



ESCUELA TÉCNICA SUPERIOR DE INGENIERÍA DE  
TELECOMUNICACIÓN

MASTER OFICIAL EN VISIÓN ARTIFICIAL

**Master thesis**

# Visual people tracking with deep learning detection and feature tracking

Author: Marcos Pieras Sagardoy

Tutor: José María Cañas Plaza

Academic course 2016/2017



©2017 Marcos Pieras Sagardoy

Esta obra está distribuida bajo la licencia de  
“Reconocimiento-CompartirIgual 4.0 Internacional (CC BY-SA 4.0)”

de Creative Commons.

Para ver una copia de esta licencia, visite  
<http://creativecommons.org/licenses/by-sa/4.0/> o envíe  
una carta a Creative Commons, 171 Second Street, Suite 300,  
San Francisco, California 94105, USA.

# Acknowledgement

# Abstract

Deep learning has rised by drastic improvements over reigning approaches towards the hardest problems in Artificial intelligence (AI), massive investments from industry giants, and exponential growth in research publications. Deep learning is a tool inside the machine learning toolbox, the goal is to make machines learn.

In some areas of artifical vision, deep learning techniques have been very succesful, however, in the field of visual tracking are not yet mature, therefore we have developed the multiple people tracking algorithm with deep learning techniques. Thus, in this work we have designed and build a software component that uses the paradigm tracking-by-detection. We mixed deep learning techniques, with feature tracking, using the Lucas-Kanade method. Combining these techniques, we make use of their advantages and reducing the effect of their drawbacks. In addition, the software component, utilize a mechanism of person reidentification.

Finally, the software component, has been validated experimentally and tested on a well-known database, Multiple object tracking dataset.

# Resumen

Deep learning ha surgido por sus grandes mejoras respecto a las técnicas reinantes en los problemas más complicados en Inteligencia Artificial, inversiones masivas de gigantes industriales y por un crecimiento exponencial en el número de publicaciones científicas. Deep learning es una herramienta más dentro del conjunto de herramientas de Machine Learning, cuyo propósito es hacer aprender a las máquinas.

En ciertas áreas de la visión artificial han sido muy exitosas, sin embargo, en el campo del seguimiento visual aún están por desarrollar, por eso hemos abordado el problema del seguimiento visual de múltiples peatones con técnicas de deep learning. Así, en este trabajo se ha diseñado y construido un componente software que usa el paradigma de *tracking by detection*. Empleando técnicas de deep learning, con *tracking by matching*, usando el algoritmo de Lucas-Kanade. Combinando estas dos técnicas, recogemos sus ventajas, minimizando el efecto de sus inconvenientes. Además, el componente, también incorpora un mecanismo de reidentificación de peatones que mejora el seguimiento. Finalmente, el componente desarrollado se ha validado experimentalmente y se ha probado en la conocida base de datos de seguimiento visual *Multiple object tracking*.

# Contents

<b>List of Figures</b>	<b>VI</b>
<b>List of tables</b>	<b>VIII</b>
<b>Acronyms</b>	<b>IX</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Computer vision . . . . .	2
1.2 Object tracking . . . . .	4
1.3 Deep learning in computer vision . . . . .	7
<b>2 Objectives</b>	<b>11</b>
2.1 Objective . . . . .	11
2.1.1 Requirements . . . . .	12
2.2 Methodology . . . . .	12
<b>3 Theoretical background</b>	<b>13</b>
3.1 Tracking . . . . .	13
3.1.1 Detection in tracking . . . . .	14
3.1.2 Feature tracking . . . . .	18
3.1.2.1 Features . . . . .	18
3.1.2.2 Motion estimation . . . . .	20
3.2 Person re-Identification . . . . .	24
3.2.1 Siamese networks . . . . .	25
<b>4 Software implementation</b>	<b>27</b>
4.1 System overview . . . . .	27
4.2 System design . . . . .	30
4.2.1 Object detector by neural networks . . . . .	30
4.2.2 Feature tracking . . . . .	30
4.2.3 Data association . . . . .	35

<b>5 Datasets and evaluation procedures</b>	<b>37</b>
5.1 Datasets for object detection . . . . .	37
5.1.1 Pascal Visual Objects Classes . . . . .	37
5.1.2 VOC07 . . . . .	38
5.1.3 VOC12 . . . . .	38
5.1.4 ImageNet . . . . .	38
5.1.5 COCO . . . . .	39
5.1.6 Dataset comparison . . . . .	39
5.1.7 Evaluation of object detection algorithms . . . . .	41
5.2 Datasets for multiple object tracking . . . . .	43
5.2.1 PETS . . . . .	43
5.2.2 MOT challenge . . . . .	44
5.2.3 Evaluation of multiple people tracking algorithms . . . . .	44
5.3 Datasets for pedestrian identification . . . . .	46
5.3.1 Evaluation for pedestrian identification . . . . .	47
<b>6 Experiments</b>	<b>48</b>
6.1 Validation experiments . . . . .	48
6.1.1 Object detector . . . . .	48
6.1.2 Tracking module . . . . .	49
6.1.2.1 Feature extraction . . . . .	50
6.1.2.2 Matching module . . . . .	52
6.1.3 Siamese networks . . . . .	52
6.2 Algorithm evaluation . . . . .	58
6.2.1 Performance . . . . .	58
6.2.2 Timing . . . . .	59
<b>7 Conclusions</b>	<b>62</b>
7.1 Future work . . . . .	63
<b>Bibliografía</b>	<b>64</b>
<b>8 Annex</b>	<b>73</b>

