based on detected individuals**

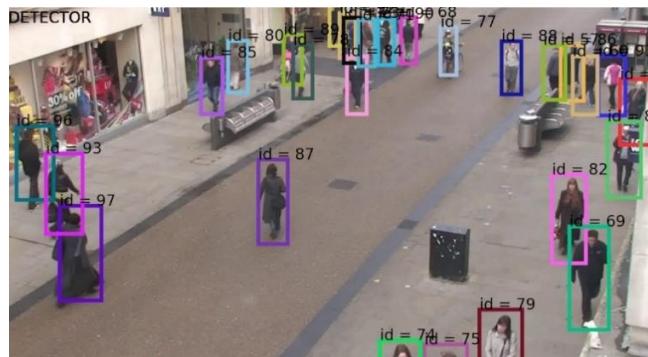
Centroid-based object tracking can be divided into four main steps that allow the assignment of IDs, as it requires multi-step procedures;

- Determine centroids by identifying bounding box coordinates.
- Calculates the distance between the newly formed bounding boxes and the already ones in Euclidean space.
- The coordinates of the existing objects would constantly be updated.
- The registration of the new objects (assigning IDs)

### 2.3.2. *SORT (Kalman filters)*

SORT (Simple Online and Realtime Tracker) is a framework that is based off Kalman filters [33],[43]. It is an online tracking algorithm which make predictions regarding the current frame, and it takes into account information from the current and prior frames [34],[35]. Associating items in two adjacent frames provides the rectification information; by using attributes and functions like IOU (intersection over union). It's intended for Realtime tracking using only previous and current frames, and instantaneously generate object identities. Each bounding box that represents an item is allocated a unique track ID; confidence value greater than a certain threshold [35].

SORT would not be able issues like occlusion and re-entering of objects [42], resulting in the development of a more complicated model that utilizes deep learning, known as DeepSort [44],[45]. However, unlike SORT, DeepSort requires high CPU and GPU usage. Hence, the hardware necessary to perform a more complicated technique at a higher frame rate is prohibitively expensive [35]. For the purpose of this project, SORT has the potential to get object tracking sorted out due to its light-weighted nature, hence easier integration with a more hardware intensive object detection module like Faster R-CNN [36].

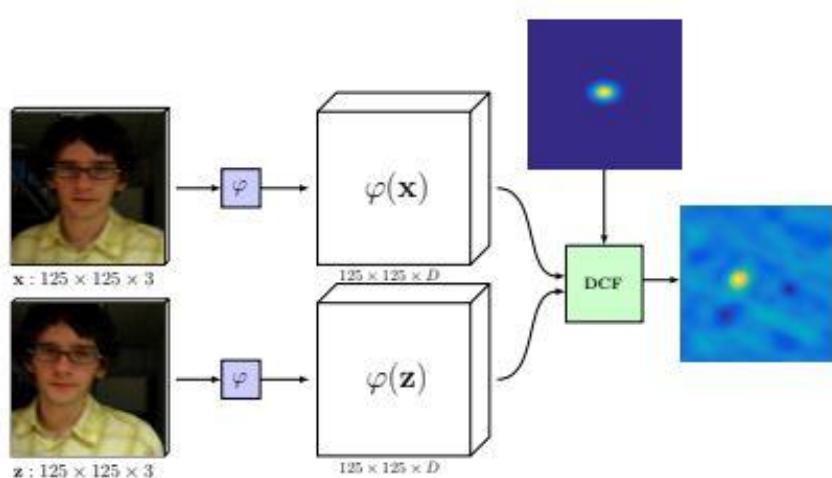


**FIGURE 2.17 – Example of the implementation of SORT**

### 2.3.3. DCFnet

DCFnet is a network architecture that is light in weight, allowing the tracker to run at a high speed of above 60 frames per seconds theoretically. Besides being highly compact, dense and quick. The term DCF is an acronym for Discriminant Correlation Filters, introduced as a particular correlation filter in a Siamese Network [36]. The problem with this deep learning online algorithm is, that much more conceptually sophisticated and complicated than its competitors of the same range [36].

The Discriminant Correlation Filter (DCF) is then derived backpropagation via it by characterising the network output as a likelihood heatmap of a particular object position. The effectiveness attribute of DCF is intact since the derivation is still done in the Fourier frequency domain.

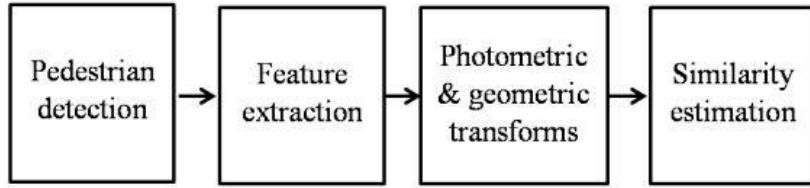


**FIGURE 2.18 – DCFnet architecture**

## 2.4 Siamese Neural Network (One-shot Image Recognition)

The Siamese Neural Network, allows the build of models that are not particularly achievable with the use of pure classification. This deep learning neural network is made up of two subnetworks that each produce two embeddings. Those extracted features are then sent into a loss function as inputs. In terms of functionality, the Siamese Model is barely distinct than most convolutional neural networks that exists; inputs images, and converts their characteristics into a numerical representation. The distinction is in the post-processing of the result [41].

Hence, from the characteristics and features collected in the first step, the metric learning stage would then determine the similarities of the two images. If the distance is smaller between these two compared images [13], it can be predicted that the chance of detecting the same person is higher. In [14], a conventional person re-identification framework however, would first include pedestrian detection, followed by extracting the features of individuals, some level of photometric and geometric transforms, and to increase the accuracy, similarity estimation would be concluded, as seen in the figure below.



**FIGURE 2.19** – Typical flow of person re-identification [14]

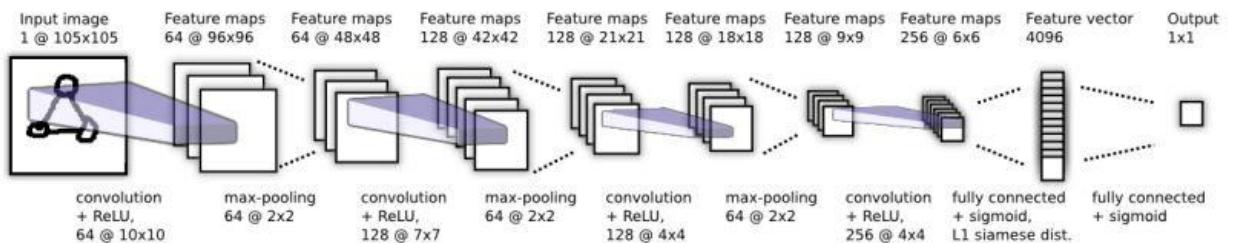
For further optimising coverage areas, surveillance systems in public places typically employ networks of cameras. However, due to financial and infrastructure constraints, these cameras frequently have limited or no functionality; no overlap within the area of view, hence identifying individuals across various other cameras is a hefty process [9]. Some of the pre-existing frameworks developed for Person Re-identification can generally be divided into two contrary subcategories; traditional and deep learning methods. Hence, this is where a model like the Siamese neural network comes in handy.

The Siamese model can be summarised as follows [40]; create a model that is able to distinguish between a set of similar and dissimilar pairings of images, and for verification, generalise to identify new classifications based on learned feature mappings. This deep learning approach, captures attributes to transformations in the input vector by employing a framework with a set of parameters and layers of non-linearities. The figure below represent the following statements, where the model should be able to identify the closest image related to, from an input image.



**FIGURE 2.20 – Verification of tasks during training [40]**

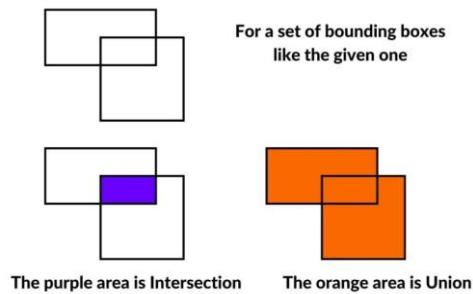
The framework is made up of a series of fully connected layers, each of which employs a single channel with different-sized filters and a constant stride of one [40]. Rectified linear activation function (ReLU) is embedded in the neural network. It is a nonlinear linear function that would output the input if the detection were positive, else it outputs zero [42]. The framework proposed in [40] suggests the input image to undergo convolution, ReLU and max-pooling for three consecutive times, before fully connecting the feature maps with sigmoid and L1 Siamese distance.



**FIGURE 2.21 – One-shot Image Recognition [40]**

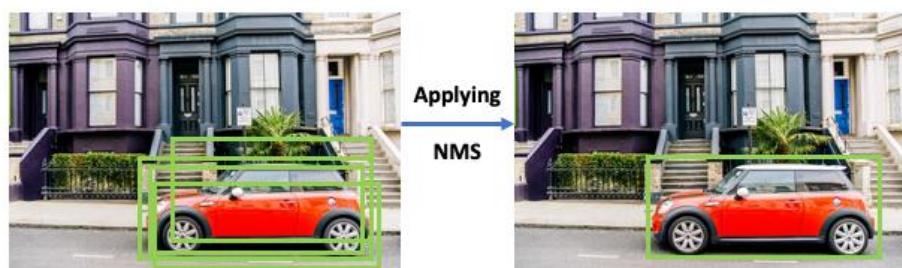
## 2.5 Non-Maximum Suppression (NMS)

NMS (Non-Maximum Suppression) is a key tool in a myriad of object detection algorithms, which drastically increases the accuracy of object localization. In object detection output, there are a significant number of bounding box candidates, and these redundant bounding boxes impair detection accuracy [37]. NMS is used by most object detection algorithms to reduce a large number of detected bounding boxes to only a few. Windowing is used by most object detectors at the most basic level. Hundreds of thousands of windows of varying sizes and shapes are created. These windows are said to contain only one object, and each class is assigned a probability/score by a classifier. After the detector generates a huge number of bounding boxes, the best ones must be filtered away. Hence, the most often used algorithm for this task is NMS.



**FIGURE 2.22 – Working principle of the NMS**

NMS is often employed at the end of the process to reduce the number of detections and establish the bounding box with the highest confidence level [38]. The Jaccard index, also known as the Intersection over Union (IoU) metric, is a method for quantifying the percentage overlap between both the ground truth Bounding Box and the predicted Bounding boxes [39]. With the following steps, there should be only a single bounding box that appears, which is also the most accurate one.

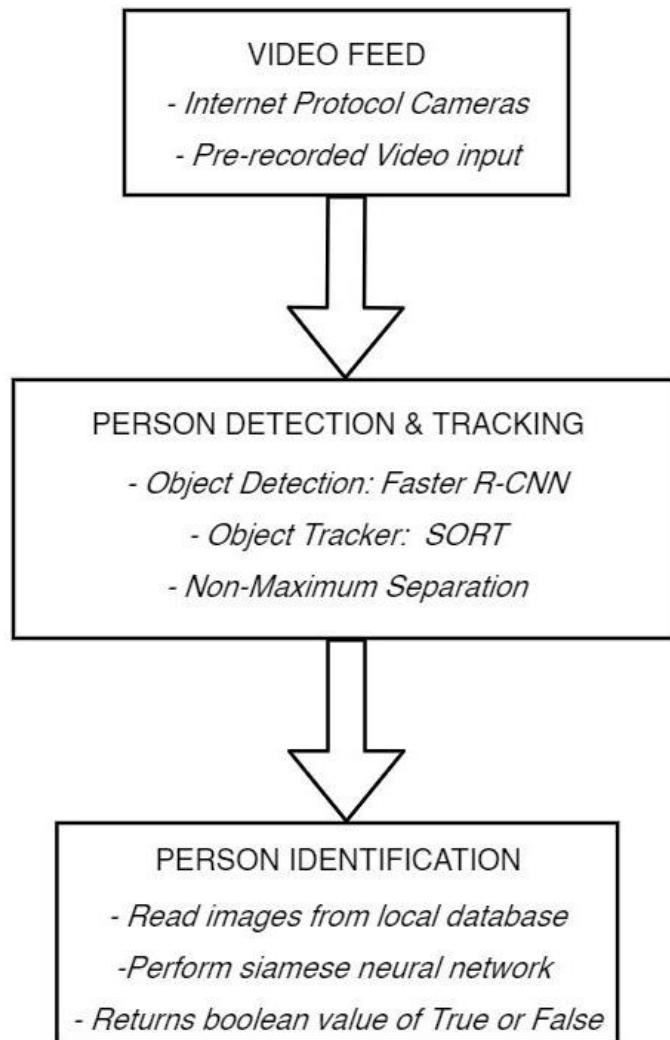


**FIGURE 2.23 – Removal of noise (additional bounding boxes)**

## CHAPTER 3: METHODOLOGY

### 3.1 Project Proposed Method

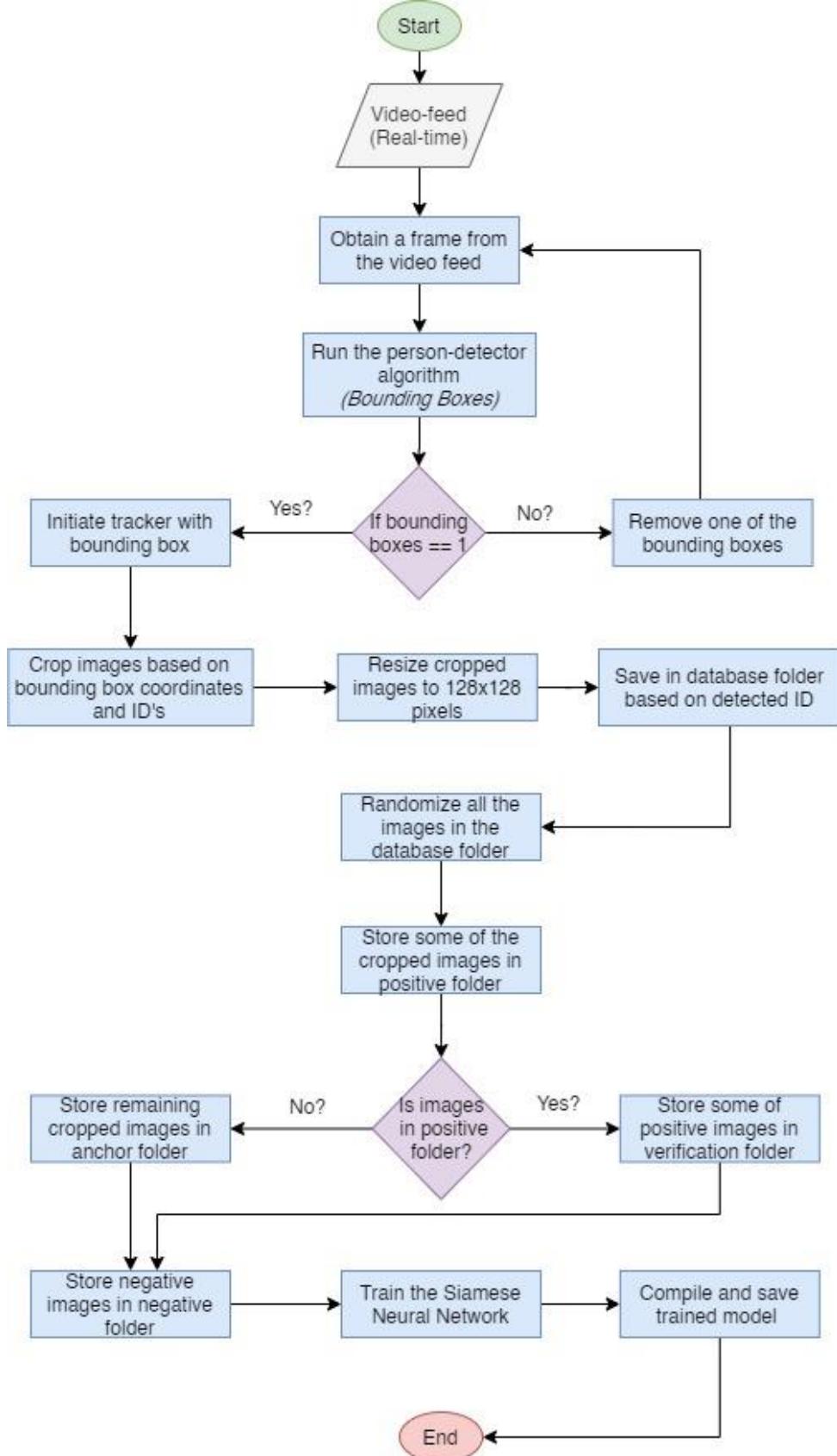
The project, Distributed Video Tracking will take real-time video feed from the internet protocol cameras. It is then divided into two main algorithms; the first algorithm is catered for Pedestrian detection and tracking whilst the second algorithm focuses in reidentifying the pedestrian.



**FIGURE 3.1 – Project Block Diagram**

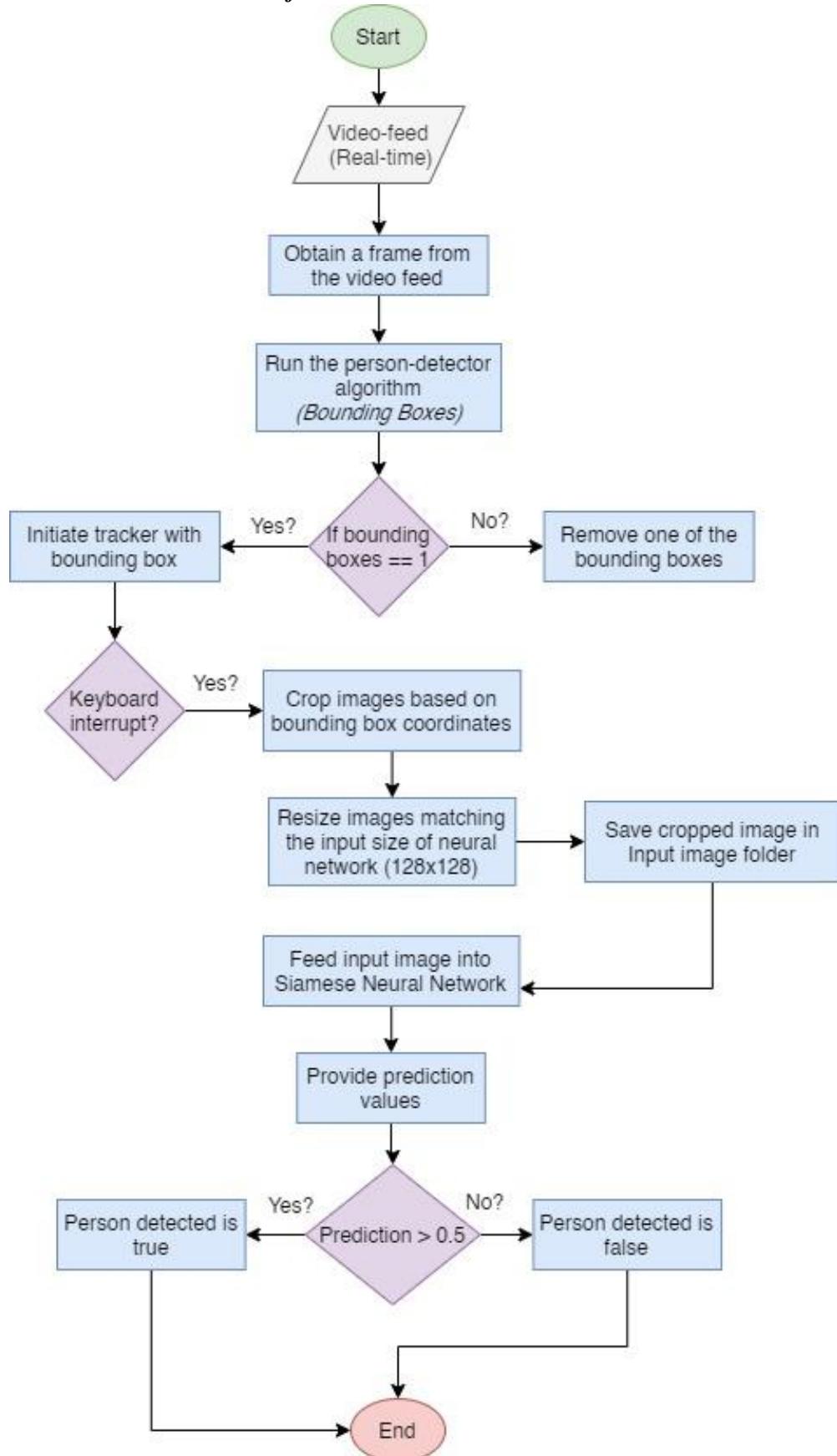
### 3.2 Project Methodology

#### i. Algorithm 1: Pedestrian Detection and Tracking



**FIGURE 3.2 – Pedestrian Detection and Tracking**

*ii. Algorithm 2: Pedestrian Identification*



**FIGURE 3.3 – Pedestrian Identification**

### **3.2.1 Pedestrian Detection and Tracking**

The *Pedestrian Detection and Tracking* algorithm compromises of three main subsections; object-detection, non-maximum suppression, and object tracking. The system would obtain a frame from the video feed (real-time footage from the Internet Protocol cameras) to run the person-detector algorithm. For this, the Faster R-CNN object detection API will be utilized; only selecting the “*person*” class. This would allow the object detector to only identify individual(s) in a continuous frame, ignoring other object classes during detection. Once the person is detected within the video feed, the non-maximum suppression step would be performed. From the obtained bounding boxes, if the number of bounding box is more than one, the remaining or the additional (not required) bounding boxes would be removed, and the step of obtaining a frame from the video feed is repeated. This is a significant step as it removes overlapping bounding boxes or noises when the “*person*” class is detected.

Next, the tracker, Simple Online and Realtime Tracking (SORT) would be initiated and through SORT, the identified bounding box coordinates and their respective IDs would allow for to crop the images of individuals across the frame. It would be resized into 128x128 pixels, and saved in a database folder based on the detected ID. The images would then be randomized and some of the images would be stored into the positive folder. The remaining images (must not be the same as the ones in positive folder), would then be stored into the anchor and verification folder respectively. Then, the negative images, which are images retrieved from a random dataset (CUHK01 or PKUv1a\_128x48) would be stored in the negative folder. Hereafter, the data is prepared, and can be used for training in the Siamese Neural Network. The trained model would then be compiled and saved for future purposes.

In Siamese Neural networks, the input image of a detected person is considered to be the anchor image. Assuming there is the anchor and positive image pair, and the anchor and negative pairs, convolutional neural network would then obtain the similarities between these two images. The positive images, are supposed to be the images of the same person (preferably with a different feature), whilst the negative images are images of a different person. Through this, the Siamese network can easily identify whether or not there are the same person, by encoding the feature representation of the given two images (negative and positive). This is also known as the triple loss function.

### **3.2.2 Pedestrian Identification**

The *Pedestrian identification* algorithm has relatively short procedures, which are more sophisticated to perform. Similarly, the system would obtain a frame from the video feed (real-time footage from the Internet Protocol cameras) to run the person-detector algorithm. Once the person is detected within the video feed, the non-maximum suppression step would be performed. From the obtained bounding boxes, if the number of bounding box is more than one, the remaining or the additional (not required) bounding boxes would be removed, and the step of obtaining a frame from the video feed is repeated.

The tracker, SORT would then be initiated, and the bounding box coordinates will be identified. The system would then continue to display the video input, whilst waiting for a keyboard interrupt. If there was a keyboard interrupt, the system would crop images based on the bounding box coordinates, resize them to the input size of the neural network of 128x128 pixels. The images are then stored into the input image folder, and fed into the Siamese Neural Network. For this project, the implementation was at 128x128 pixels; square dimensions due to the way the model was built upon. As mentioned previously, with an input image of 128x128 pixels, the Siamese model would resize it to 105x105 pixels, before performing the typical steps involved in a neural network; *convolutional, maxpooling, flatten* and *dense*.

The input image would then be compared with the verification images, and hence provide the prediction values. Once the input image instance is fed into the network, the prediction probabilities are made. For probabilities that have any values that are above 0.5, predictions are classified to 1 and 0 for probabilities lesser than 0.5. This ensures that the person of interest detected in a different camera or field of view is the same exact person; returns a Boolean value of **True or False**.

### 3.3 Software, Hardware and Tools

Software, Hardware and Tools	Name	Description
	TP-LINK R1043ND	<ul style="list-style-type: none"> <li>In charge of internet traffic direction</li> <li>Allows the configuration of internet protocol cameras (static IP addresses)</li> </ul>
	D-LINK DCS5030L	<ul style="list-style-type: none"> <li>Obtains the video footages in real-time</li> <li>Supports HTTP stream instead of RTSP</li> <li>Connected to a router</li> <li>Configured to static IP addresses</li> </ul>
	VISUAL STUDIO CODE (PYTHON)	<ul style="list-style-type: none"> <li>An IDE that utilizes Python, which allows to program the whole system.</li> <li>The ability to transition between Python environments with ease</li> </ul>

	OPENCV	<ul style="list-style-type: none"> <li>• Would be used to configure various models from TensorFlow</li> <li>• Take advantage of hardware acceleration and multi-core systems to deploy</li> <li>• Strength are in the deployment side</li> </ul>
	TENSORFLOW	<ul style="list-style-type: none"> <li>• TensorFlow to build, test and train the proposed model</li> <li>• Scalable; allows its users to develop any kind of system using TensorFlow</li> <li>• Strength are in the training side</li> </ul>
	KERAS	<ul style="list-style-type: none"> <li>• Various backend and modularity</li> <li>• Multiple GPU support</li> <li>• Has pretrained models</li> </ul>

**TABLE 1 - Software, Hardware and Tools**

### 3.4 Gantt Chart

FYP I												
Activities/Week	W1	W2	W3	W4	W5	W6	W7	W8	W9	W10	W11	W12
Obtain & understand the topic advised	■											
Perform background study, objectives, and problem statement		■	■	■								
Literature review on the different ML/DNN models					■							
Receive required hardware; IP Cam/Router/Switch					■	■						
Selecting the tools and ML/DNN Models					■	■	■					
Deploy Object Detection API and Tracking.(Algorithm I)							■	■				
Testing on the different object detection models (Faster R-CNN and SSD MobileNet )							■	■				
Testing on the different trackers (Centroid tracker/DCFnet/DeepSORT)									■	■	■	
Developed simple model to crop images from video input based on the coordinates of detected bounding boxes									■			
Result Analysis and further testing									■	■	■	
Documentation									■	■	■	

FYP II												
Activities/Week	W1	W2	W3	W4	W5	W6	W7	W8	W9	W10	W11	W12
Identify approaches to develop re-identification network												
Switching from Centroid tracker to SORT for object tracking												
Testing on transfer-learning approach on Faster R-CNN model for custom object detection												
Develop the Siamese Neural Network (Algorithm II)												
Integrate the two developed models ( Algorithm I and Algorithm II)												
Result Analysis and further testing												
Improvement and optimisation to reduce false detections												
Documentation												

## CHAPTER 4: RESULT & DISCUSSION

### 4.1 Initial Phase

The Faster-RCNN Inception v2 object detection API from TensorFlow was used for in detecting individuals or person in a particular frame. The Inception V2 architecture is used since it is one of the most accurate Convolutional Neural Network (CNN) architectures that exists. Kernels of various sizes are required for successful identification of such a variable-sized feature. That is exactly what Inception accomplishes. It expands rather than just becoming deeper in terms of the number of levels. Within the same layer, many kernels of various sizes are implemented. Despite requiring high graphical power usage, it is able to provide extremely accurate results in identifying and classifying objects detected from a video input.

Comparison					
Network	Year	Salient Feature	top5 accuracy	Parameters	FLOP
AlexNet	2012	Deeper	84.70%	62M	1.5B
VGGNet	2014	Fixed-size kernels	92.30%	138M	19.6B
Inception	2014	Wider - Parallel kernels	93.30%	6.4M	2B
ResNet-152	2015	Shortcut connections	95.51%	60.3M	11B

**FIGURE 4.11 – Comparison between the different Network architectures**

In this project, object detection APIs were imported from TensorFlow for them to be deployed on OpenCV. However, in order to integrate object detection models from TensorFlow, OpenCV requires an additional configuration file. Similarly, as in TensorFlow serialised graph formatted as protocol buffers (protobuf) would be used, OpenCV relies on a text version of it. Besides, Cuda and cuDNN was installed for high-performance GPU acceleration, based on their respective versions, to ensure compatibility with TensorFlow. For Tensorflow 2.8.0, cuDNN 8.1 and CUDA 11.2 were required, while ensuring python 3.x was already installed as base environment. This ensures GPU hardware support during rendering.

Version	Python version	Compiler	Build tools	cuDNN	CUDA
tensorflow-2.7.0	3.7-3.9	GCC 7.3.1	Bazel 3.7.2	8.1	11.2
tensorflow-2.6.0	3.6-3.9	GCC 7.3.1	Bazel 3.7.2	8.1	11.2
tensorflow-2.5.0	3.6-3.9	GCC 7.3.1	Bazel 3.7.2	8.1	11.2
tensorflow-2.4.0	3.6-3.8	GCC 7.3.1	Bazel 3.1.0	8.0	11.0
tensorflow-2.3.0	3.5-3.8	GCC 7.3.1	Bazel 3.1.0	7.6	10.1
tensorflow-2.2.0	3.5-3.8	GCC 7.3.1	Bazel 2.0.0	7.6	10.1

**FIGURE 4.12 – TensorFlow versions and their compatibilities**

The Internet Protocol Cameras (D-LINK DCS5030L) obtains the video footages in real-time. It supports HTTP stream instead of RTSP, and needed to be connected to a router for the static IP addresses configuration. Initially, the camera were connected via Wi-Fi of the same network, streamed with HTTP. This resulted with the issue of reassigning new IP addresses, each time either the cameras or router restarts, as several additional steps were needed to obtain the live stream on the system for further processing. To eliminate this, the cameras were configured to a specific IP address and mapped to their corresponding MAC Addresses in the router configuration settings. For B0-C5-54-2E-2E-E7 it was bind to 192.168.0.30, while B0-C5-54-2D-41-94 to 192.168.0.40. These two cameras were connected to TP-LINK R1043ND router; In charge of internet traffic direction and allows the configuration of internet protocol cameras (static IP addresses).



**FIGURE 4.13 – DLINK-DCS5030L configuration page**

ID	MAC Address	IP Address	Bind	Modify
1	B0-C5-54-2E-2E-E7	192.168.0.30	<input checked="" type="checkbox"/>	<a href="#">Modify</a> <a href="#">Delete</a>
2	B0-C5-54-2D-41-94	192.168.0.40	<input checked="" type="checkbox"/>	<a href="#">Modify</a> <a href="#">Delete</a>

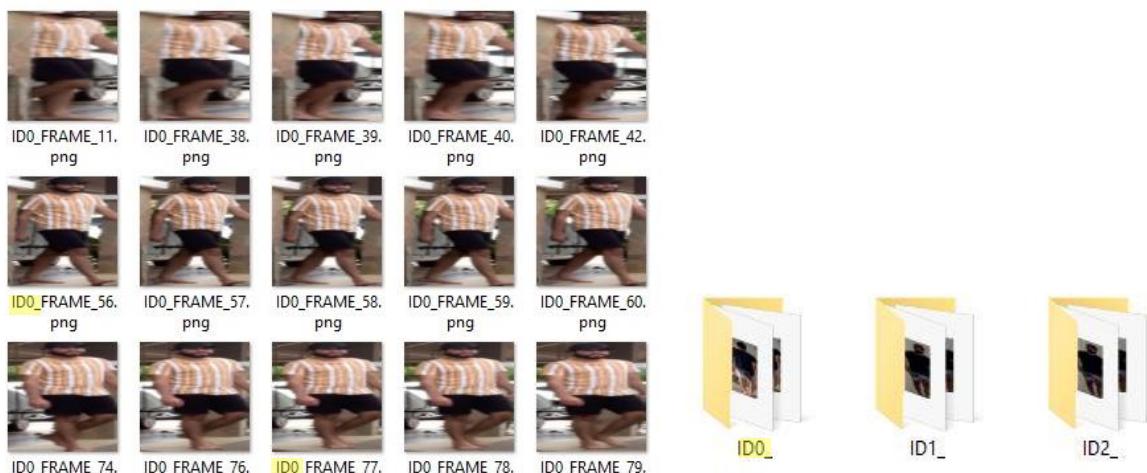
**FIGURE 4.14 – TP-LINK R1043ND router configuration page**

## 4.2 Person Detection and Tracking

The hardware specifications of the computer used throughout the testing were; NVIDIA GTX 1060 with 6GB (gigabytes) of VRAM. This was made possible due to CUDA, which enables applications to use specific types of graphics processing units. Three inputs were used to test the overall performance of the Faster-RCNN Inception v2 in detecting “*person*” class; two pre-recorded videos with distinct recording time at 60fps, and a live stream through HTTP at 30fps with the D-LINK DCS5030L IP cameras. Two of these video inputs all have an individual appearing in and out of the frames, albeit from different angles, light setting, and locations, which challenges the effectiveness of the object detection model.

The Faster-RCNN Inception V2, accompanied by the Simple Online and Realtime Tracking (SORT) would work simultaneously to detect a person moving across the frame. This would work regardless on the type of input, whether it is from a pre-recorded video, or a HTTP stream. As mentioned, the identified bounding box coordinates and their respective IDs would allow to crop the images of individuals as he moves across the frames in the video. It would then be resized into 128x128 pixels, and saved in a folder based on the detected ID. It is also to be reminded that the names of the cropped images would be titled as follows; “*path + str(i) + '\_FRAME\_' + cur\_strg +'.png'*”.

The number of frames and detected ID play a significant role in segregating the cropped images based on the ID identified. If there are more than a person detected in the video input, then the ID would be classified based on the detected ID and placed into their corresponding folders. The figure 4.15 below, demonstrates the working principle of the algorithm described, highlighted in yellow.



**FIGURE 4.15** - Demonstrates the working principle of the algorithm

Figures 4.16 and 4.17, SORT detects a person as ID\_0, which would result the images that are cropped to be placed in a folder respective to the detected ID. So, in this particular case, the cropped images are all stored in a folder named “ID\_0”. The cropped images are then resized to 128x128 pixels to ensure the standardization of all images (positive, anchor and negative), before being fed to the Siamese model. The neural network was built to only receive square images of 105x105 pixels, which quickly translates to being able to identify images that have sides that are parallel and equal in dimensions.

```
f.default_graph.close()

__main__":
path = 'D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/FYP2/AIComputerVision-master/SORT/faster_rcnn_inception_v2/frozen_inference_graph.pb'
DetectorAPI(path_to_ckpt=model_path)
ld = 0.5
cv2.VideoCapture('http://admin:0BtHmuSP@192.168.0.30/VIDEO.CGI')
cv2.VideoCapture('test_video.mp4')
cv2.VideoCapture('IMG_0286.MOV')
int(cap.get(cv2.CAP_PROP_FRAME_HEIGHT))
int(cap.get(cv2.CAP_PROP_FRAME_WIDTH))
Height:, height)
Width:, width)

= cap.read()
imutils.resize(img, width=600)
cv2.resize(img, (1280, 720))

scores, classes, num = odapi.preprocessing.visualization_of_the_results_of_a_detections()
for i in range(len(boxes)):
    # Class 1 represents human
    if classes[i] == 1 and scores[i] > ld:
        box = boxes[i]
        objectID = "ID_{}".format(i)
        cv2.rectangle(img, (box[1], box[0]-5), cv2.FONT_HERSHEY_COMPLEX_SMALL, 0.8, (208,224,64), 1)
        cv2.putText(img,objectID,(box[1],box[0]-5),cv2.FONT_HERSHEY_COMPLEX_SMALL, 0.8, (208,224,64),1)
        cropped_image = img[box[0]+2:box[2]-2,box[1]+2:box[3]-2]
        path = 'D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/frames/ID'+str(i)+'_FRAME_'+ cur_strg+'.png'
        #cv2.imwrite(path,cropped_image)
```