# List of Figures

1.1	Frontal view of autonomous car.	2
1.2	Tractography map.	3
1.3	Camera for inspection.	3
1.4	Augmented reality image.	4
1.5	Control room.	5
1.6	Visual tracking for science.	6
1.7	Visual tracking for sports analysis.	6
1.8	Visual tracking for art.	7
1.9	Representation of a Convolutional neural network.	8
1.10	Classification error in ImageNet challenge.	10
3.1	Mean average precision over the years in PASCAL dataset.	14
3.2	SSD detector scheme.	16
3.3	SSD architecture.	16
3.4	Comparison architectures.	17
3.5	Types of patches.	18
3.6	Optical flow example.	20
3.7	Optical flow with pyramids.	23
3.8	Siamese CNN topologies.	25
4.1	Block diagram of the component.	28
4.2	Timing of the component.	28
4.3	Flow chart of the system.	29
4.4	Detections of the algorithm.	30
4.5	Shi-Tomasi points on a person.	31
4.6	Detections with feature points.	32
4.7	Matched feature points.	33
4.8	Displacement of each ROI.	33
4.9	Tracking failure.	34
4.10	Tracking failure.	34

4.11	Wrong trajectory.	35
4.12	Spatio-temporal data association.	36
4.13	Spatio-temporal data association.	36
5.1	Distribution of VOC07 dataset.	38
5.2	Distribution of pascal.	40
5.3	Distribution of pascal.	41
5.4	Comparision interpolated and normal curve.	42
5.5	Example of Pets.	44
5.6	Example of measures.	45
6.1	Barplot of the timming.	48
6.2	Barplot of the timming.	49
6.3	Sequence of translational movement.	50
6.4	Sequence of no translational movement.	50
6.5	Barplot of the timming.	51
6.6	Different preprocessings.	51
6.7	Ilustration forward backward error.	52
6.8	Siamese CNN topologies.	53
6.9	Final layers.	54
6.10	Data augmentation.	56
6.11	Results training.	57
6.12	CMC plot.	57
6.13	Performance-timing comparision.	58
6.14	Barplot of the timming.	60
6.15	Zoom in of the barplot.	60
6.16	Time versus points and size of the ROI.	61
7.1	Differences texture examples.	63

# List of tables

3.1	Summarize of the object detectors.	17
5.1	Datasets tables	41
5.2	Statistical comparision datasets.	46
6.1	Results algorithm global.	58
6.2	Results algorithm by sequences.	59

# Acronyms

**CNN** Convolutional neural networks.

**COCO** Common Objects in Context.

**ReLU** Rectified Linear Unit.

**VOC** Visual Objects Classes.

# Chapter 1

## Introduction

As engineers we want to build systems that are better than our brain, to point out the difficulty of this endeavor, we can summarize the brain characteristics as follows: it has 100 billion computing elements, processing and memory are performed by the same components, works as parallel recurrent paradigm, it solves problems not soluble by previous machines, and it only requires 20 watts of power.

There are several computation challenges very interesting but, despite recent success in most of them, it struggles us to reach brain performance and efficiency. Machines have beaten us in extracting information for large collection of data, they can process larger amount of data than the humans. Also, in memory, they beat us, they can store more information and access faster than humans. It is not a matter of speed computation, they also exceed us in reasoning tasks, like playing chess or Go. However, in low-level sensorimotor skills, like seeing or walking our brains performs better than machines. These kinds of tasks, that humans perform unconsciously, for a computer is really complex to achieve it.

This fact is called the Moravec's paradox, this paradox came out during the dawn of Artificial Intelligence back in the 80s when M.Minsky, R.Brooks, and H. Moravec tried to mimic human skills by reverse engineer the brain. This paradox consists in, contrary to traditional assumptions, high-level reasoning requires very little computation, but tasks involving with perception, attention, visualization, motor, and social skills require enormous computational resources and are difficult to transfer to machines.

One possible explanation of this paradox, is based on evolution, humans skills are implemented biologically, improvement over years of natural selection. The older a skill is, the more time natural selection has had to improve the design. In contrast, abstract

thought developed only very recently is easy to implement due this shorter developed. These categories of intelligence will take much more time to implement in machines, but research keep going.

## 1.1 Computer vision

In the late 60s, computer vision began at universities that were pioneering artificial intelligence. It was meant to mimic the human visual system, as a stepping stone to endowing robots with intelligent behaviour. In 1966, it was believed that this could be achieved through a summer project, by attaching a camera to a computer and having it *describe what it saw*.

This describes the excitement of that time and their underestimation of the field as we stated in the previous section. Although, it is a complex area of studying, it has been a lot of developments during fifty years, real world application have been developed and are part of daily use.

One recent application of computer vision, is the usage of these technique on robotics, in particular on autonomous cars. These cars developed by technological giants are being used in some States of US. In these systems, the surrounding information is extracted by cameras. As we can observe in figure 4.2 are used to detect other types of vehicles on the road.

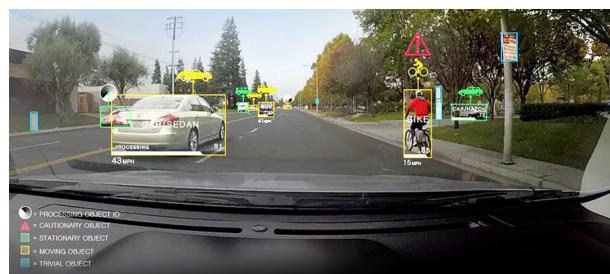


Figure 1.1: Frontal view of autonomous car.

Another computer vision's area is medical imaging, this area studies the techniques and process of creating visual representations of the interior of a body for clinical analysis. One example is the tractography map, like the one in figure 1.2 created from a diffusion weighted images, it allows us to establish connections between different areas of the brain.

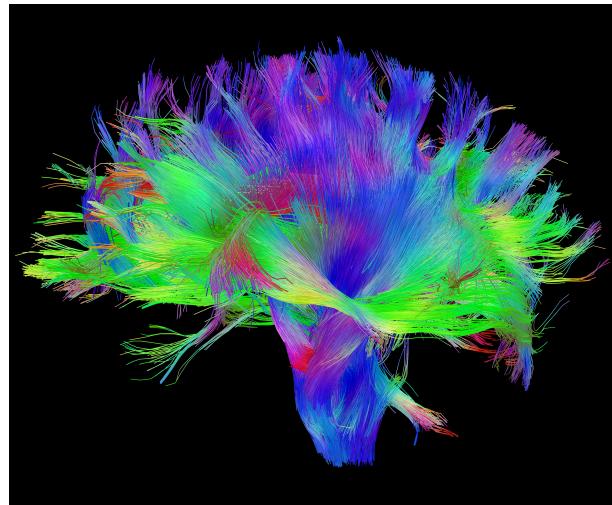


Figure 1.2: Tractography map.

One pioneer in technology is the automation industry, technology used to manufacture quickly and better. Computer vision is deeply used in factories, their usage is for inspection quality of manufactured products, they can check whether a product fulfill quality characteristics, as we can observe in figure 1.3.

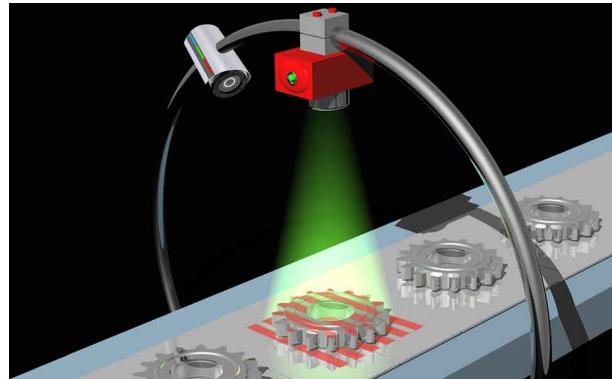


Figure 1.3: Camera for inspection.

Finally, another application of computer vision, is to mix the information provided by the camera with graphics, this is called, augmented reality. One example of it, is the work of the company Snapchat, it allows you to render different artifacts on an image and share it with your friends. We can observe one example on figure 1.4.



Figure 1.4: Augmented reality image.

## 1.2 Object tracking

In the widely computer vision field there are several study areas, one of them is Object tracking. It consists in estimating target state over time from image sequences. As a state we can embed the position, velocity, shape, appearance or any interesting characteristics. It is very challenging field due to:

- Variations due to geometric changes, some targets might be deformed as they move in the scene which would change their structure.
- Variations due to photometric factors, the appearance of the targets might change due to changes in illumination.
- Occlusions, targets might mix with other elements of the scene from the camera perspective.
- Image quality, the image sequences could incorporate noise or a low resolution.
- Similar objects in the scene, this could cause problems to maintain the identity of targets.

To solve these problems the community has used the following paradigms [1]:

- **Tracking using matching**, these kind of methods performs a matching of the representation between the current and the possible candidates in the next frame. Key points of these methods are the representation and the similarity measurement to perform the matching. The most famous methods are Normalized Cross-Correlation [2], Lucas-Kanade tracker [3], Kalman appearance tracker [4] and Mean shift tracking [5].
- **Tracking-by-detection**, these kinds of methods build a classifier to distinguish target pixels from the background. Once you have the detection, you need a data association method to link those detections. Traditionally the community has used kernel methods with support vector machines [6] to perform the detections, but in the recent years people are shifting to neural networks. In the data association algorithms graph theory techniques are dominant [7] [8].
- **Tracking learning and detection**, this is an extension of the previous category. It includes a mechanism to update the classifier during the execution of the system. This learning procedure allows the algorithm to be invariant to changes in the target. The most famous algorithms are the Predator [9] and the Alien [10].

These kinds of algorithms are quite mature and are deployed in real life applications. Like all the informatics process it allows us to process a huge quantity of information really quickly.

In video surveillance, it allows us to track all the targets without human intervention and notify when there are dangerous situations.



Figure 1.5: Control room.

For science, it allows us to study the environment, in the case of the figure 1.6, for humans

it will be difficult to do not miss the correct identity of these ants. With information supplied by the tracking, scientist can study how animals move and interact with others.

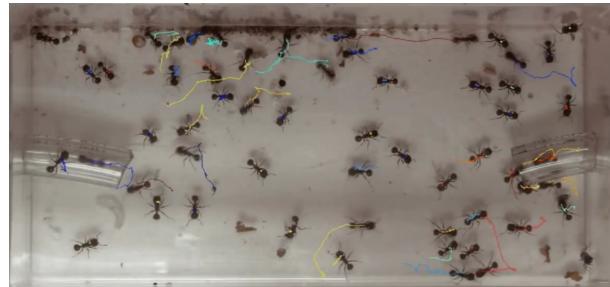


Figure 1.6: Visual tracking for science.

These algorithms are deeply used in all kind of sports like the NBA and NFL. In this situations the algorithm tracks the players during the game and allows to analysis his performance and the strategy of the team, like in figure 4.3 .



(a) Input image

(b) Layout

Figure 1.7: Visual tracking for sports analysis.

In artistic performance, these algorithms are used to track the subject and render some graphics in the scene, like an extension to video of augmented reality. In the example of the figure 1.8, the systems tracks the singer's head and projects a visualization of his voice.

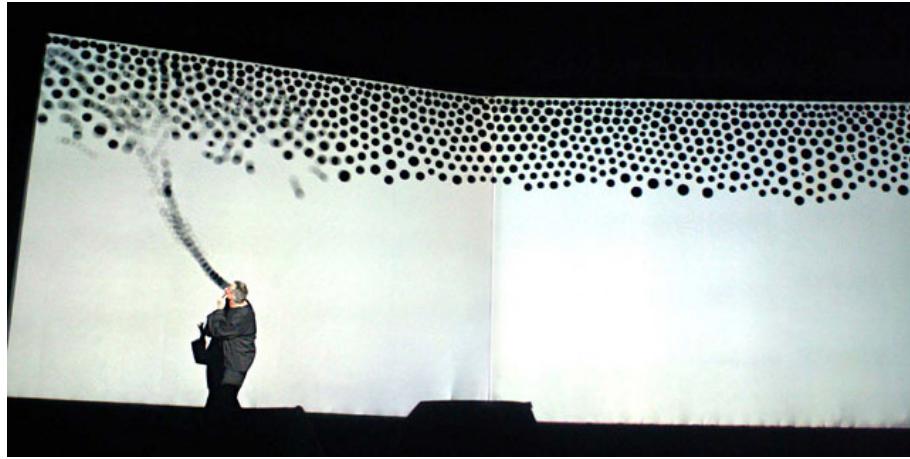


Figure 1.8: Visual tracking for art.