**FIGURE 4.16 – Video input from a pre-recorded video(1)**

```
f.default_graph.close()

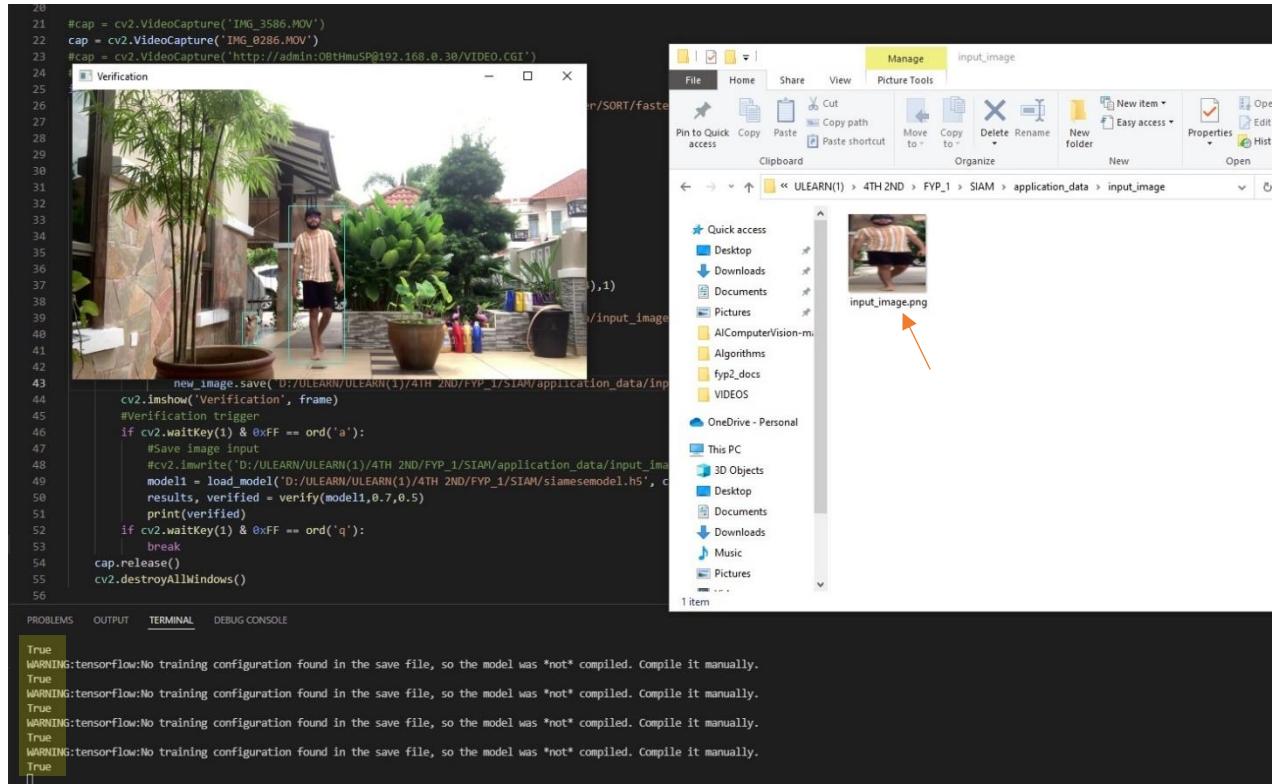
__main__":
path = 'D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/FYP2/AIComputerVision-master/SORT/faster_rcnn_inception_v2/frozen_inference_graph.pb'
DetectorAPI(path_to_ckpt=model_path)
ld = 0.5
cv2.VideoCapture('http://admin:0BtHmuSP@192.168.0.30/VIDEO.CGI')
cv2.VideoCapture('test_video.mp4')
cv2.VideoCapture('IMG_0286.MOV')
int(cap.get(cv2.CAP_PROP_FRAME_HEIGHT))
int(cap.get(cv2.CAP_PROP_FRAME_WIDTH))
Height:, height)
Width:, width)

true:
img = cap.read()
imutils.resize(img, width=600)
g = cv2.resize(img, (1280, 720))

scores, classes, num = odapi.preprocessing.visualization_of_the_results_of_a_detections()
for i in range(len(boxes)):
    # Class 1 represents human
    if classes[i] == 1 and scores[i] > ld:
        box = boxes[i]
        objectID = "ID_{}".format(i)
        cv2.rectangle(img, (box[1], box[0]-5), cv2.FONT_HERSHEY_COMPLEX_SMALL, 0.8, (208,224,64), 1)
        cv2.putText(img,objectID,(box[1],box[0]-5),cv2.FONT_HERSHEY_COMPLEX_SMALL, 0.8, (208,224,64),1)
        cropped_image = img[box[0]+2:box[2]-2,box[1]+2:box[3]-2]
        path = 'D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/frames/ID'+str(i)+'_FRAME_'+ cur_strg+'.png'
        #cv2.imwrite(path,cropped_image)
        #image = Image.open(path)
        #new_image = image.resize((128, 128))
```

**FIGURE 4.17 – Video input from a pre-recorded video(2)**

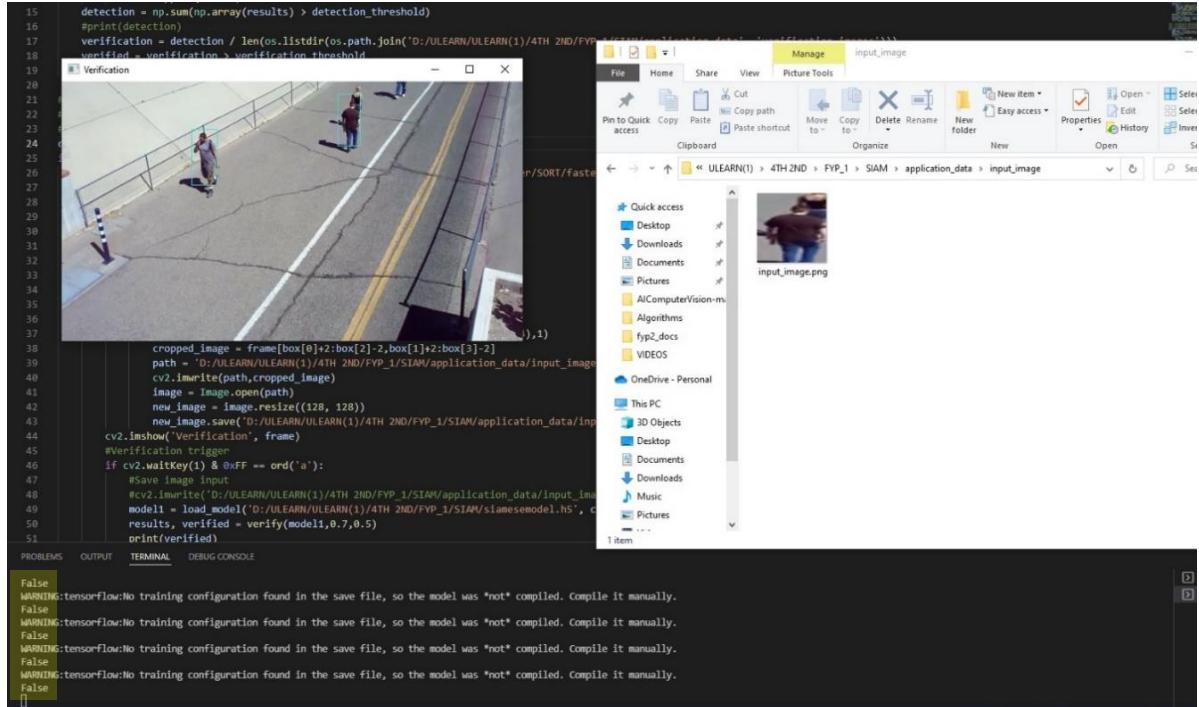
Once the training is complete, another recorded video, preferably from a different angle was used to test the legitimacy of the training results. But first, it would be ideal to perform a simple experiment to test the Siamese Model on the video used for the training itself. Figure 4.18 shows the video that was used for training in performing the verification, and that the results detected are mostly positive (seen in the bottom left, highlighted in yellow). This means that the Siamese Model was able to identify the person as the same person, that was used as the anchor and positive images.



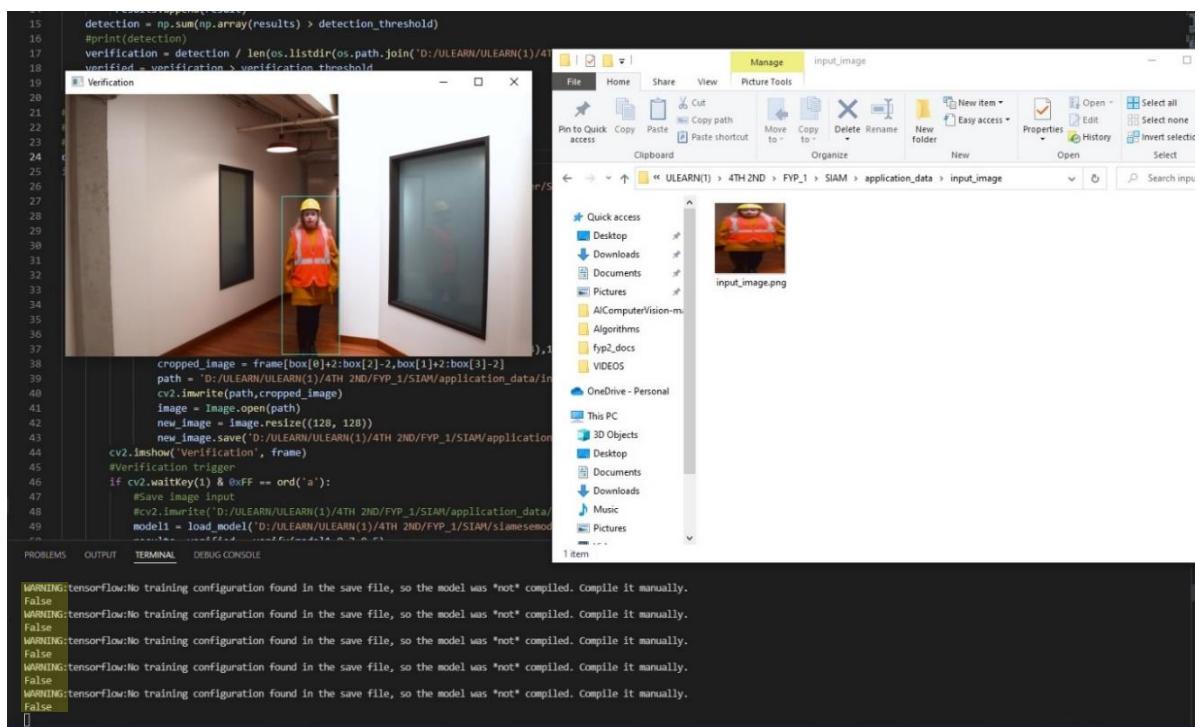
**FIGURE 4.18 –To test the Siamese Model on the video used for the training itself**

On the right of Figure 4.18, shows the input image that would be used for verification. This is also the cropped images that are based off the bounding box coordinates. The input image would then be saved as ***input\_image.png***, which can be triggered by a keyboard interrupt. This is simply an instance of the image, meaning, that each time there is a keyboard interrupt, the next detected bounding box would replace the current cropped image. The ***input\_image.png*** is basically the image that would be fed into the neural network where the probabilities would be predicted. The test were done with a verification threshold of 0.7; any predicted probabilities of over 70% would return a Boolean value of **True**, and anything below as **False**. As can be seen, the model does a pretty decent job in identifying the same person.

The test was then mimicked on two random videos, simply to confirm that the Siamese Neural Network developed did not fake the results obtained. The figures 4.19 and 4.20 below shows that the person detected is not the same person that was used in training, or in other words the anchor and positive images. Consecutively, the result would be **False**, classifying that this particular person is not recognised.



**FIGURE 4.19 – Random video input (1)**



**FIGURE 4.20 – Random video input (2)**

This time, however, the second video input simulates a different camera angle, of which the same person walks across the frame. Take a look in the figures 4.21 and 4.22 below. Unsurprisingly, the model was able to detect the individual as the same person detected from the first video input, and hence returns a Boolean value of **True**. The Siamese Model satisfactorily produces results that was able to classify that this particular person, was indeed the same exact person from the first video input, albeit from a different camera perspective.



```
18 verified = verification > verification_threshold
19 return results, verified
20
21 cap = cv2.VideoCapture('IMG_0285.MP4') # Verification
22 #cap = cv2.VideoCapture('IMG_0286.MP4')
23 #cap = cv2.VideoCapture('testvideo.mp4')
24 if __name__ == '__main__':
25     model_path = 'D:/ULEARN/ULEARN'
26     odapi = DetectorAPI(path_to_ckpt= model_path)
27     threshold = 0.5
28     while True:
29         ret, frame = cap.read()
30         frame = imutils.resize(frame, width=600)
31         boxes, scores, classes, num = odapi.detect(frame)
32         for i in range(len(boxes)):
33             if classes[i] == 1 and scores[i] > threshold:
34                 box = boxes[i]
35                 objectID = "ID: {}"
36                 cv2.rectangle(frame, box, (0, 255, 0))
37                 cropped_image = frame[box[1]:box[1]+box[3], box[0]:box[0]+box[2]]
38                 path = 'D:/ULEARN/ULEARN/1/4TH 2ND/FYP_1/SIAM/application_data/input_image/input_image.png'
39                 cv2.imwrite(path, cropped_image)
40                 image = Image.open(path)
41                 new_image = image.resize((128, 128))
42                 new_image.save('D:/ULEARN/ULEARN/1/4TH 2ND/FYP_1/SIAM/application_data/input_image/input_image.png')
43             cv2.imshow("Verification", frame)
44     #Verification trigger
45     if cv2.waitKey(1) & 0xFF == ord('a'):
46         #Save image input
47         #cv2.imwrite('D:/ULEARN/ULEARN/1/4TH 2ND/FYP_1/SIAM/application_data/input_image/input_image.png',frame)
48         model = load_model('D:/ULEARN/ULEARN/1/4TH 2ND/FYP_1/SIAM/siamsemamodel.h5', custom_objects={'L1Dist':L1dist, 'BinaryCrossentropy':tf.losses.BinaryCrossentropy})
49         results, verified = verify(model,0.7,0.5)
50
51 o.a
```

PROBLEMS OUTPUT TERMINAL DEBUG CONSOLE

WARNING:tensorflow:No training configuration found in the save file, so the model was \*not\* compiled. Compile it manually.  
True  
WARNING:tensorflow:No training configuration found in the save file, so the model was \*not\* compiled. Compile it manually.  
False  
WARNING:tensorflow:No training configuration found in the save file, so the model was \*not\* compiled. Compile it manually.  
True  
WARNING:tensorflow:No training configuration found in the save file, so the model was \*not\* compiled. Compile it manually.  
True  
WARNING:tensorflow:No training configuration found in the save file, so the model was \*not\* compiled. Compile it manually.  
True  
WARNING:tensorflow:No training configuration found in the save file, so the model was \*not\* compiled. Compile it manually.  
True  
WARNING:tensorflow:No training configuration found in the save file, so the model was \*not\* compiled. Compile it manually.  
False

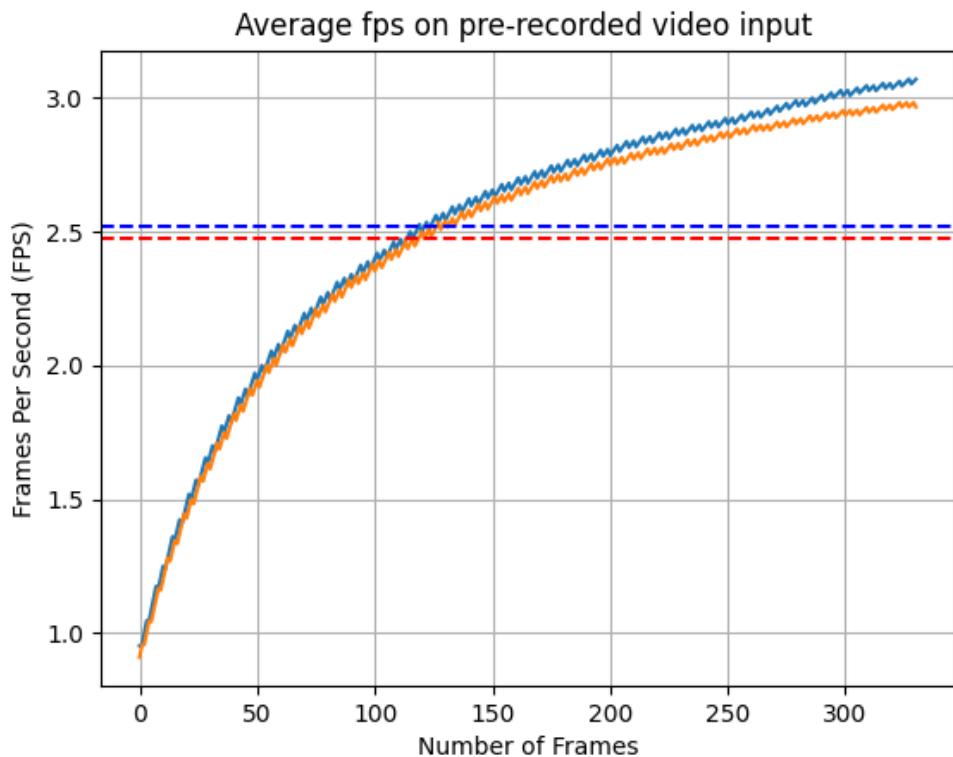
**FIGURE 4.21** – A pre-recorded video, from a different angle (1)

```
18 verified -> verification > verification_threshold
19
20 return results, verified
21
22 cap = cv2.VideoCapture('IMG_0285.MP4') Verification
23 #cap = cv2.VideoCapture('IMG_0286.MP4')
24 #cap = cv2.VideoCapture('testvideo.mp4')
25 if __name__ == "__main__":
26     model_path = "D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/STIAW/application_data/input_image/input_image.pb"
27     odapi = DetectorAPI(path_to_ckpt=model_path)
28     threshold = 0.5
29
30     while True:
31         ret, frame = cap.read()
32         frame = imutils.resize(frame, width=600)
33         boxes, scores, classes, num = odapi.detect(frame)
34         for i in range(len(boxes)):
35             if classes[i] == 1 and scores[i] > threshold:
36                 box = boxes[i]
37                 objectID = "ID {}".format(i)
38                 cv2.rectangle(frame, (box[0], box[1]), (box[2], box[3]), (0, 255, 0), 2)
39                 cropped_image = frame[box[1]:box[3], box[0]:box[2]]
40                 path = "D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/STIAW/application_data/input_image/input_image.png"
41                 cv2.imwrite(path, cropped_image)
42                 image = Image.open(path)
43                 new_image = image.resize((128, 128))
44                 new_image.save("D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/STIAW/application_data/input_image/input_image.png")
45         cv2.imshow('Verification', frame)
46         #Verification trigger
47         if cv2.waitKey(1) & 0xFF == ord('a'):
48             #Save image input
49             cv2.imwrite("D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/STIAW/application_data/input_image/input_image.png",frame)
50             model1 = load_model("D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/STIAW/siammodel.h5", custom_objects={'L1Dist':L1Dist, 'BinaryCrossentropy':tf.losses.BinaryCrossentropy})
51             results, verified = verify(model1,0.7,0.5)
52             print(results,verified)
53
54
55
56
57
58
59
59
```

**FIGURE 4.22** –A pre-recorded video, from a different angle (2)

The frame rate is the number of pictures collected or shown in a row. The average frames per second recorded from these two-video inputs were approximately, 2.521 and 2.477 respectively. Undeniably, these results ridiculously slow, despite the fact that there were from inputs of already recorded video. The primary reason to this is due to the fact that both the object detection model, Faster-RCNN Inception V2 and the tracker, SORT are extremely GPU intensive. The system was run on a satisfactory and relatively mid-range NVIDIA GPU, causing the overall running speed to be time-consuming.