### 1.3 Deep learning in computer vision

Deep learning has raised by drastic improvements over reigning approaches towards the hardest problems in Artificial intelligence (AI), massive investments from industry giants, and exponential growth in research publications. Deep learning is a tool inside the machine learning toolbox, the goal is to make machines learn.

The first incursion was made by Frank Rosenblatt, the Percetron [?]. Rosenblatt conceived of the Percetron as a simplified mathematical model of how the neurons in our brains operate. This model of the neuron built on the work of McCulloch-Pitts [11], who showed that a neuron model could model the basis OR/AND/NOT functions. Which in the early days of Artificial intelligence was great, because the predominant thought at that time was that making computers able to perform formal logical reasoning would essentially solve AI. However, the McCulloch-Pitts model lacked a mechanism for learning, which was crucial for it to be usable for AI. This is where the Perceptron exceeded, Rosenblatt came up with a way to make such artificial neurons learn, inspired by the Hebb's Rule. This learning method was as follows: if the output of the perceptron was low, increase the weights, otherwise decrease the weights if the output is too high. Also, another researchers came with ADALINE [12] learning procedure, they used the signal before the activation function to compute the derivative, how much the error changes when each weight is changed can be used to drive the error down and find the optimal weights values. Similar as we train the networks nowadays.

Researchers were really excited about this idea of Connectionism: that networks of such

simple computational units could be vastly much more powerful and solve the hard problems of AI. But in 1969, Minsky and Papert published an analysis on the limitations of perceptrons [13]. The biggest criticism was that a perceptron could not learn the simple boolean function XOR because it is not linearly separable. However, they stated that it could be learnt with multiple layers perceptron but the learning procedure did not work for multiple layers. After this book, the interest on Neural networks decreased, and initialize a period called *AI winter*, AI shifted to logic programming and common sense reasoning.

This period lasted till 1986 when Rummelhart, Hinton, and Williams published the algorithm of backpropagation [14], which specifically addressed the problems discussed by Minsky in Perceptrons, a method to train multiple layer neural nets. With this discovery in 1989, LeCun showed a real world application, recognized handwritten digits [15]. The architecture of this model was a convolutional neural network. It was inspired by the Neurocognitron [16] of Fukushima, which took ideas from studies of the brain. In particular the studies of Hubel and Wiesel, they propose that the visual cortex is formed by a hierarchical model, primarily for simple cells who respond for simple structures and then complex cells respond to a more complicated feature. As we can observe in figure 4.2, in the lower layer, the network learns Gabor like features, and while going upwards, the networks learns more abstract concepts.

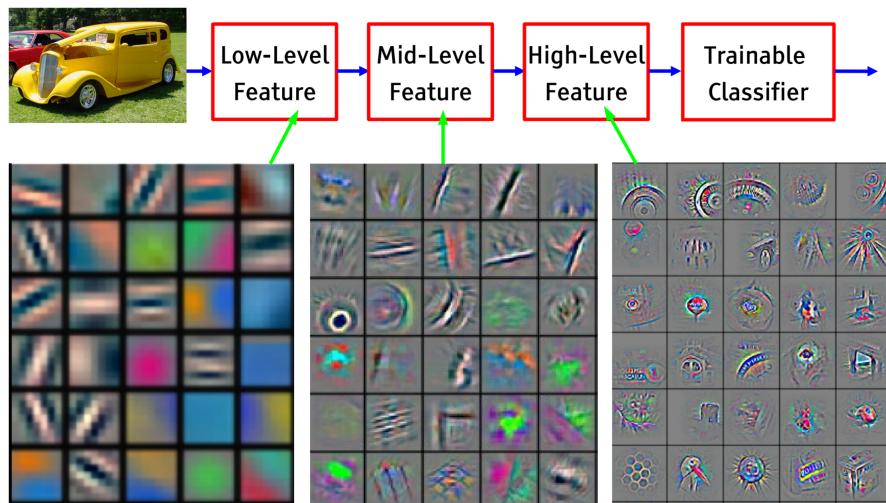


Figure 1.9: Representation of a Convolutional neural network.

But this approach did not scale to larger problems, the biggest source of problems was the vanishing gradients problem. When the backpropagation gradients backpropagates

trough the network, in some nodes the local gradient is very low ( in the extreme of the sigmoid functions ) the signal vanishes or saturates. By the 90s other techniques became the method of choice, like support vector machine (SVM). Although some progress was made for other kinds of problems.

- Unsupervised learning. This type of architectures is used to find a smaller representation of some data from which the original data can be reconstructed, it is useful for compression, visualization, and classification. One example of this architecture with neural networks is the Restricted Boltzmann machine [17], developed by Hinton.
- Reinforcement learning. The goal of this type of learning is to learn how to make good decisions, it requires rewards, not labels. One example of this sort of systems, is the TD-Gammon [18], a neural network that learned to be a backgammon player.
- Recurrent neural networks. Plain neural networks could not process sequences due to they do not have memory, they need mechanism to remember the pasts outputs. With memory, it can process sequences like audio or text. One approach to this, by Waibel [19] in 1989.

In 2006, there was a breakthrough [20], Hinton realized that a neural network with many layers really could be trained well, if the weights are initialized in a clever way. The basic idea was to train each layer one by one with unsupervised training ( like an autoencoder architecture ) and finally stack all together and train it in a supervised way.

Although, these improvements, the big step forward came in 2012, when AlexNet [21] beat the state of the art in the ImageNet challenge, an image classification challenge, where the error rate was 15.3% whereas the winner of the previous year was 26.3%. In the figure 1.10 we can observe the advance in the state of the art of the ImageNet challenge with the inclusion of deep learning techniques.