The figure 4.23 below represents the average fps on both the pre-recorded video input. The frames spikes up from the very beginning and level out evenly around the average of 2.5fps. It is also to be clarified that the data were collected from a mere 330 total number of frames, due to the inconsistent video duration from the inputs.



**FIGURE 4.23 – Average fps on pre-recorded video input from two distinct angles**

The receiver operating characteristic curve (ROC curve) is a graph that shows how well a classification model performs across all categorization levels. There are only two parameters are shown on this curve, true positive rate, and false positive rate. Hence, the ROC curve is simply the plot of the true positive rate (TPR) against the false positive rate (FPR). As the classification threshold is lowered, more items would be classified as positive resulting in an increase in both the false positives and true positives.

*The True Positive Rate can be defined by the following formula;*

$$TPR = \frac{TP}{TP + FN}$$

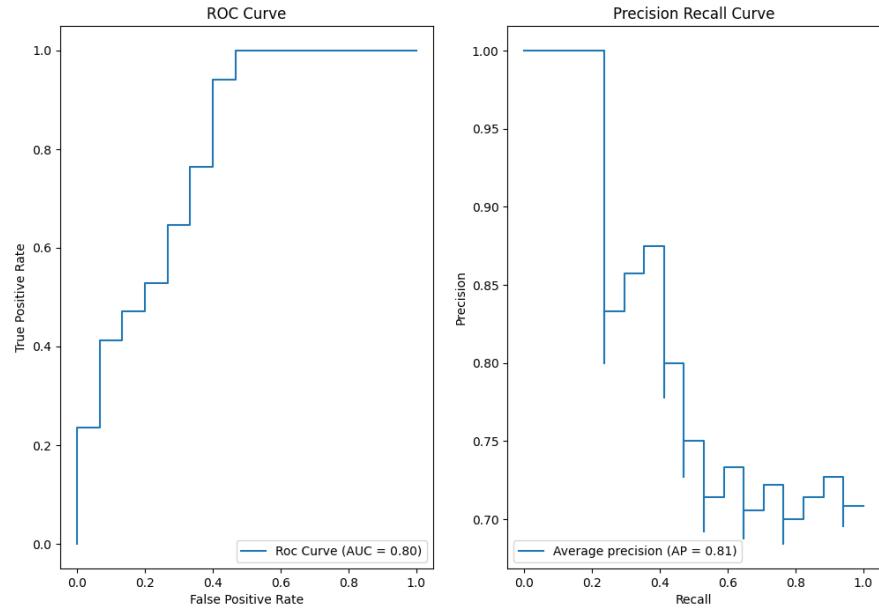
*Whereas the False Negative Rate is defined as follows;*

$$FPR = \frac{FP}{FP + TN}$$

The area under the curve (AUC) is a summary of the ROC curve that measures a classifier's ability to differentiate between classes. The AUC indicates how well the model differentiates between positive and negative classes. The greater the AUC, the better the model has done in predicting. For instance, for the case when AUC is 1, the classifier is capable of successfully distinguishing between all Positive and Negative class points. Contrarily, if the AUC is closer to 0, the classifier would classify all Negatives as Positives and all Positives as Negatives.

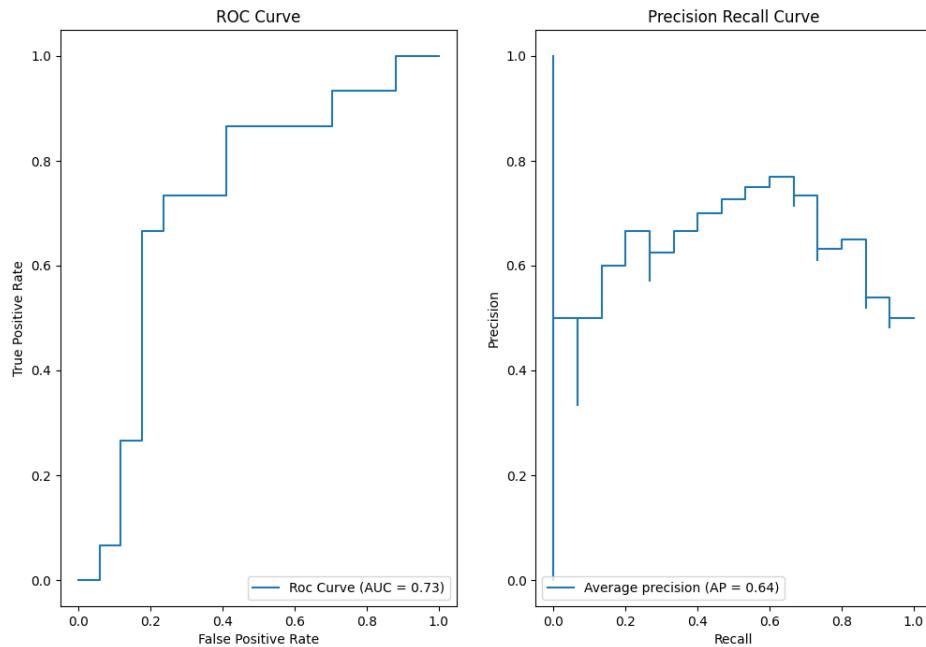
The precision recall curve (PRC), is used to monitor the efficiency and performance of binary classification systems. It is widely used in circumstances where the classes are severely disproportionate. The PRC plot provides a graphical depiction of the classifier's performance over a range of thresholds, where the values of precision are plotted against recall. Ideally, the plot should have a decreasing value and if the average precision is 1, means that the system can correctly classify the true positive and negative cases. Average precision closer to zero indicates that a particular developed system does not perform well overall. Both the ROC and PR curves for the trained model in the first test, can be represented in the following figures; Figure 4.24 and 4.25.

Each batch has a set of 32 images which would be passed for training, which were then consecutively trained for 10 epochs. To be fair, these results were obtained on two random batches, rather than the whole test sample. On the first batch, the area under the curve, AUC was 0.81, whilst the average precision in the PR curve was at 0.81.



**FIGURE 4.24 – ROC Curve and Precision Recall curves from a pre-recorded video (1)**

On a second random batch, the area under the curve, AUC was 0.73, with an average precision in the PR curve of 0.64. This comes to show that this particular model, on average, does a decent job in identifying the person of interest against other people.



**FIGURE 4.25 – ROC Curve and Precision Recall curves from a pre-recorded video (2)**

The second test was performed by receiving a live stream input from the DCS5030L Internet Protocol camera. The figures 4.26, 4.27, 4.28 and 4.29 shows the test held in a bright room with minor occlusion. This time, the video input oversees an individual lingering across the frame, long enough to obtain enough images for training. Similarly, to collect images for training, the coordinates of the bounding boxes were used.



```

83     for i in range(len(boxes)):
84         # Class 1 represents human
85         if classes[i] == 1 and scores[i] > threshold:
86             box = boxes[i]
87             objectID = "ID {}".format(i)
88             cv2.rectangle(img,(box[1],box[0]),(box[3],box[2]),(208,224,64),1)
89             cv2.putText(img,objectID,(box[1],box[0]-5),cv2.FONT_HERSHEY_COMPLEX_SMALL, 0.8, (208,224,64),1)
90             cropped_image = img[box[0]+2:box[2]-2,box[1]+2:box[3]-2]
91             path = 'D:/ULEARN/ULEARN1/4TH 2ND/FYP_1/frames/ID'+str(i)+'/FRAME '+ cur_strg+'.png'
92             cv2.imwrite(path,cr)
93             image = Image.open(path)
94             new_image = image.resize((128,128))
95             new_image.save('D:/')
96             currentFrame = cap.read()
97             cur_strg = str(currentFrame)
98             cur_strg = cur_strg
99
100            cv2.imshow("IMG", img)
101            key = cv2.waitKey(1)
102            if key & 0xFF == ord('q'):
103                dir_list = os.listdir(path)
104                for f in dir_list:
105                    dir_name = f[:4]
106                    print(dir_name)
107                    dir_path = path + '/' + dir_name
108                    if not os.path.exists(dir_path):
109                        os.makedirs(dir_path)
110                    if os.path.exists(dir_path):
111                        file_path = path + '/' + f
112                        shutil.move(file_path, dir_path)
113                break
114
115
116
117
PROBLEMS OUTPUT TERMINAL DEBUG CONSOLE
Height: 720
Width: 1280
2022-03-17 15:31:15.303007: W tensorflow/core/grappler/costs/op_level_cost_estimator.cc:690] Error in PredictCost() for the op: op: "CropAndResize" attr { key: "T" value { type: DT_FLOAT } } attr { key: "extrapolation_value" value { f: 0 } } attr { key: "method" value { s: "bilinear" } } inputs { dtype: DT_FLOAT shape { dim { size: -364 } dim { size: -366 } dim { size: -367 } dim { size: -365 } } } inputs { dtype: DT_FLOAT shape { dim { size: -24 } dim { size: 4 } } } inputs { dtype: DT_INT32 shape { dim { size: -24 } } } inputs { dtype: DT_INT32 shape { dim { size: 2 } } value { dtype: DT_INT32 tensor_shape { dim { size: 2 } } } }
int_val: 14 } } device { type: "GPU" vendor: "NVIDIA" model: "NVIDIA GeForce GTX 1060 6GB" frequency: 1759 num_cores: 10 environment { key: "architecture" value: "6.1" } environment { key: "cuda" value: "11020" } environment { key: "cudnn" value: "8100" } num_registers: 65536 l1_cache_size: 24576 l2_cache_size: 1572864 shared_memory_size_per_multiprocessor: 98304 memory_size: 4857331712 bandwidth: 192192000 } output { s { dtype: DT_FLOAT shape { dim { size: -24 } dim { size: 14 } dim { size: 14 } dim { size: -365 } } }
2022-03-17 15:31:19.593233: I tensorflow/stream_executor/cuda/cuda_dnn.cc:368] Loaded cuDNN version 8100

```

**FIGURE 4.26 – Video input from live stream HTTP (1)**



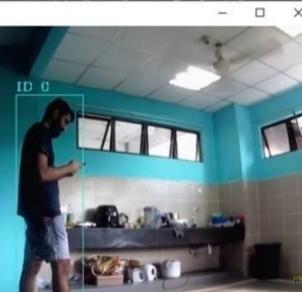
```

83     for i in range(len(boxes)):
84         # Class 1 represents human
85         if classes[i] == 1 and scores[i] > threshold:
86             box = boxes[i]
87             objectID = "ID {}".format(i)
88             cv2.rectangle(img,(box[1],box[0]),(box[3],box[2]),(208,224,64),1)
89             cv2.putText(img,objectID,(box[1],box[0]-5),cv2.FONT_HERSHEY_COMPLEX_SMALL, 0.8, (208,224,64),1)
90             cropped_image = img[box[0]+2:box[2]-2,box[1]+2:box[3]-2]
91             path = 'D:/ULEARN/ULEARN1/4TH 2ND/FYP_1/frames/ID'+str(i)+'/FRAME '+ cur_strg+'.png'
92             cv2.imwrite(path,cr)
93             image = Image.open(path)
94             new_image = image.resize((128,128))
95             new_image.save('D:/')
96             currentFrame = cap.read()
97             cur_strg = str(currentFrame)
98             cur_strg = cur_strg
99
100            cv2.imshow("IMG", img)
101            key = cv2.waitKey(1)
102            if key & 0xFF == ord('q'):
103                dir_list = os.listdir(path)
104                for f in dir_list:
105                    dir_name = f[:4]
106                    print(dir_name)
107                    dir_path = path + '/' + dir_name
108                    if not os.path.exists(dir_path):
109                        os.makedirs(dir_path)
110                    if os.path.exists(dir_path):
111                        file_path = path + '/' + f
112                        shutil.move(file_path, dir_path)
113                break
114
115
116
117
PROBLEMS OUTPUT TERMINAL DEBUG CONSOLE
Height: 720
Width: 1280
2022-03-17 15:31:15.303007: W tensorflow/core/grappler/costs/op_level_cost_estimator.cc:690] Error in PredictCost() for the op: op: "CropAndResize" attr { key: "T" value { type: DT_FLOAT } } attr { key: "extrapolation_value" value { f: 0 } } attr { key: "method" value { s: "bilinear" } } inputs { dtype: DT_FLOAT shape { dim { size: -364 } dim { size: -366 } dim { size: -367 } dim { size: -365 } } } inputs { dtype: DT_FLOAT shape { dim { size: -24 } dim { size: 4 } } } inputs { dtype: DT_INT32 shape { dim { size: -24 } } } inputs { dtype: DT_INT32 shape { dim { size: 2 } } value { dtype: DT_INT32 tensor_shape { dim { size: 2 } } } }
int_val: 14 } } device { type: "GPU" vendor: "NVIDIA" model: "NVIDIA GeForce GTX 1060 6GB" frequency: 1759 num_cores: 10 environment { key: "architecture" value: "6.1" } environment { key: "cuda" value: "11020" } environment { key: "cudnn" value: "8100" } num_registers: 65536 l1_cache_size: 24576 l2_cache_size: 1572864 shared_memory_size_per_multiprocessor: 98304 memory_size: 4857331712 bandwidth: 192192000 } output { s { dtype: DT_FLOAT shape { dim { size: -24 } dim { size: 14 } dim { size: 14 } dim { size: -365 } } }
2022-03-17 15:31:19.593233: I tensorflow/stream_executor/cuda/cuda_dnn.cc:368] Loaded cuDNN version 8100

```

**FIGURE 4.27 – Video input from live stream HTTP (2)**

Since this was a test to challenge the accuracy of the predicted model, the individual in the video input purposely lingered longer than usual, at least if compared to a typical input from a live stream. This was to ensure that enough images were collected and resized before being trained into the network.



```

83     for i in range(len(boxes)):
84         # Class 1 represents human
85         if classes[i] == 1 and scores[i] > threshold:
86             box = boxes[i]
87             objectID = "ID {}".format(i)
88             cv2.rectangle(img,(box[1],box[0]),(box[3],box[2]),(208,224,64),1)
89             cv2.putText(img,objectID,(box[1],box[0]-5),cv2.FONT_HERSHEY_COMPLEX_SMALL, 0.8, (208,224,64),1)
90             cropped_image = img[box[0]:box[2]:box[1]-2,box[1]:box[1]+2:box[3]-2]
91             path = 'D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/images/ID'+str(i)+'.png'
92             cv2.imwrite(path, cropped_image)
93             image = Image.open(path)
94             new_image = image.resize((250, 250))
95             new_image.save('D:/I')
96             currentFrame = cap.read()
97             cur_strg = str(cur)
98             cur_strg = cur_strg + cur_strg
99
100            cv2.imshow("IMG", img)
101            key = cv2.waitKey(1)
102            if key & 0xFF == ord('q'):
103                dir_list = os.listdir(path)
104                for f in dir_list:
105                    dir_name = f[:4]
106                    print(dir_name)
107                    dir_path = path1 + dir_name
108                    if not os.path.exists(dir_path):
109                        os.makedirs(dir_path)
110                    if os.path.exists(dir_path):
111                        file_path = path1 + f
112                        shutil.move(file_path, dir_path)
113                break
114
115
116
117

```

PROBLEMS OUTPUT TERMINAL DEBUG CONSOLE

Height: 720  
Width: 1280  
2022-03-17 15:31:15.303007: W tensorflow/core/grappler/costs/op\_level\_cost\_estimator.cc:690] Error in PredictCost() for the op: op: "CropAndResize" attr { key: "T" value { type: DT\_FLOAT } } attr { key: "extrapolation\_value" value { f: 0 } } attr { key: "method" value { s: "bilinear" } } inputs { dtype: DT\_FLOAT shape { dim { size: -364 } dim { size: -366 } dim { size: -367 } dim { size: -365 } } } inputs { dtype: DT\_FLOAT shape { dim { size: -24 } dim { size: 4 } } } inputs { dtype: DT\_INT32 shape { dim { size: -24 } } } inputs { dtype: DT\_INT32 shape { dim { size: 2 } } value { dtype: DT\_INT32 tensor\_shape { dim { size: 2 } } int\_val: 14 } } device { type: "GPU" vendor: "NVIDIA" model: "NVIDIA GeForce GTX 1060 6GB" frequency: 1759 num\_cores: 10 environment { key: "architecture" value: "6.1" } environment { key: "cuda" value: "11020" } environment { key: "cudnn" value: "8100" } num\_registers: 65536 l1\_cache\_size: 24576 l2\_cache\_size: 1572864 shared\_memory\_size\_per\_multiprocessor: 98304 memory\_size: 4857331712 bandwidth: 192192000 } output s { dtype: DT\_FLOAT shape { dim { size: -24 } dim { size: 14 } dim { size: 14 } dim { size: -365 } } }  
2022-03-17 15:31:19.593233: I tensorflow/stream\_executor/cuda/cuda\_dnn.cc:368] Loaded cuDNN version 8100

**FIGURE 4.28 – Video input from live stream HTTP (3)**



```

82     # visualization of img / results of a detection.
83     for i in range(len(boxes)):
84         # Class 1 represents human
85         if classes[i] == 1 and scores[i] > threshold:
86             box = boxes[i]
87             objectID = "ID {}".format(i)
88             cv2.rectangle(img,(box[1],box[0]),(box[3],box[2]),(208,224,64),1)
89             cv2.putText(img,objectID,(box[1],box[0]-5),cv2.FONT_HERSHEY_COMPLEX_SMALL, 0.8, (208,224,64),1)
90             cropped_image = img[box[0]:box[2]:box[1]-2,box[1]:box[1]+2:box[3]-2]
91             path = 'D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/images/ID'+str(i)+'.png'
92             cv2.imwrite(path, cropped_image)
93             image = Image.open(path)
94             new_image = image.resize((250, 250))
95             new_image.save('D:/I')
96             currentFrame = cap.read()
97             cur_strg = str(cur)
98             cur_strg = cur_strg + cur_strg
99
100            cv2.imshow("IMG", img)
101            key = cv2.waitKey(1)
102            if key & 0xFF == ord('q'):
103                dir_list = os.listdir(path)
104                for f in dir_list:
105                    dir_name = f[:4]
106                    print(dir_name)
107                    dir_path = path1 + dir_name
108                    if not os.path.exists(dir_path):
109                        os.makedirs(dir_path)
110                    if os.path.exists(dir_path):
111                        file_path = path1 + f
112                        shutil.move(file_path, dir_path)
113                break
114
115
116
117

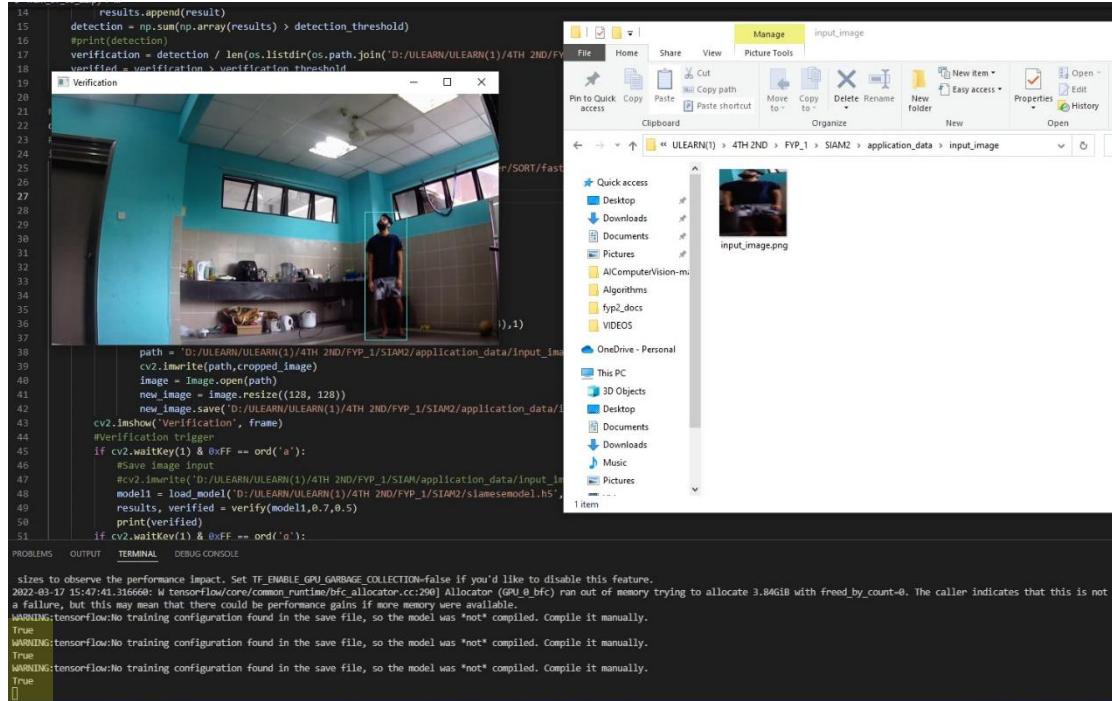
```

PROBLEMS OUTPUT TERMINAL DEBUG CONSOLE

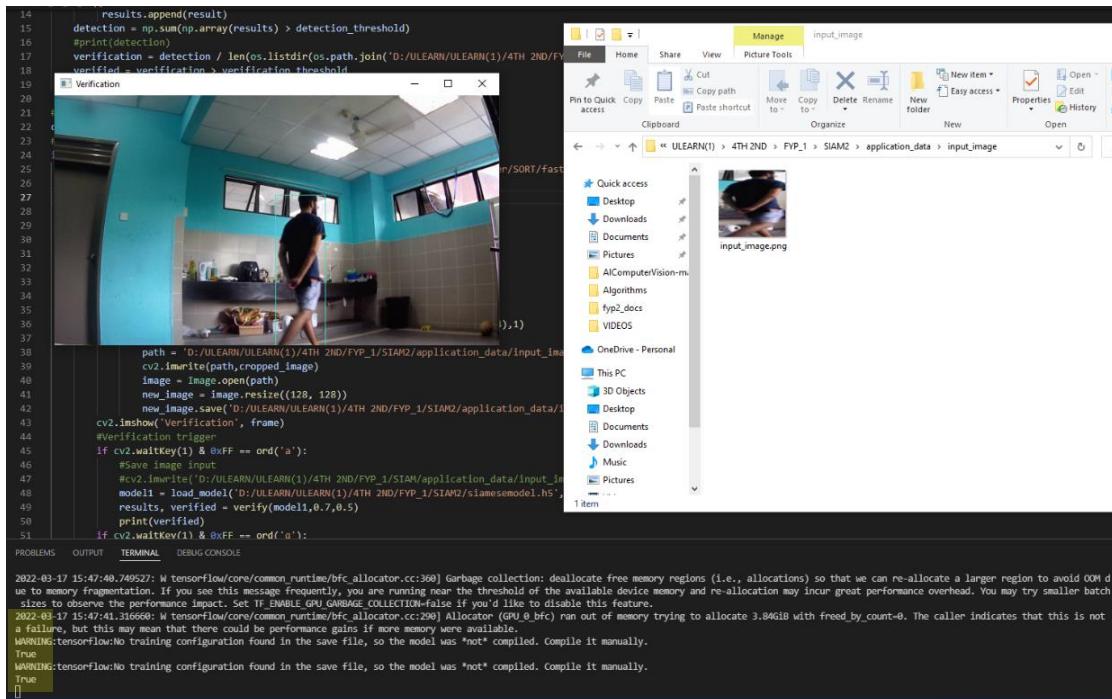
Height: 720  
Width: 1280  
2022-03-17 15:31:15.303007: W tensorflow/core/grappler/costs/op\_level\_cost\_estimator.cc:690] Error in PredictCost() for the op: op: "CropAndResize" attr { key: "T" value { type: DT\_FLOAT } } attr { key: "extrapolation\_value" value { f: 0 } } attr { key: "method" value { s: "bilinear" } } inputs { dtype: DT\_FLOAT shape { dim { size: -364 } dim { size: -366 } dim { size: -367 } dim { size: -365 } } } inputs { dtype: DT\_FLOAT shape { dim { size: -24 } dim { size: 4 } } } inputs { dtype: DT\_INT32 shape { dim { size: -24 } } } inputs { dtype: DT\_INT32 shape { dim { size: 2 } } value { dtype: DT\_INT32 tensor\_shape { dim { size: 2 } } int\_val: 14 } } device { type: "GPU" vendor: "NVIDIA" model: "NVIDIA GeForce GTX 1060 6GB" frequency: 1759 num\_cores: 10 environment { key: "architecture" value: "6.1" } environment { key: "cuda" value: "11020" } environment { key: "cudnn" value: "8100" } num\_registers: 65536 l1\_cache\_size: 24576 l2\_cache\_size: 1572864 shared\_memory\_size\_per\_multiprocessor: 98304 memory\_size: 4857331712 bandwidth: 192192000 } output s { dtype: DT\_FLOAT shape { dim { size: -24 } dim { size: 14 } dim { size: 14 } dim { size: -365 } } }  
2022-03-17 15:31:19.593233: I tensorflow/stream\_executor/cuda/cuda\_dnn.cc:368] Loaded cuDNN version 8100

**FIGURE 4.29 – Video input from live stream HTTP (4)**

As mentioned previously, by having enough images the Siamese Model would be able to make predictions much more error-free and perfect by classifying that this person of interest is indeed the same person detected in a different camera, most likely from a different view. But first, the test was run on the camera from the same angle used to collect the train images. As can be seen in figures 4.30 and 4.31, the model detects the same person as **True**.

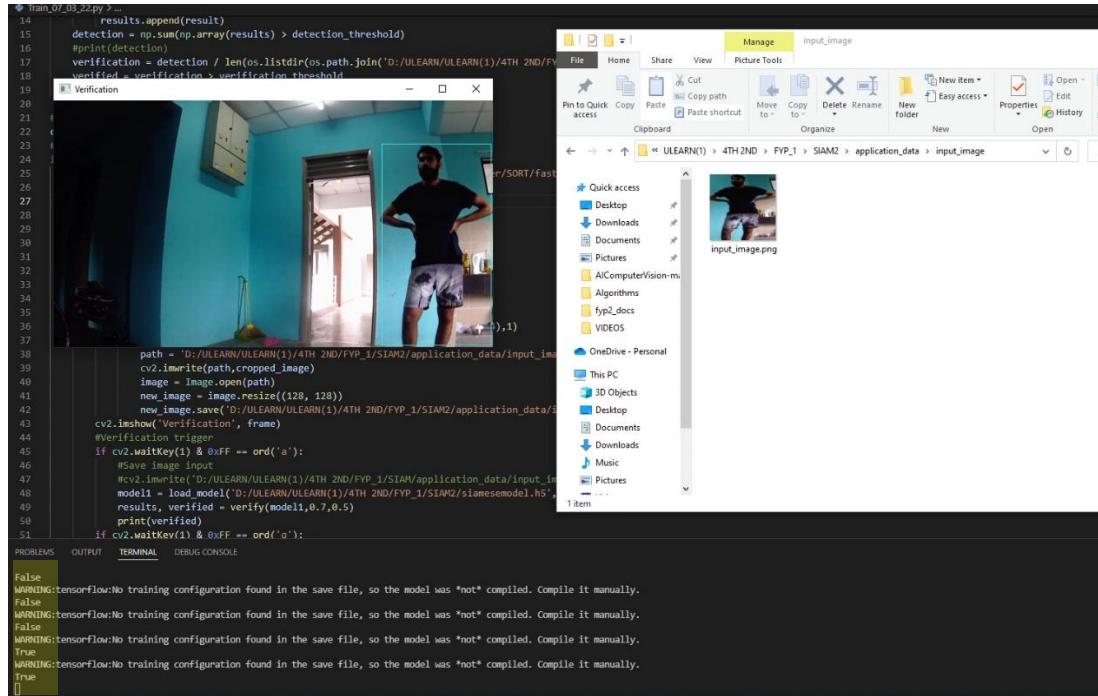


**FIGURE 4.30 – Video input from live stream, for testing (1)**

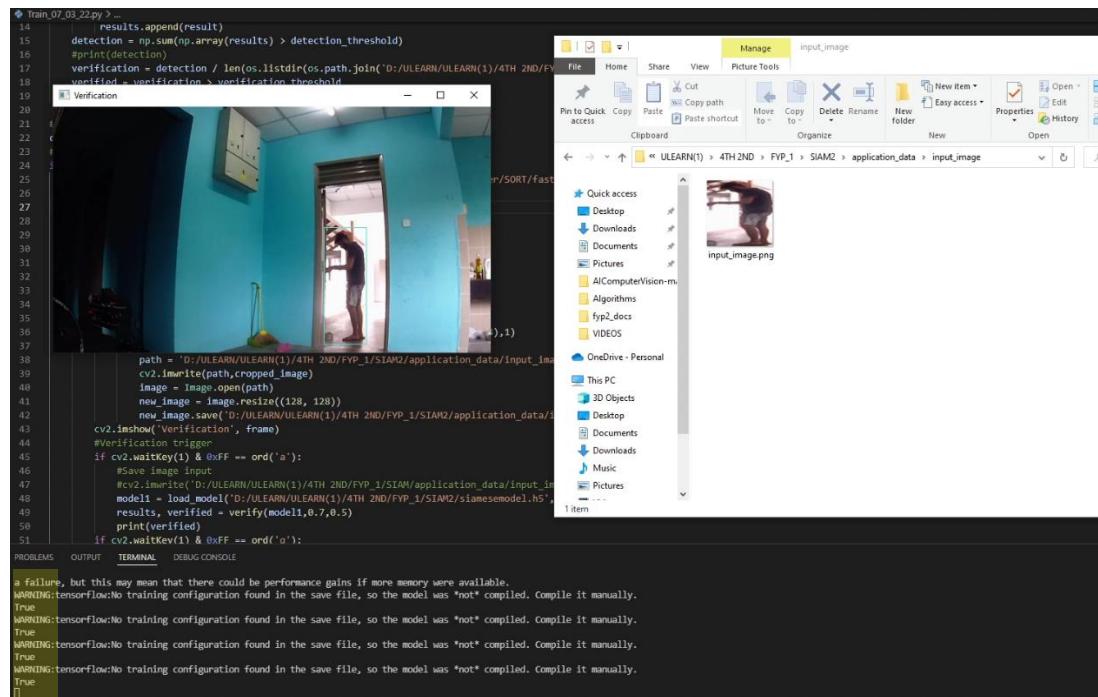


**FIGURE 4.31 – Video input from live stream, for testing (2)**

Now, the test was done by the same input, but from a different field of view. There are some false detections, but the overall performance is satisfactory. To tweak things up, the door was purposely let open to determine whether or not the Siamese Model trained can conclude that the person detected is the exact person used in training. The results were highlighted in yellow. Take a look in the figures 4.32 and 4.33 below.



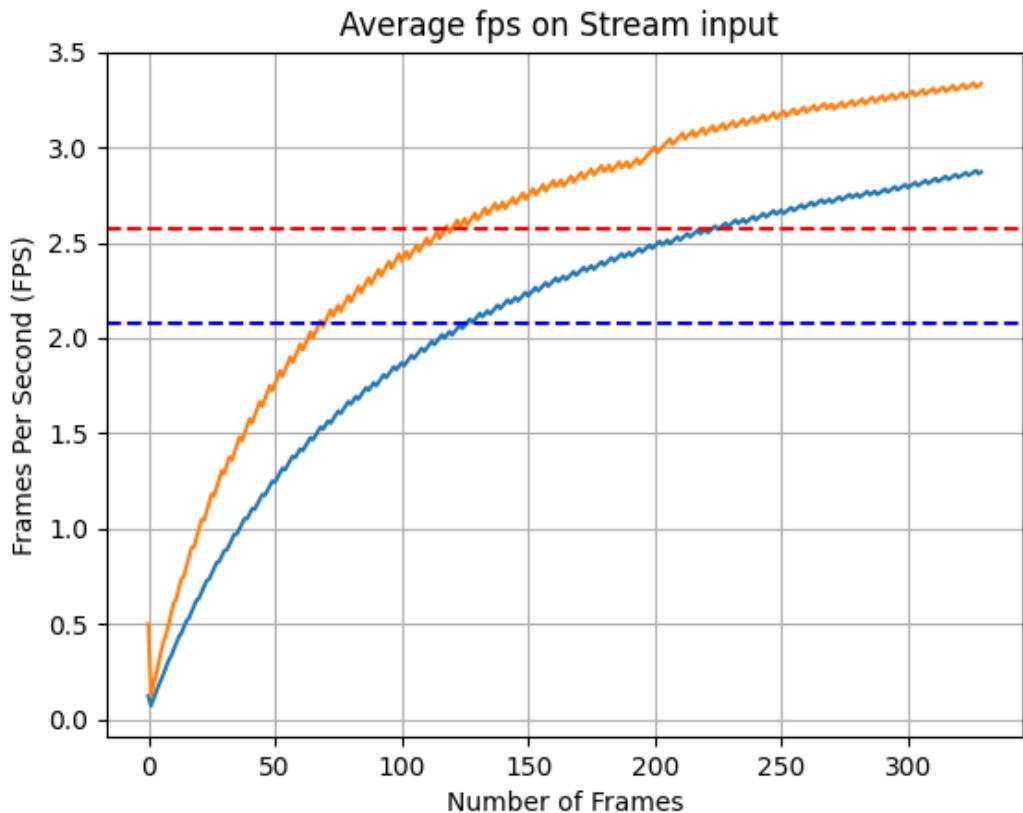
**FIGURE 4.32 – Video input from live stream, from a different angle (1)**



**FIGURE 4.33 – Video input from live stream, from a different angle (2)**

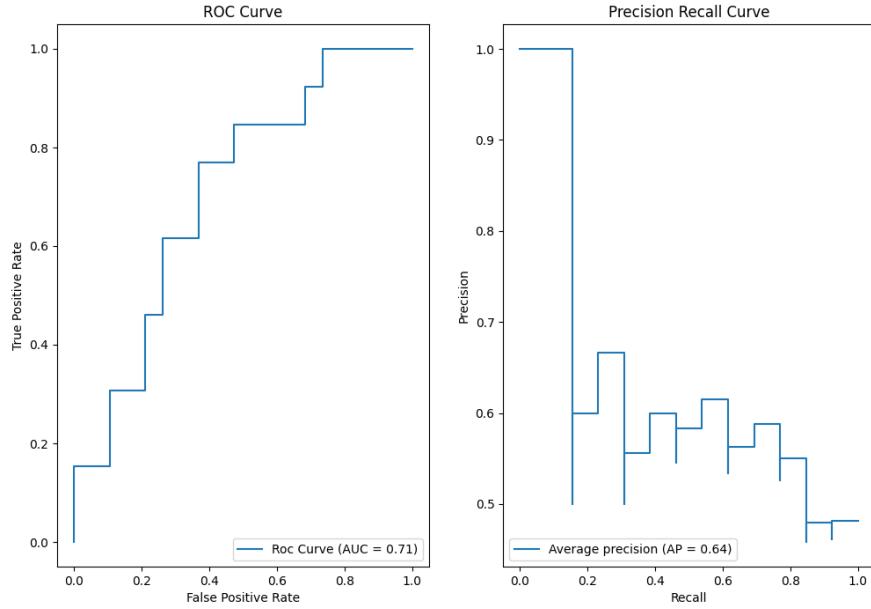
The average frames per second recorded from these two stream inputs were approximately, 2.08 and 2.569 respectively. Although, these results does not differ as much as the average frames per second recorded in the video input, there were some latency or delay that exist. This “*delay*” was not present in the pre-recorded video inputs. Hence it is safe to assume that this is due to the fact that the stream runs over the internet, accompanied by the use of hardware intensive object detection API and tracker that may have resulted this unprecedented delay. Since the project does not cover these aspects or issues, the root cause of the problem was left unexplored.