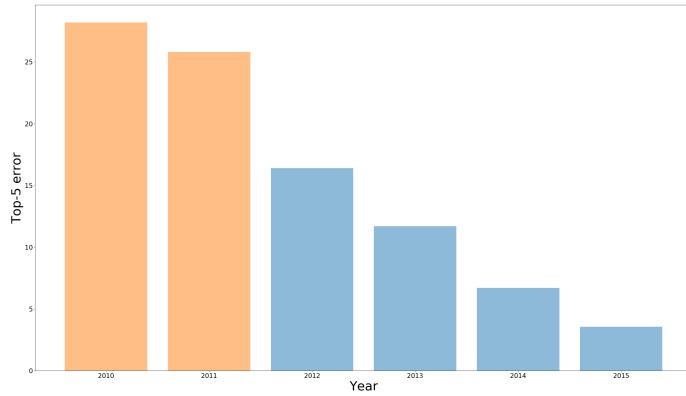


Figure 1.10: Classification error in ImageNet challenge.

The emergence of these techniques were the culmination of decades of research but the step forward was due by three aspects:

- **Appearance of large and high quality dataset**, The increasing size and quality of the dataset helps the networks to converge easily.
- **Parallel computation**, The increasing of computing capabilities helped train larger models in less time.
- **Optimization details**. With the discovery of the proper initialization and activation functions larger networks can be trained.

# Chapter 2

## Objectives

Once we put in context our work, in this chapter we explain the objectives of this thesis, its requirements and the methodology to accomplish it.

### 2.1 Objective

The main objective of this thesis is to develop and characterize an algorithm of multiple people tracking with two different methods: deep learning techniques and feature tracking. The essence of this work is study how combine it, and reach real time operation and a high performance. Finally, validate our solution with a international dataset, the Multiple Object Tracking dataset. We divided this target in several sub-objectives:

- **Object detector using deep learning.** Study the fundamentals of object detectors with deep learning techniques. We analyzed the performance on the main datasets and finally, we chose one to our task.
- **Development of a tracking module.** We studied which tracking technique would fit our problem, when it was selected, we implemented on code.
- **Join these two techniques.** We integrated this two techniques to perform a complete tracking algorithm.
- **Test the component on an international databases.** We validate our solution on well-known database.

### 2.1.1 Requirements

In addition to the previous objectives, our solution must satisfy the following requirements:

- The solution will make use of the JdeRobot framework, release 5.5, which is the developing environment of the *Grupo de Robótica* of the *Universidad Rey Juan Carlos*.
- The software will run on the GNU/Linux Ubuntu 16.04 environment.
- The algorithm will only make use of video sequences not other information.
- The algorithm must achieve an execution on real time and guarantee a precision.

## 2.2 Methodology

To achieve our objectives we used several tools that helped to monitoring the project for all members of the team. It allowed to comment or correct the task.

The main tool, it has been the videoconference, we established a weekly meeting with all the members of the team. In this meetings we showed the results and we shared our feedback with the other members of the team.

As complementary tools, we used a website and Github repository, these tools helped to control the development of our work. The website was developed using the wiki of JdeRobot [22], this shows the weekly tasks and results. The Git Hub repository [23] allows to access to the code by all the members of the team.

Our development plan was based on the Spiral model. It consists in four steps. In the first step, we determine the objectives of our project, in the second one, we analyze the risks and evaluate which problems we will face, then we develope and test our prototype and the last step consists to evaluate the results. We apply several iteration of this process till we get a satisfactory project.

# Chapter 3

## Theoretical background

In this chapter we explain the theoretical concepts of our work, these include the theory of tracking and person re-identification.

### 3.1 Tracking

As we explained in previous chapters, there are a traditional family of methods to solve the tracking problem. But with the incursion of deep learning techniques, they have been adapted to it and create new paradigms. The main ways to apply deep learning techniques to tracking are the following [24]:

- **Tracking-by-detection.** These methods use a specific class classifier and there is not need to train it online. So, these methods use a neural network to extract instances of the frames and then linked with temporal restrictions.
- **Tracking learning and detection.** Starting from the first frame of a video, a tracker will sample patches near the target object, and they are used to train a foreground-background classifier, and this classifier is used to score patches from the next frame to estimate the new location of the target object. These methods showed a state-of-the-art performance results. Unfortunately, neural networks are slow to train, therefore the speed of the method is reduced [25] [26].
- **Siamese based tracking.** In this approach, many candidate patches are passed through the network, and the patch with the highest matching score is selected as the tracking output [27].

- **Tracking as regression**, these methods are an extension of object localization using neural networks, these methods given an image containing an object predict the bounding box which contain the object in every frame [24]. They are restricted to one object.
- **Tracking with RNN**. From the output of an object detector, these tracking algorithms model the sequence of movement of objects using an recurrent neural network [28]. These methods represent the current state of art in tracking.

For this thesis we chose the tracking-by-detection paradigm. We get the detections with a object detector based on deep learning networks, and we link those detections with a tracking by matching, particularly, tracking by feature.

### 3.1.1 Detection in tracking

In object detection too, the emergence of the neural networks has supposed a turning point. As we can observe in 3.1, the mean average precision, has almost doubled since the appearance of deep neural networks.

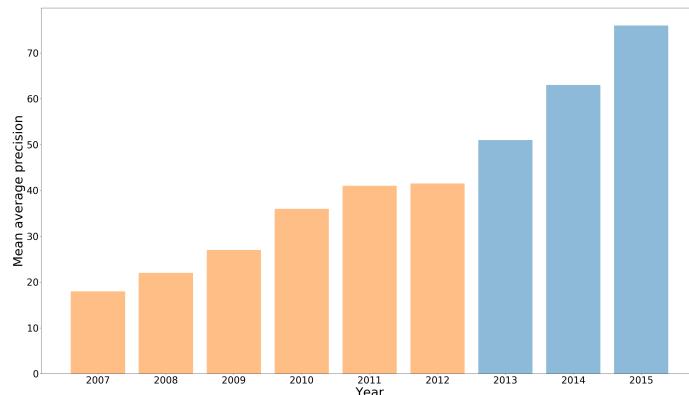


Figure 3.1: Mean average precision over the years in PASCAL dataset.