The figure 4.34 below represents the average fps on both the stream inputs. The average frames per second recorded was anywhere around 2.1 to 2.6fps. Similarly, the data were collected from 330 total number of frames, just to keep the measurements consistent for all scenarios. Since these were live streams from Hypertext Transfer Protocol (HTTP), some level of latency/delay was bound to happen, which may result in the decreased frame rate, if compared to an input from a pre-recorded video



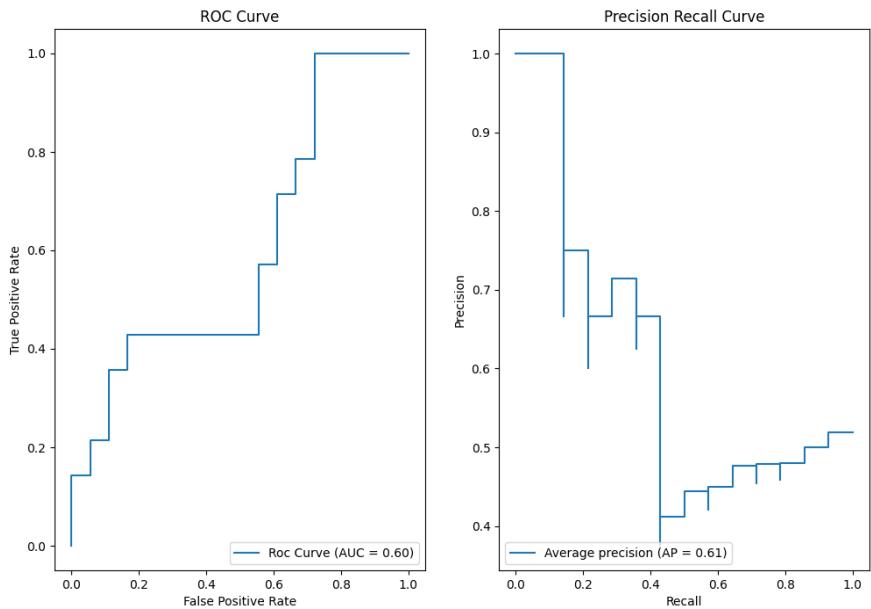
**FIGURE 4.34 – Average fps on stream input from two distinct angles**

Likewise, these results were obtained on two random batches, rather than the whole test sample. On the first batch, the area under the curve, AUC was 0.71, with an average precision in the PR curve of 0.64.



**FIGURE 4.35 – ROC Curve and Precision Recall curves from a stream (1)**

On the second random batch, the area under the curve, AUC was 0.60 with an average precision in the PR curve of 0.61. This model does a satisfactory job in identifying the person of interest, albeit not as accurate as the previously trained model. There could a few factors that could have resulted this; amount of light, occlusion, and the difference in the video resolution itself.



**FIGURE 4.36 – ROC Curve and Precision Recall curves from a stream (2)**

The third test was performed again by receiving a recorded video input from a phone camera. The figures 4.37 and 4.38 shows the test held in a bright open area with minor occlusion. Much similar to the second test, the third test was done along the hostel's corridor, accompanied with good lighting and ambience. This represents the first camera angle.

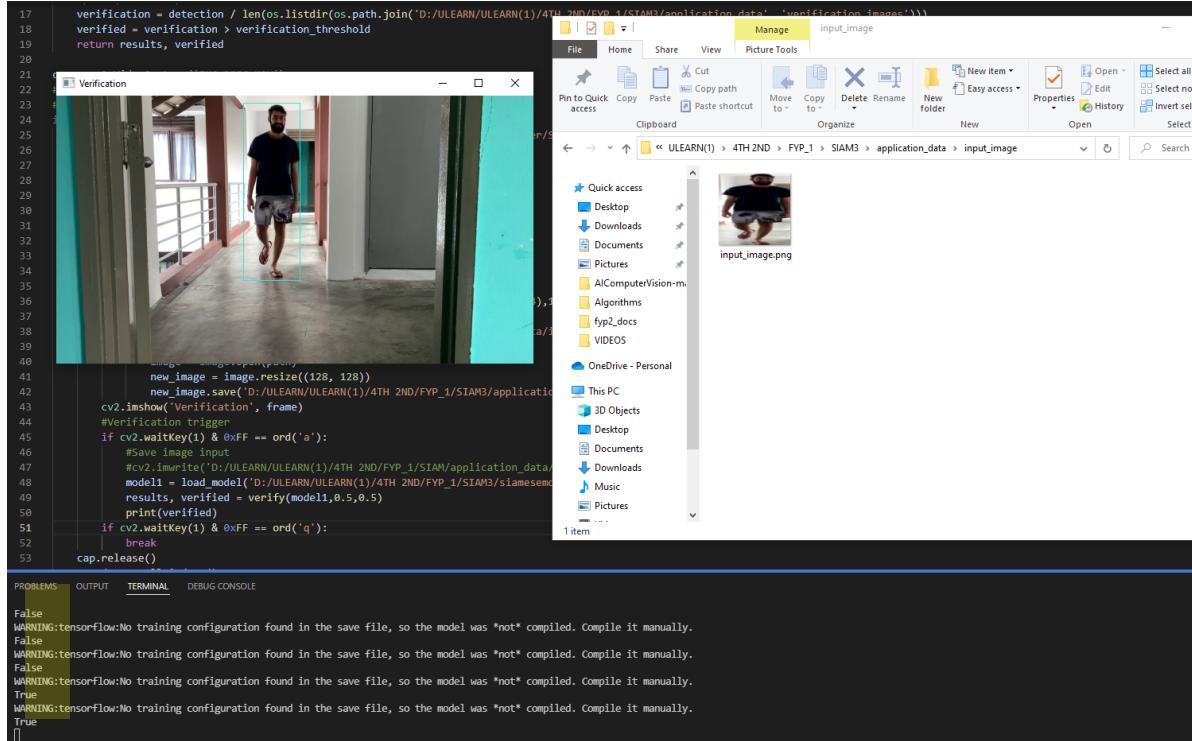


```

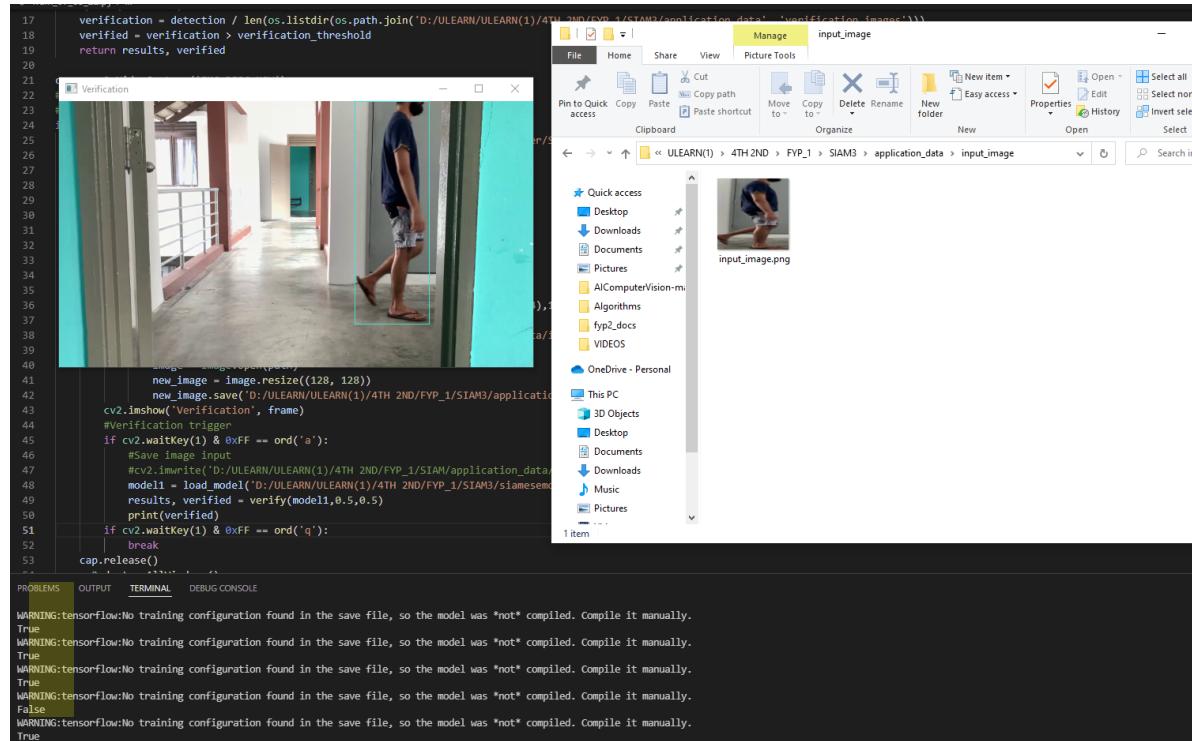
67     fps_start_time = datetime.datetime.now()
68     fps = 0
69     total_frames = 0
70     fps2 = []
71
72     model_path = 'D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/FYP2/AIComputerVision-master/SORT/faster_rcnn_inception_v2/frozen_inference_graph.pb'
73     odapi = DetectorAPI(path_to_cptk=model_path)
74     threshold = 0.5
75     #cap = cv2.VideoCapture('http://admin:0BthMuSP@192.168.0.30/VIDEO.CGI')
76     #cap = cv2.VideoCapture('test_video.mp4')
77     cap = cv2.VideoCapture('IMG_3585.MOV')
78     height = int(cap.get(cv2.CAP_PROP_FRAME_HEIGHT))
79     width = int(cap.get(cv2.CAP_PROP_FRAME_WIDTH))
80     print("Height:", height)
81     print("Width", width)
82
83
84     while True:
85         r, img = cap.read()
86         img = imutils.resize(img, width=width)
87         total_frames = total_frames + 1
88         fps_end_time = datetime.datetime.now()
89         time_diff = fps_end_time - fps_start_time
90
91         if time_diff.seconds == 0:
92             fps = 0.0
93         else:
94             fps = (total_frames / time_diff.seconds)
95         fps2.append(fps)
96
97         boxes, scores, classes, num
98
99         # Visualization of the results of a detection.
100        for i in range(len(boxes)):
101            # Class 1 represents human
102            if classes[i] == 1 and scores[i] > threshold:
103                box = boxes[i]
104                objectID = "ID {}".format(i)
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
709
710
711
712
713
714
715
716
717
718
719
719
720
721
722
723
724
725
726
727
728
729
729
730
731
732
733
734
735
736
737
737
738
739
739
740
741
742
743
744
745
745
746
747
747
748
749
749
750
751
752
753
753
754
755
755
756
756
757
757
758
758
759
759
760
760
761
761
762
762
763
763
764
764
765
765
766
766
767
767
768
768
769
769
770
770
771
771
772
772
773
773
774
774
775
775
776
776
777
777
778
778
779
779
780
780
781
781
782
782
783
783
784
784
785
785
786
786
787
787
788
788
789
789
790
790
791
791
792
792
793
793
794
794
795
795
796
796
797
797
798
798
799
799
800
800
801
801
802
802
803
803
804
804
805
805
806
806
807
807
808
808
809
809
810
810
811
811
812
812
813
813
814
814
815
815
816
816
817
817
818
818
819
819
820
820
821
821
822
822
823
823
824
824
825
825
826
826
827
827
828
828
829
829
830
830
831
831
832
832
833
833
834
834
835
835
836
836
837
837
838
838
839
839
840
840
841
841
842
842
843
843
844
844
845
845
846
846
847
847
848
848
849
849
850
850
851
851
852
852
853
853
854
854
855
855
856
856
857
857
858
858
859
859
860
860
861
861
862
862
863
863
864
864
865
865
866
866
867
867
868
868
869
869
870
870
871
871
872
872
873
873
874
874
875
875
876
876
877
877
878
878
879
879
880
880
881
881
882
882
883
883
884
884
885
885
886
886
887
887
888
888
889
889
890
890
891
891
892
892
893
893
894
894
895
895
896
896
897
897
898
898
899
899
900
900
901
901
902
902
903
903
904
904
905
905
906
906
907
907
908
908
909
909
910
910
911
911
912
912
913
913
914
914
915
915
916
916
917
917
918
918
919
919
920
920
921
921
922
922
923
923
924
924
925
925
926
926
927
927
928
928
929
929
930
930
931
931
932
932
933
933
934
934
935
935
936
936
937
937
938
938
939
939
940
940
941
941
942
942
943
943
944
944
945
945
946
946
947
947
948
948
949
949
950
950
951
951
952
952
953
953
954
954
955
955
956
956
957
957
958
958
959
959
960
960
961
961
962
962
963
963
964
964
965
965
966
966
967
967
968
968
969
969
970
970
971
971
972
972
973
973
974
974
975
975
976
976
977
977
978
978
979
979
980
980
981
981
982
982
983
983
984
984
985
985
986
986
987
987
988
988
989
989
990
990
991
991
992
992
993
993
994
994
995
995
996
996
997
997
998
998
999
999
1000
1000
1001
1001
1002
1002
1003
1003
1004
1004
1005
1005
1006
1006
1007
1007
1008
1008
1009
1009
1010
1010
1011
1011
1012
1012
1013
1013
1014
1014
1015
1015
1016
1016
1017
1017
1018
1018
1019
1019
1020
1020
1021
1021
1022
1022
1023
1023
1024
1024
1025
1025
1026
1026
1027
1027
1028
1028
1029
1029
1030
1030
1031
1031
1032
1032
1033
1033
1034
1034
1035
1035
1036
1036
1037
1037
1038
1038
1039
1039
1040
1040
1041
1041
1042
1042
1043
1043
1044
1044
1045
1045
1046
1046
1047
1047
1048
1048
1049
1049
1050
1050
1051
1051
1052
1052
1053
1053
1054
1054
1055
1055
1056
1056
1057
1057
1058
1058
1059
1059
1060
1060
1061
1061
1062
1062
1063
1063
1064
1064
1065
1065
1066
1066
1067
1067
1068
1068
1069
1069
1070
1070
1071
1071
1072
1072
1073
1073
1074
1074
1075
1075
1076
1076
1077
1077
1078
1078
1079
1079
1080
1080
1081
1081
1082
1082
1083
1083
1084
1084
1085
1085
1086
1086
1087
1087
1088
1088
1089
1089
1090
1090
1091
1091
1092
1092
1093
1093
1094
1094
1095
1095
1096
1096
1097
1097
1098
1098
1099
1099
1100
1100
1101
1101
1102
1102
1103
1103
1104
1104
1105
1105
1106
1106
1107
1107
1108
1108
1109
1109
1110
1110
1111
1111
1112
1112
1113
1113
1114
1114
1115
1115
1116
1116
1117
1117
1118
1118
1119
1119
1120
1120
1121
1121
1122
1122
1123
1123
1124
1124
1125
1125
1126
1126
1127
1127
1128
1128
1129
1129
1130
1130
1131
1131
1132
1132
1133
1133
1134
1134
1135
1135
1136
1136
1137
1137
1138
1138
1139
1139
1140
1140
1141
1141
1142
1142
1143
1143
1144
1144
1145
1145
1146
1146
1147
1147
1148
1148
1149
1149
1150
1150
1151
1151
1152
1152
1153
1153
1154
1154
1155
1155
1156
1156
1157
1157
1158
1158
1159
1159
1160
1160
1161
1161
1162
1162
1163
1163
1164
1164
1165
1165
1166
1166
1167
1167
1168
1168
1169
1169
1170
1170
1171
1171
1172
1172
1173
1173
1174
1174
1175
1175
1176
1176
1177
1177
1178
1178
1179
1179
1180
1180
1181
1181
1182
1182
1183
1183
1184
1184
1185
1185
1186
1186
1187
1187
1188
1188
1189
1189
1190
1190
1191
1191
1192
1192
1193
1193
1194
1194
1195
1195
1196
1196
1197
1197
1198
1198
1199
1199
1200
1200
1201
1201
1202
1202
1203
1203
1204
1204
1205
1205
1206
1206
1207
1207
1208
1208
1209
1209
1210
1210
1211
1211
1212
1212
1213
1213
1214
1214
1215
1215
1216
1216
1217
1217
1218
1218
1219
1219
1220
1220
1221
1221
1222
1222
1223
1223
1224
1224
1225
1225
1226
1226
1227
1227
1228
1228
1229
1229
1230
1230
1231
1231
1232
1232
1233
1233
1234
1234
1235
1235
1236
1236
1237
1237
1238
1238
1239
1239
1240
1240
1241
1241
1242
1242
1243
1243
1244
1244
1245
1245
1246
1246
1247
1247
1248
1248
1249
1249
1250
1250
1251
1251
1252
1252
1253
1253
1254
1254
1255
1255
1256
1256
1257
1257
1258
1258
1259
1259
1260
1260
1261
1261
1262
1262
1263
1263
1264
1264
1265
1265
1266
1266
1267
1267
1268
1268
1269
1269
1270
1270
1271
1271
1272
1272
1273
1273
1274
1274
1275
1275
1276
1276
1277
1277
1278
1278
1279
1279
1280
1280
1281
1281
1282
1282
1283
1283
1284
1284
1285
1285
1286
1286
1287
1287
1288
1288
1289
1289
1290
1290
1291
1291
1292
1292
1293
1293
1294
1294
1295
1295
1296
1296
1297
1297
1298
1298
1299
1299
1300
1300
1301
1301
1302
1302
1303
1303
1304
1304
1305
1305
1306
1306
1307
1307
1308
1308
1309
1309
1310
1310
1311
1311
1312
1312
1313
1313
1314
1314
1315
1315
1316
1316
1317
1317
1318
1318
1319
1319
1320
1320
1321
1321
1322
1322
1323
1323
1324
1324
1325
1325
1326
1326
1327
1327
1328
1328
1329
1329
1330
1330
1331
1331
1332
1332
1333
1333
1334
1334
1335
1335
1336
1336
1337
1337
1338
1338
1339
1339
1340
1340
1341
1341
1342
1342
1343
1343
1344
1344
1345
1345
1346
1346
1347
1347
1348
1348
1349
1349
1350
1350
1351
1351
1352
1352
1353
1353
1354
1354
1355
1355
1356
1356
1357
1357
1358
1358
1359
1359
1360
1360
1361
1361
1362
1362
1363
1363
1364
1364
1365
1365
1366
1366
1367
1367
1368
1368
1369
1369
1370
1370
1371
1371
1372
1372
1373
1373
1374
1374
1375
1375
1376
1376
1377
1377
1378
1378
1379
1379
1380
1380
1381
1381
1382
1382
1383
1383
1384
1384
1385
1385
1386
1386
1387
1387
1388
1388
1389
1389
1390
1390
1391
1391
1392
1392
1393
1393
1394
1394
1395
1395
1396
1396
1397
1397
1398
1398
1399
1399
1400
1400
1401
1401
1402
1402
1403
1403
1404
1404
1405
1405
1406
1406
1407
1407
1408
1408
1409
1409
1410
1410
1411
1411
1412
1412
1413
1413
1414
1414
1415
1415
1416
1416
1417
1417
1418
1418
1419
1419
1420
1420
1421
1421
1422
1422
1423
1423
1424
1424
1425
1425
1426
1426
1427
1427
1428
1428
1429
1429
1430
1430
1431
1431
1432
1432
1433
1433
1434
1434
1435
1435
1436
1436
1437
1437
1438
1438
1439
1439
1440
1440
1441
1441
1442
1442
1443
1443
1444
1444
1445
1445
1446
1446
1447
1447
1448
1448
1449
1449
1450
1450
1451
1451
1452
1452
1453
1453
1454
1454
1455
1455
1456
1456
1457
1457
1458
1458
1459
1459
1460
1460
1461
1461
1462
1462
1463
1463
1464
1464
1465
1465
1466
1466
1467
1467
1468
1468
1469
1469
1470
1470
1471
1471
1472
1472
1473
1473
1474
1474
1475
1475
1476
1476
1477
1477
1478
1478
1479
1479
1480
1480
1481
1481
1482
1482
1483
1483
1484
1484
1485
1485
1486
1486
1487
1487
1488
1488
1489
1489
1490
1490
1491
1491
1492
1492
1493
1493
1494
1494
1495
1495
1496
1496
1497
1497
1498
1498
1499
1499
1500
1500
1501
1501
1502
1502
1503
1503
1504
1504
1505
1505
1506
1506
1507
1507
1508
1508
1509
1509
1510
1510
1511
1511
1512
1512
1513
1513
1514
1514
1515
1515
1516
1516
1517
1517
1518
1518
1519
1519
1520
1520
1521
1521
1522
1522
1523
1523
1524
1524
1525
1525
1526
1526
1527
1527
1528
1528
1529
1529
1530
1530
1531
1531
1532
1532
1533
1533
1534
1534
1535
1535
1536
1536
1537
1537
1538
1538
1539
1539
1540
1540
1541
1541
1542
1542
1543
1543
1544
1544
1545
1545
1546
1546
1547
1547
1548
1548
1549
1549
1550
1550
1551
1551
1552
1552
1553
1553
1554
1554
1555
1555
1556
1556
1557
1557
1558
1558
1559
1559
1560
1560
1561
1561
1562
1562
1563
1563
1564
1564
1565
1565
1566
1566
1567
1567
1568
1568
1569
1569
1570
1570
1571
1571
1572
1572
1573
1573
1574
1574
1575
1575
1576
1576
1577
1577
1578
1578
1579
1579
1580
1580
1581
1581
1582
1582
1583
1583
1584
1584
1585
1585
1586
1586
1587
1587
1588
1588
1589
1589
1590
1590
1591
1591
1592
1592
1593
1593
1594

```

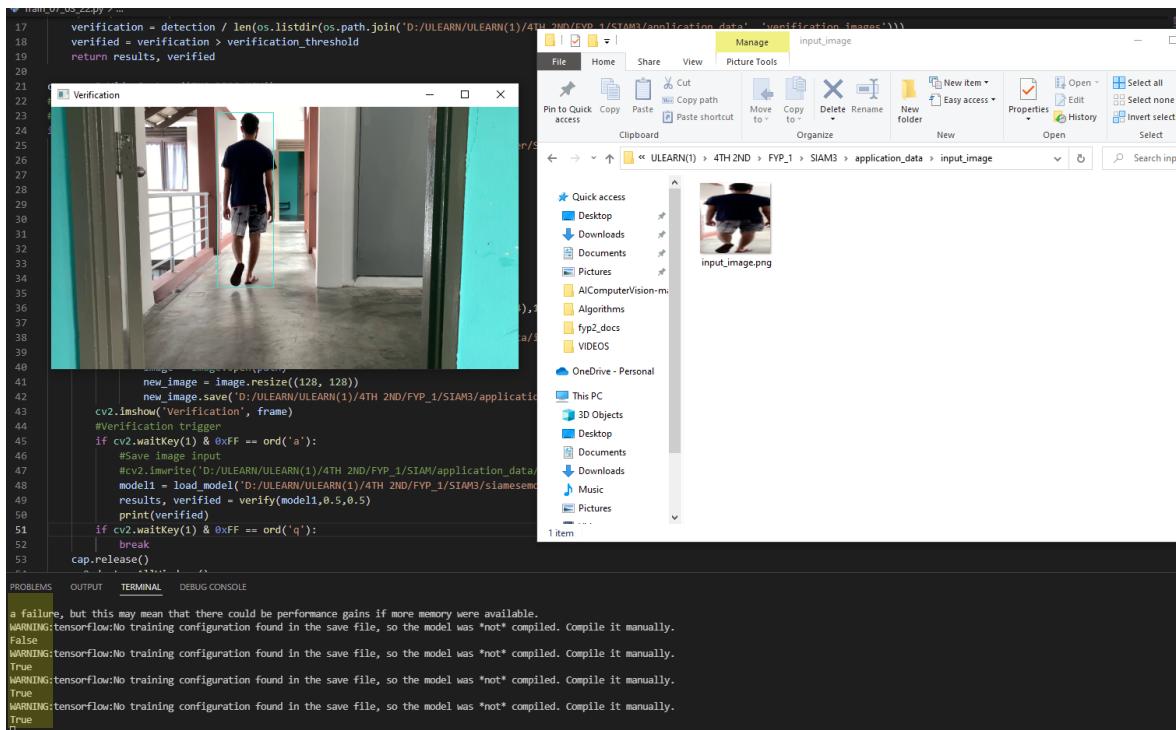
Now, since the video input was taken in the corridor (first camera angle), it would consecutively be used for training. Have a look at the figures 4.39, 4.40 and 4.41. The developed Siamese Model can identify the person of interest as the same person detected from the first video input, although this time the camera was placed on a different field of view.



**FIGURE 4.39 –Input from recorded video, from a different angle (1)**

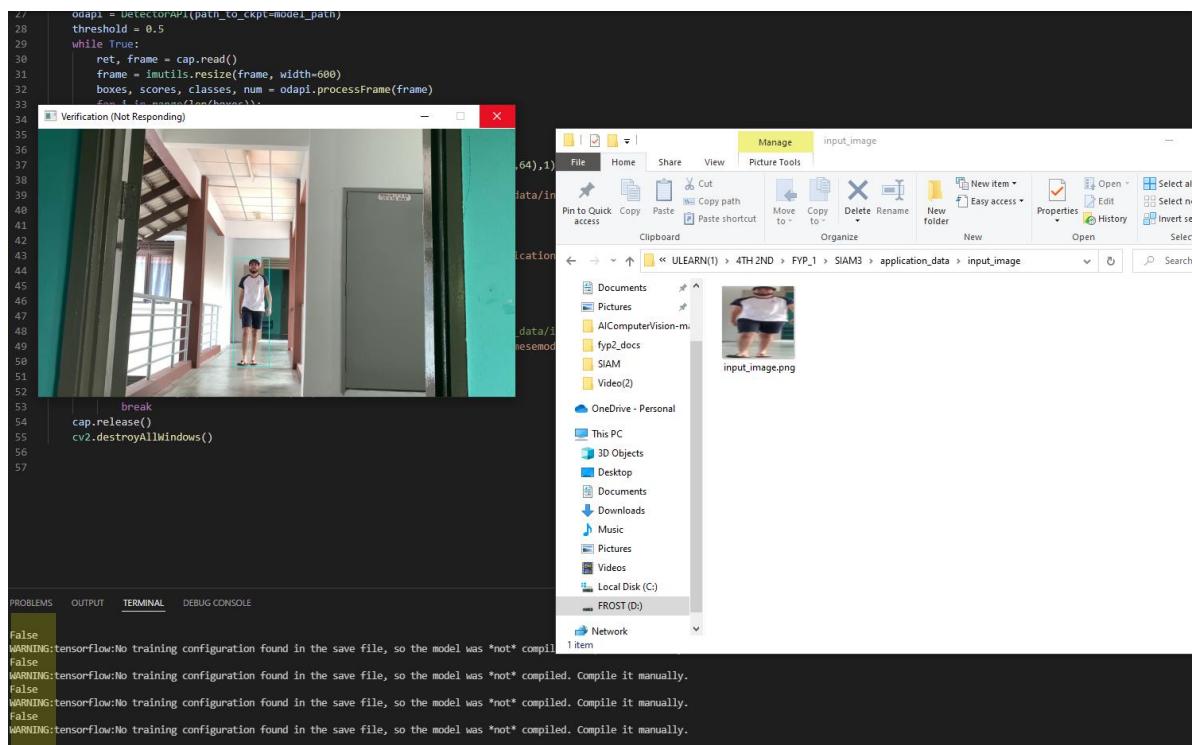


**FIGURE 4.40 –Input from recorded video, from a different angle (2)**



**FIGURE 4.41 – Input from recorded video, from a different angle (3)**

Generally, the Siamese model did a good job in predicting that this person, in the second video input is the same person detected in the first video input. In simpler words, the network was able to re-identify the person from a non-overlapping camera, that is of a different angle. Look at figures 4.42 and 4.43; was the same field of view, but with a different person present.



**FIGURE 4.42 – Input from recorded video, with a different person (1)**

The screenshot shows a code editor with Python code for a Siamese model. The code reads frames from a video, processes them with a detector API, and checks for verification. A terminal window shows TensorFlow compilation warnings. To the right, a file explorer window displays an image named 'input\_image.png' which is a screenshot of a hallway scene. A second window titled 'Verification (Not Responding)' shows a similar hallway scene with a bounding box around a person's legs.

```

21 odapi = DetectorAPI(path_to_ckpt=model1_path)
22 threshold = 0.5
23 while True:
24     ret, frame = cap.read()
25     frame = imutils.resize(frame, width=600)
26     boxes, scores, classes, num = odapi.processFrame(frame)
27     for i in range(len(bboxes)):
28         if scores[i] > threshold:
29             print("Verification trigger")
30             if cv2.waitKey(1) & 0xFF == ord('a'):
31                 #Save image input
32                 #cv2.imwrite('Verification', frame)
33                 #model1 = load_model('D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/SIAM/application_data/model1.h5')
34                 #results, verified = verify(model1,0.7,0.5)
35                 #print(verified)
36
37             break
38         cap.release()
39         cv2.destroyAllWindows()

```

**FIGURE 4.43 – Input from recorded video, from a different angle (2)**

With a different person present, in the same camera angle, the Siamese model was able to make good predictions and concluded that this was not the same individual detected in the first camera. Now, have a look at figures 4.44, 4.45 and 4.46 instead.

The screenshot shows a code editor with Python code for a Siamese model. The code reads frames from a video, processes them with a detector API, and checks for verification. A terminal window shows TensorFlow compilation warnings. To the right, a file explorer window displays an image named 'input\_image.png' which is a screenshot of a hallway scene with a person walking away from the camera. A second window titled 'Verification (Not Responding)' shows a similar hallway scene with a bounding box around a person's legs.

```

14     results.append(result)
15     detection = np.sum(np.array(results) > detection_threshold)
16     #print(detection)
17     verification = detection / len(os.listdir(os.path.join('D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/SIAM/application_data', 'input_image')))
18     verified = verification > verification_threshold
19     return results, verified
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41     image = Image.open(path)
42     new_image = image.resize((128, 128))
43     new_image.save('D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/SIAM/application_data/input_image.png')
44     cv2.imshow('Verification', frame)
45     #Verification trigger
46     if cv2.waitKey(1) & 0xFF == ord('a'):
47         #Save image input
48         #cv2.imwrite('Verification', frame)
49         #model1 = load_model('D:/ULEARN/ULEARN(1)/4TH 2ND/FYP_1/SIAM/application_data/model1.h5')
50         #results, verified = verify(model1,0.7,0.5)
51         #print(verified)

```

**FIGURE 4.44 – Input from recorded video, with a different lighting (1)**

To further question the legitimacy of the predicted model, the test was done on the same angle, with an unfamiliar person, once again but at a different lighting situation. The return results from the verification were **False**, indicating that this is not the person detected from the first camera input. Conclusively, the neural network has done an excellent job in reidentifying whether the individual tracked in the second camera angle, is the same person or not.

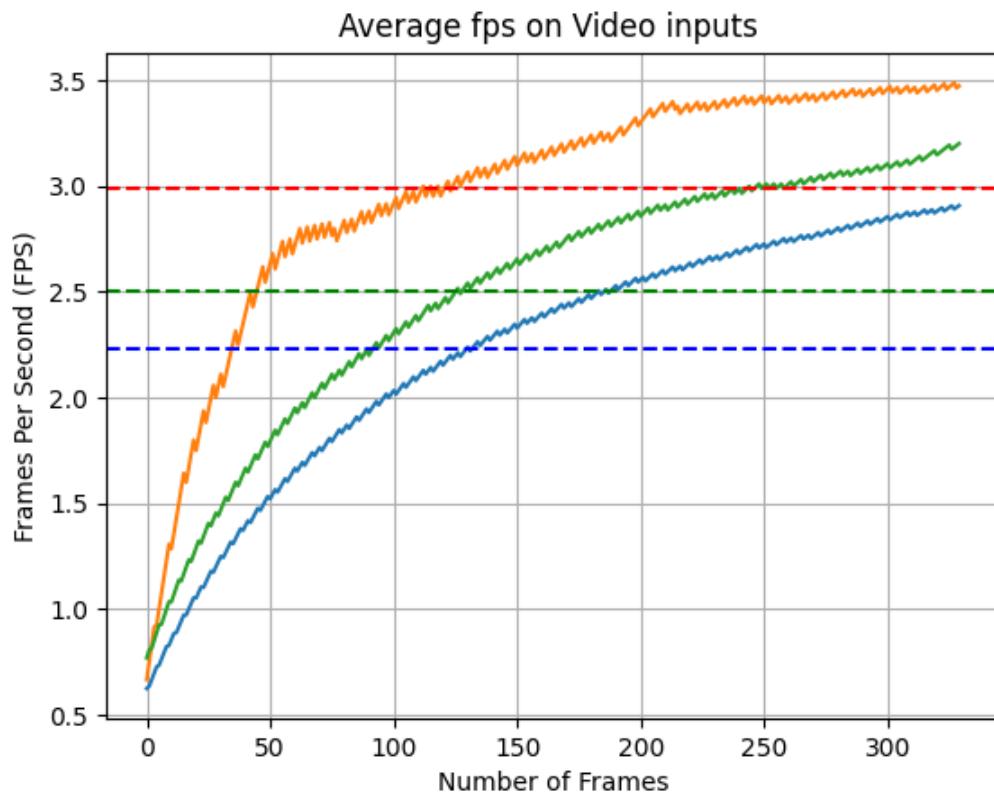
The screenshot shows a development environment with several windows open:

- Code Editor:** Displays Python code for a verification script. The code includes imports for numpy, cv2, and tensorflow, as well as logic for reading an image, performing a verification threshold check, and saving the result.
- File Explorer:** Shows a folder structure under 'input\_image'. It contains subfolders like 'Downloads', 'Documents', 'Pictures', and 'OneDrive - Personal', and files such as 'input\_image.png'.
- Terminal:** Displays TensorFlow logs indicating that training configuration was not found in the save file, so the mode was not compiled and needs to be done manually.

**FIGURE 4.45 – Input from recorded video, with a different lighting (2)**

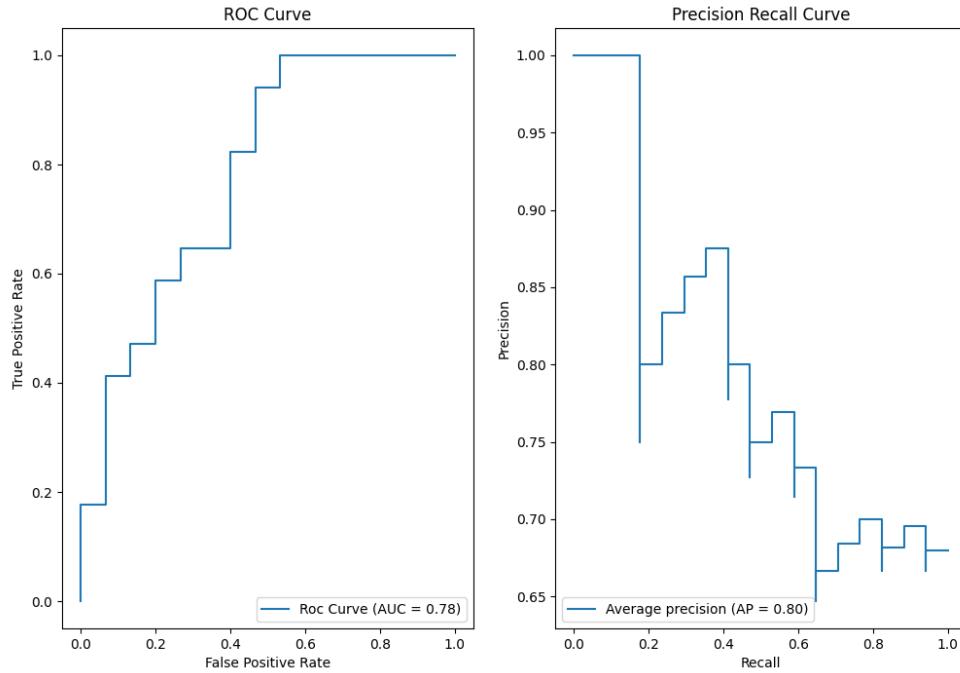
**FIGURE 4.46** – Input from recorded video, with a different person (3)

The figure 4.45 below shows the average fps on all three of the pre-recorded video inputs. The frames spikes up from the very beginning and level out evenly around the average of 2.22, 2.98 and 2.50 fps respectively. As for all the other inputs, these data were also collected from a 330 total number of frames.