Present deep learning object detectors are based on three main family of architectures [29], named by the reference algorithm of the category: FasterRCNN, SSD, and RFCN, the characteristics of these systems are:

- Faster RCNN [30], it is the last output of a trilogy of detectors developed by R. Girshick and his team. Which are called Region-Based object detectors. They work

as follows: Use some mechanism to extract region of an image that are probable to be an object and then classify those proposals with a CNN. The first paper to do so, was [31], and suppose a breakthrough in the field, increasing the precision of the state of the art of those days. But, it had a messy pipeline, slow and difficult to train. Later on, they developed [32], in this paper they applied the region proposal algorithm in the cnn feature map, so, they avoid to compute the features for each proposal. They increase the speed and it could be trained much easily. Finally, they showed FasterRCNN [30], in this algorithm, they eluded the external region proposal algorithm and they implemented a CNN to compute those proposals. This CNN share parameters with the main net and they saved a lot of time. This network, has become the standard object detector with CNN. With the association of novel net architecture like ResNet [33], Inception [34], and [35] they have won all the contests.

- SSD, it stands for Single shot multibox detector. These family of method differs from previous ones considering that these treats the problem of object detection as a regression problem. So, they are called Regression-based object detector or single shot object detector due it does not have a region proposal algorithm, they classify the image with one mechanism. The maximum exponent of these algorithms are [36] and [37]. These work as follows, they discretize the image at the features level in a fixed grid and for each grid it predicts a class and some number of bounding boxes with different shapes and sizes. It merges all, and apply a Non-Maximum suppression algorithm to obtain a set of detections. We can observe this process in 3.2. In addition, they apply this process in a multiresolution scheme as we can observe in 3.3 to deal with objects of different sizes.

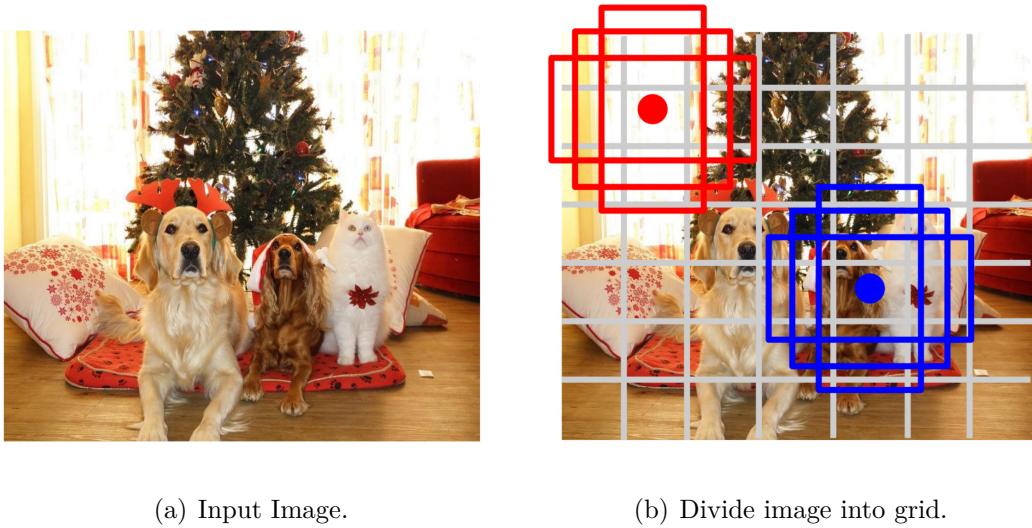


Figure 3.2: SSD detector scheme.

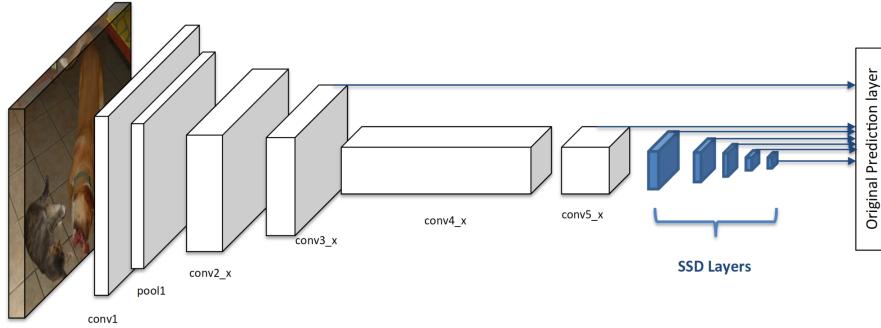


Figure 3.3: SSD architecture.

- RFCN [38], it stands for Region-based fully convolutional network and it was developed by the same authors of SSD. They noticed the lacks of the SSD, the SSD algorithm computes the object detector on the feature map, and at this level the features have a low spatial resolution, this involves do not detect small objects. So the authors inspired by the fully convolutional architectures, upsample those feature maps and compute the object detector like the SSD algorithm.

In the survey [29], they compared the different methods including changing the features extractors ( ResNet, Inception, VGG ) and they measured the precision ( mean average precision ) and computing time. This results are showed in 3.4. The conclusion are as follows, SSD is the fastest detector, RFCN it has the best balance between speed-accuracy, and FasterRCNN, is the most accurate detector although is slower than the other ones.