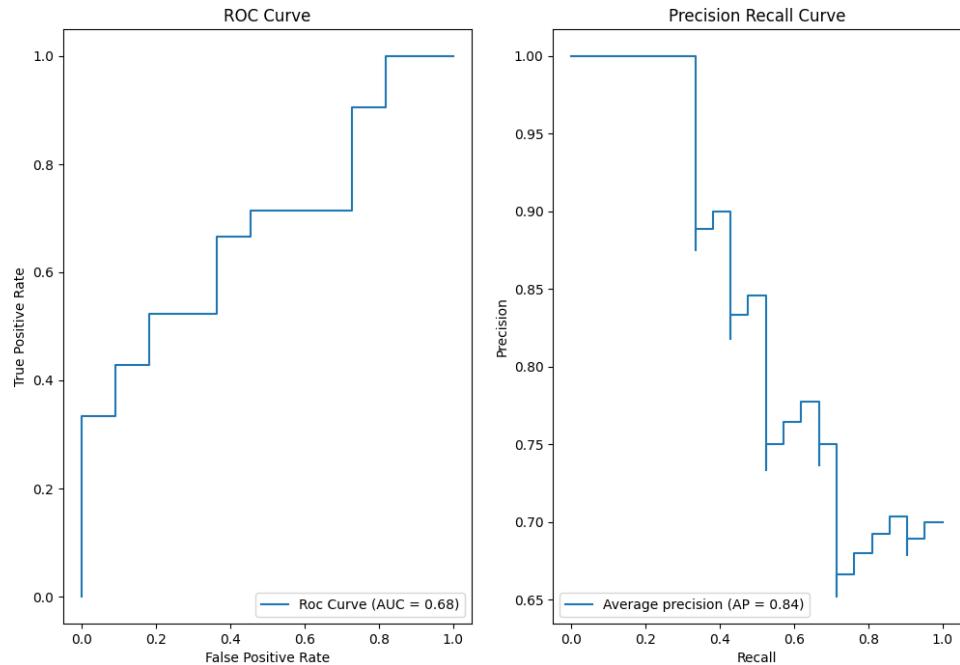
**FIGURE 4.47** - Average fps on pre-recorded video inputs from two distinct angles

On the first random batch, AUC was 0.78, with an average precision in the PR curve of 0.80. This shows that the number of false positives and negatives were little.



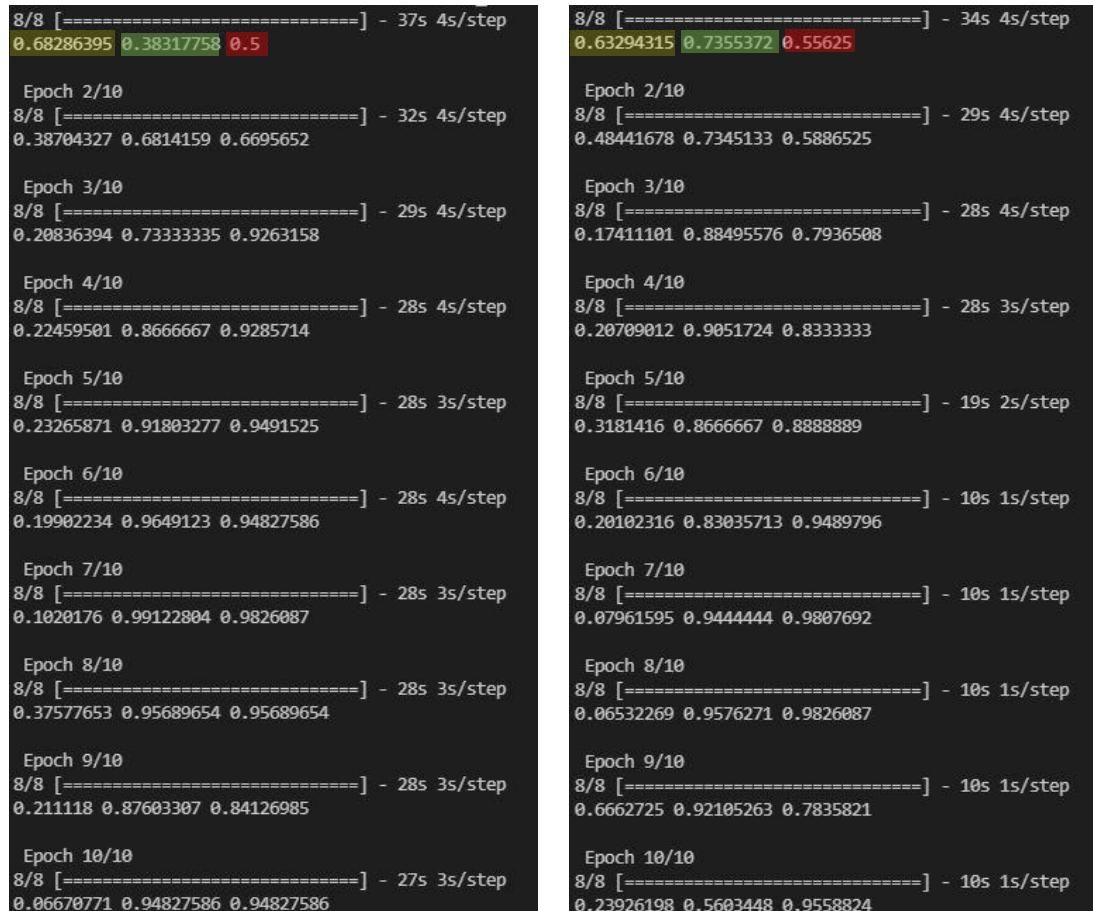
**FIGURE 4.48 – ROC Curve and Precision Recall curves from the pre-recorded videos (1)**

On a second random batch, the area under the curve, AUC was 0.68, with an average precision in the PR curve of 0.84. Similarly, this model also performed well in identifying the person of interest against other people.



**FIGURE 4.49 – ROC Curve and Precision Recall curves from the pre-recorded videos (2)**

Since all the models take exactly 10 epochs for training, the difference between the time taken does not differ as much, regardless the type of input given. That being said, the average time taken to train the neural networks were approximately 190 to 310 seconds. As the number of epochs increases, the values of precision and recall, tends to get closer to 1, while loss, would get closer to zero. Loss is highlighted in yellow, green for precision and red for recall. Have a look at figure 4.50 below.



**FIGURE 4.50** - Average time taken to run 10 epochs

# CHAPTER 5: CONCLUSION AND RECOMMENDATIONS

## 5.1 Conclusion

This project represents a framework that implements a pedestrian re-identification network, which by itself may provide a wide range of advantages. With the urgent demand for public safety and an increasing number of surveillance cameras, the re-identification of people is of great importance. The development of such system is undeniably cost-effective; doesn't require the need of actual humans to do constant monitoring at all times. Hence, with the development of distributed video tracking the various methods for capturing parameters for the person of interest could be comprehended, and an existing network in a real-world scenario and on a new data can be verified and evaluated.

Thus far the project can only perform verification on a single person, from a given video input. This is due to the way the model was built upon, which lacks in the number of input images it can verify. Since the input image are detected by using the coordinates of the bounding box, the issue arrives when there more than a single person in a particular frame. In other words, during verification, the system cannot crop multiple images of multiple person in any given frame, to perform the verification. Hence, it is virtually impossible to point out the exact IDs of individuals, across two non-overlapping camera angles.

Hence, to further prevent complications as such, the whole project was built to ensure that only a single person is detected in the first, and second camera angles. This developed project would also work if there were multiple people in the frames as well, but instead of returning the IDs of the corresponding individuals, the system can only retrieve Boolean values of **True** or **False**. This means that, the neural network can identify whether or not all these people detected in the first camera angle, has reappeared in the second camera angle. However, it cannot distinguish the individuals of the already detected IDs and reassign them respectively.

Perhaps with more time and extensive testing, the system would be able to perform much better, by differentiating IDs which were assigned in the first camera, to simultaneously be detected in the second camera angle.

## 5.2 Recommendation

For developing such systems, it's recommended to have the best in-line hardware, to be able to cope up with the extensive computation. This is especially useful if the training and deployment were to be done locally, and not over the internet. Contrarily, with NVIDIA Tesla T4 GPU with 16GB (gigabytes) of video memory on Google Colab would allow for quicker response time, and hence faster performance and higher frame rate. Besides, instead of using a streaming protocol like RTSP or HTTP, it is better to have cameras physically connected to a particular system to minimize overall delay.

It is also important to have footages recorded or streamed in a particularly bright area, with minor occlusion as it allows for far more excellent model training. The neural network design learns data point similarity using Siamese networks. Since the Siamese Neural Network make predictions based on images, it is best to ensure that the resolution and overall quality of an input video to be high, be it a live stream or from a recorded video.

## CHAPTER 6: REFERENCES

- [1] Klingler, N., 2021. *Deep Learning for Person Re-Identification - viso.ai*. [online] viso.ai. Available at: <<https://viso.ai/deep-learning/deep-learning-for-person-re-identification/>> [Accessed 1 October 2021].
- [2] Sharabok, G., 2021. Why we need person re-identification?. [online] Medium. Available at: <<https://towardsdatascience.com/why-we-need-person-re-identification-3a45d170098b>> [Accessed 1 October 2021].
- [3] PETER DANIELSON, M., 2021. *Monitoring and Surveillance / Encyclopedia.com*. [online] Encyclopedia.com. Available at: <<https://www.encyclopedia.com/science/encyclopedias-almanacs-transcripts-and-maps/monitoring-and-surveillance>> [Accessed 12 November 2021].
- [4] A. McHale, Director, and Memoori Research, “2019 surpassed expectations but demand structure in video surveillance is unbalanced,” *asmag.com*. [Online]. Available: <https://www.asmag.com/showpost/30652.aspx>. [Accessed: 12-Nov-2021].
- [5] S. Feldstein, The global expansion of AI surveillance, vol. 17. Washington, DC: Carnegie Endowment for International Peace, 2019.
- [6] J. Brownlee, “A gentle introduction to object recognition with deep learning,” Machinelearningmastery.com, 21-May-2019. [Online]. Available: <https://machinelearningmastery.com/object-recognition-with-deep-learning/>. [Accessed: 12-Nov-2021].
- [7] L. Wu, C. Shen and A. van den Hengel, "PersonNet: Person Reidentification with Deep Convolutional Neural Networks", Computer Vision and Pattern Recognition, 2016
- [8] L. Zhang, T. Xiang, and S. Gong, "Learning a discriminative null space for person re-Identification," CVPR, 2016.
- [9] F. Xiong, M. Gou, O. Camps and M. Sznajer, "Person Re-Identification Using Kernel-Based Metric Learning Methods", Computer Vision – ECCV 2014, pp. 1-16, 2014.

- [10] S. Liao, Y. Hu, Xiangyu Zhu and S. Li, "Person re-identification by Local Maximal Occurrence representation and metric learning", 2015 IEEE International Conference on Signal and Image Processing Applications (IEEE ICSIPA 2019), IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2015.
- [11] T. Matsukawa, T. Okabe, E. Suzuki, and Y. Sato, "Hierarchical gaussian descriptor for person re-identification," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 1363–1372
- [12] E. H. Land and J. J. McCann, "Lightness and retinex theory," *J. Opt. Soc. Am.*, vol. 61, no. 1, pp. 1–11, 1971.13)
- [13] S. Salehian, P. Sebastian, and A. B. Sayuti, "Framework for pedestrian detection, tracking and re-identification in video surveillance system," in *2019 IEEE International Conference on Signal and Image Processing Applications (ICSIPIA)*, 2019.
- [14] W. Li, R. Zhao, T. Xiao, and X. Wang, "DeepReID: Deep filter pairing neural network for person re-identification," in 2014 IEEE Conference on Computer Vision and Pattern Recognition, 2014..
- [15] Y. Li, S. Pirk, H. Su, C. R. Qi, and L. J. Guibas, "FPNN: Field probing Neural Networks for 3D data," *arXiv [cs.CV]*, 2016
- [16] "What is Object Tracking? - An Introduction (2021) - viso.ai," Viso.ai, 10-Jul-2021. [Online]. Available: <https://viso.ai/deep-learning/object-tracking/>. [Accessed: 12-Nov-2021].
- [17] "Papers with code - object tracking," Paperswithcode.com. [Online]. Available: <https://paperswithcode.com/task/object-tracking>. [Accessed: 12-Nov-2021].
- [18] A. Lukežič, T. Vojíř, L. Čehovin Zajc, J. Matas, and M. Kristan, "Discriminative correlation filter tracker with channel and spatial reliability," *Int. J. Comput. Vis.*, vol. 126, no. 7, pp. 671–688, 2018.
- [19] A. Yelisetty, "Understanding object detection and R-CNN," Towards Data Science, 08-Jul-2020. [Online]. Available: <https://towardsdatascience.com/understanding-object-detection-and-r-cnn-e39c16f37600>. [Accessed: 12-Nov-2021].

- [20] “Object detection techniques in deep learning,” Inblog.in. [Online]. Available: <https://inblog.in/Object-Detection-Techniques-in-Deep-Learning-eTqP3Hneil>. [Accessed: 12-Nov-2021].
- [21] A. Yelisetty, “Understanding object detection and R-CNN,” *Towards Data Science*, 08-Jul-2020. [Online]. Available: <https://towardsdatascience.com/understanding-object-detection-and-r-cnn-e39c16f37600>. [Accessed: 12-Nov-2021].
- [22] R. Gandhi, “R-CNN, fast R-CNN, faster R-CNN, YOLO — object detection algorithms,” *Towards Data Science*, 09-Jul-2018. [Online]. Available: <https://towardsdatascience.com/r-cnn-fast-r-cnn-faster-r-cnn-yolo-object-detection-algorithms-36d53571365e>. [Accessed: 12-Nov-2021].
- [23] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards real-time object detection with region proposal networks,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [24] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards real-time object detection with region proposal networks,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [25] R. L. Galvez, A. A. Bandala, E. P. Dadios, R. R. P. Vicerra and J. M. Z. Maningo, "Object Detection Using Convolutional Neural Networks," TENCON 2018 - 2018 IEEE Region 10 Conference, 2018, pp. 2023-2027.
- [26] Kunar, A., 2021. Object Detection with SSD and MobileNet. [online] Medium. Available at: <<https://adityakunar.medium.com/object-detection-with-ssd-and-mobilenet-aeedc5917ad0>> [Accessed 26 October 2021].
- [27] He, K., Gkioxari, G., Dollar, P. and Girshick, R., 2020. Mask R-CNN. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(2), pp.386-397.
- [28] J. F. Villacrés and F. Auat Cheein, “Detection and characterization of cherries: A deep learning usability case study in Chile,” *Agronomy (Basel)*, vol. 10, no. 6, p. 835, 2020.
- [29] A. Yelisetty, “Understanding object detection and R-CNN,” *Towards Data Science*, 08-Jul-2020. [Online]. Available: <https://towardsdatascience.com/understanding-object-detection-and-r-cnn-e39c16f37600>. [Accessed: 12-Nov-2021].

- [30] W. Liu *et al.*, “SSD: Single Shot MultiBox Detector,” in *Computer Vision – ECCV 2016*, Cham: Springer International Publishing, 2016, pp. 21–37.
- [31] Rosebrock, A., 2021. Simple object tracking with OpenCV - PyImageSearch. [online] PyImageSearch. Available at: <<https://www.pyimagesearch.com/2018/07/23/simple-object-tracking-with-opencv/>> [Accessed 26 October 2021].
- [32] Medium. 2021. People Tracking using Deep Learning. [online] Available at: <<https://towardsdatascience.com/people-tracking-using-deep-learning-5c90d43774be>> [Accessed 26 October 2021].
- [33] Medium. 2021. People Tracking using Deep Learning. [online] Available at: <<https://towardsdatascience.com/people-tracking-using-deep-learning-5c90d43774be>> [Accessed 26 October 2021].
- [34] AI & Machine Learning Blog. 2021. DeepSORT: Deep Learning to track custom objects in a video. [online] Available at: <<https://nanonets.com/blog/object-tracking-deepsort/>> [Accessed 26 October 2021].
- [35] K. Host, M. Ivašić-Kos, and M. Pobar, “Tracking handball players with the DeepSORT algorithm,” in *Proceedings of the 9th International Conference on Pattern Recognition Applications and Methods*, 2020.
- [36] Q. Wang, J. Gao, J. Xing, M. Zhang, and W. Hu, “DCFNet: Discriminant Correlation Filters network for visual tracking,” *arXiv [cs.CV]*, 2017.
- [37] Liang, Feng & Yang, Shaofei & Mai, Tinrui & Yang, Yichen. (2018). The Design of Objects Bounding Boxes Non-Maximum Suppression and Visualization Module Based on FPGA. 1-5. 10.1109/ICDSP.2018.8631668.
- [38] D. Wang *et al.*, “Daedalus: Breaking Non-Maximum Suppression in Object Detection via adversarial examples,” *arXiv [cs.CV]*, 2019.
- [39] L. M. Brown, R. Feris, and S. Pankanti, “Temporal non-maximum suppression for pedestrian detection using self-calibration,” in *2014 22nd International Conference on Pattern Recognition*, 2014.
- [40] G. Koch, R. Zemel, en R. Salakhutdinov, “Siamese Neural Networks for One-shot Image Recognition”, 2015.

- [41] N. I. Mokhtari, “What are Siamese neural networks in deep learning?,” Medium, 15-Feb-2022. [Online]. Available: <https://towardsdatascience.com/what-are-siamese-neural-networks-in-deep-learning-bb092f749dcb>. [Accessed: 08-Mar-2022].
- [42] J. Brownlee, “A gentle introduction to the rectified linear unit (ReLU),” Machine Learning Mastery, 20-Aug-2020. [Online]. Available: <https://machinelearningmastery.com/rectified-linear-activation-function-for-deep-learning-neural-networks/#:~:text=The%20rectified%20linear%20activation%20function,otherwise%2C%20it%20will%20output%20zero>. [Accessed: 08-Mar-2022].
- [43] A. Bewley, Z. Ge, L. Ott, F. Ramos, en B. Upcroft, “Simple online and realtime tracking”, in 2016 IEEE International Conference on Image Processing (ICIP), 2016, bll 3464–3468.
- [44] N. Wojke, A. Bewley, en D. Paulus, “Simple Online and Realtime Tracking with a Deep Association Metric”, in 2017 IEEE International Conference on Image Processing (ICIP), 2017, bll 3645–3649.
- [45] N. Wojke en A. Bewley, “Deep Cosine Metric Learning for Person Re-identification”, in 2018 IEEE Winter Conference on Applications of Computer Vision (WACV), 2018, bll 748–756.