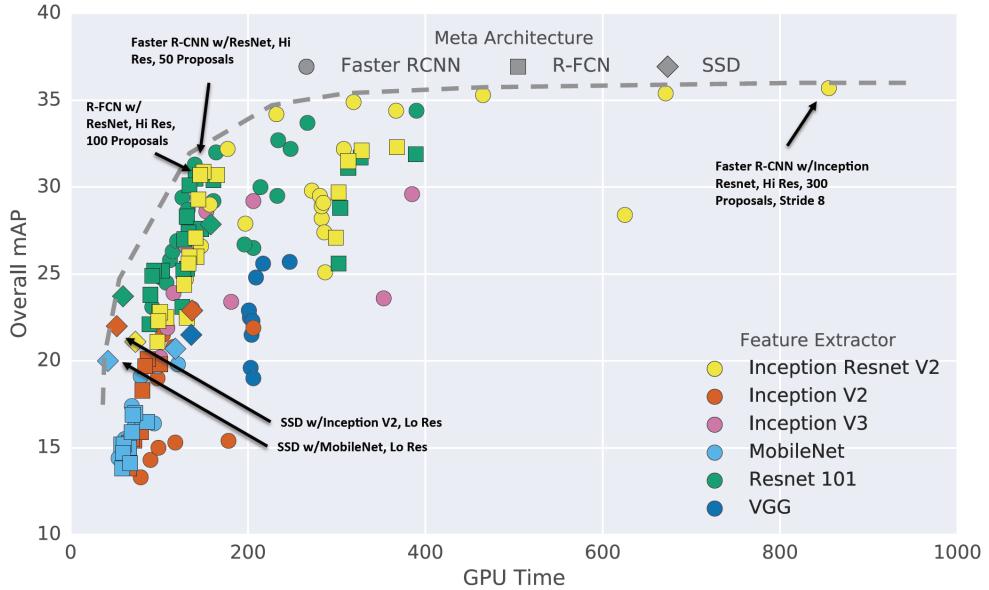


Figure 3.4: Comparison architectures.

We will finish off our review with a numeric comparison of the methods, as we can observe in the table 3.1. This information is extracted from the original papers with their implementation, all of them are trained with the union of the training set of VOC07, VOC12, and COCO, and subsequently evaluate on VOC07 test set on a Nvidia Titan X GPU. These results give us an intuition on which detector will be suitable for our task.

	mAP	mAP_person	FPS	Proposals
<i>RCNN</i>	66	64.2	0.077	2000
<i>FastRCNN</i>	70	69.9	6.7	2000
<i>FasterRCNN</i>	85.6	82.3	7	6000
<i>SSD300</i>	81.2	81.4	46	8732
<i>SSD512</i>	83.2	84.6	19	24564
<i>YOLO</i>	66.4	63.5	45	98
<i>YOLOv2</i>	78.6	81.3	40	-
<i>RFCN</i>	83.6	-	10	-
<i>PVANET</i>	84.9	-	31.3	300

Tabla 3.1: Summarize of the object detectors.

### 3.1.2 Feature tracking

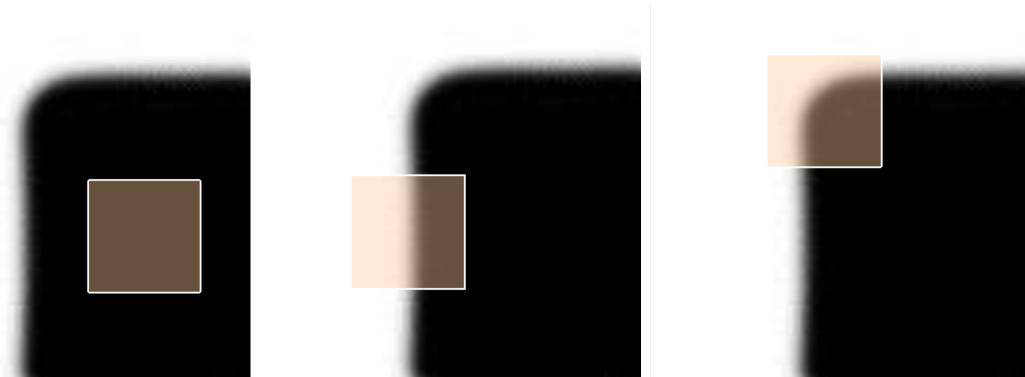
#### 3.1.2.1 Features

Our goal is to find points in an image, that can be found in other images and then compute some information, in this case, the movement. The characteristics of good features are:

- Repeatability, the same feature can be found in several images despite geometric and photometric transformations.
- Matchability, each feature has a distinctive description, thus easy to find.
- Efficiency, few features have to compact much more possible information.
- Locality, a feature occupies a relatively small area of the image, so therefore it is robust to clutter and occlusion.
- Performance, computation speed of features is a critical parameter.

Features points are used in all sort of operations in computer vision: Image alignment, 3D reconstruction, Motion Tracking, Object recognition, Index database retrieval, robot navigation and so on.

Looking at the figure 3.5, the flat patch, is a patch without texture and impossible to localize. Patches with large contrast edges are easier to localize, although straight lines segments at a single orientation suffer from the *aperture problem*, are also impossible to localize. Finally, patches with large gradients in at least two different orientations are the easiest to localize.



(a) Flat region.

(b) Edge region.

(c) Corner region.

Figure 3.5: Types of patches.

These intuitions can be formalized by looking at the simplest possible matching criterion for comparing two image patches, their weighted summed square difference:

$$E(u) = \sum_i w(x_i)[I(x_i + u) - I(x_i)]^2$$

where  $I(x)$  is the image,  $I(x + u)$  is the shifted image, and  $w(x, y)$  is a window function like a box or gaussian kernel around the pixel, and the summation  $i$  is over all the pixels in the patch. Then we are looking for points, which if we move according to  $u$  we have a change.

When performing feature detection, we do not know which other image locations the feature will end up being matched against. Therefore, we can only compute how stable this metric is with respect to small variations in positions  $\Delta u$  by comparing an image patch against itself:

$$E(\Delta u) = \sum_i w(x_i)[I(x_i + \Delta u) - I(x_i)]^2$$

Using a Taylor series expansion of the image function  $I(x_i + \Delta u) \approx I(x_i) + \nabla I(x_i) * \Delta u$  we can approximate the expression as follows:

$$E(\Delta u) \approx \sum_i w(x_i)[I(x_i) + \nabla I(x_i)\Delta u - I(x_i)]^2$$

$$E(\Delta u) = \sum_i w(x_i)[\nabla I(x_i)\Delta u]^2$$

With algebraic notation it transforms to:

$$E(\Delta u) = \Delta u^T M \Delta u$$

where  $\nabla I(x_i) = [I_x, I_y](x_i)$  is the image gradient and  $M$  is the second moment matrix:

$$M = \begin{pmatrix} I_x^2 & I_{xy}^2 \\ I_{xy}^2 & I_y^2 \end{pmatrix}.$$

Computing the eigenvalue decomposition of this matrix, shows the directions of the fastest change, thus a measure of the *cornerness*. There are several algorithms that use in different ways this eigenvalues:

- Harris [39], they propose a corner detection response function. So for each pixel, they compute a matrix  $M$  and with it, they compute the function  $R$ ,  $R = \det(M) - a \operatorname{trace}^2(M)$ . if  $R$  is large, that pixel is a corner, if  $R$  is negative with larger magnitude, it is a an edge, and if  $R$  is small it is a flat region. So the they a threshold to classify those pixels as a corner.
- Shi-Tomasi [40], they define the *cornerness* in another way. The image has a maximum value ( e.g. 255), so  $\lambda_1, \lambda_2$  also have an upper bound, then it is only necessary to check that  $\min(\lambda_1, \lambda_2)$  is large enough, this is how they define *cornerness*. This feature is called good features to track, because the authors defined a *good* features those whose motion can be estimated reliably, and they reached the same conclusions as Harris. This method is implement in the OpenCV's routine `goodFeaturesToTrack()`.

### 3.1.2.2 Motion estimation

Now, we have invariant points, we want to estimate the motion of those points. In order to do so, we compute the optical flow. This is the apparent two-dimensional motion of brightness pattern in the image. In the next figure 3.6 we visualized this idea.

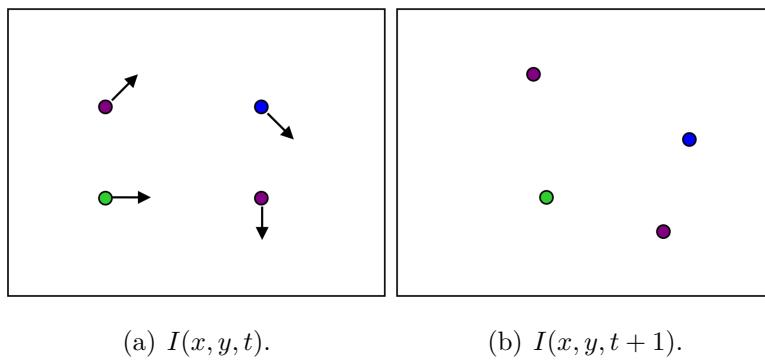


Figure 3.6: Optical flow example.

So, the question is: How do we estimate pixel motion from image  $I(x, y, t)$  to image  $I(x, y, t + 1)$ . We need to solve the pixel correspondence problem. Given a pixel in  $I(x, y, t)$ , look for nearby pixels of the same color in  $I(x, y, t + 1)$ . Solving this problem is what is referred as the optical flow problem. By nearby pixels and same colour we have two assumptions:

- Colour constancy: a point in  $I(x, y, t)$  looks the same in  $I(x', y', t+1)$ . For grayscale images, this is called *Brightness constancy constraint*. Stated in mathematical formulation:

$$I(x, y, t) = I(x + u, y + v, t + 1)$$

$$0 = I(x + u, y + v, t + 1) - I(x, y, t)$$

- Small motion: Subsequent points do not move very far, so we can estimate the motion by Taylor expansion:

$$I(x + u, y + v) \approx I(x, y) + \frac{\partial I}{\partial x}u + \frac{\partial I}{\partial y}v + \text{higher order terms}$$

Then, combining these two equations, we get:

$$0 \approx I(x, y, t + 1) + I_x u + I_y v - I(x, y, t)$$

where  $I_x = \frac{\partial I}{\partial x}$ , isolating the terms we obtain:

$$0 \approx [I(x, y, t + 1) - I(x, y, t)] + I_x u + I_y v$$

$$0 \approx I_t + I_x u + I_y v$$

In the limit of  $t$ ,  $u$  and  $v$  approaches zero ( assumption of small motion ), so it becomes, what it is called the the *brightness constancy constraint equation*:

$$0 = I_t + I_x u + I_y v$$

If we look closely, we realized that we have two unknowns  $u, v$  and one equation. This is an underdetermined system. Intuitively, this means, that locally we can only determine the component of the flow in the gradient direction, the component of the flow parallel to an edge is unknown, this is the called the aperture problem. To recover the motion we need to add some extra constraints. There are several types of constraints to solve this problem:

- **Global constraint**, adding a smooth constraint to the brightness constraint, this new constraints penalizes for changes in  $u$  and  $v$  over the images, it assumes that the motion fields vary smoothly over the image. This approach was developed by Horn and Schunk [?].
- **Local constraint**, locally the motion field is almost the same, so we add the neighbours pixels to the equation. This approach was developed by Lucas and Kanade [41].

### Local constraint

In this thesis we use the Local constraint to solve the optical flow problem. As we stated above, we add a local constraint to get more equations, this assumes that the motion field is the same in the locality. From the brightness constraint equation:

$$0 = I_t(p_i) + \nabla I(p_i) [u \ v]$$

Adding the neighborhood equations:

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ \vdots & \vdots \\ I_x(p_n) & I_y(p_n) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} -I_t(p_1) \\ \vdots \\ -I_t(p_n) \end{bmatrix}$$

Now, there are more equation than unknowns, it is an overdetermined system, we have to solve it with the least squares technique. It is based on the optimization of the function:

$$(A^T A) d = A^T b$$

Using the image notation:

$$\begin{bmatrix} \sum I_x I_x & \sum I_y I_x \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} -\sum I_x I_t \\ -\sum I_y I_t \end{bmatrix}$$

The system has a solution when  $A^T A$  is invertible, it will be invertible when is well conditioned, this is when the ratio of the great and the small eigenvalues of the matrix is large but no too much. The matrix  $A^T A$  in terms of image formulation is the second order matrix that we stated in the section 3.1.2.1 developing the *cornerness*, then in order to be solvable it should have a strong gradient in both directions. After checking the invertability, we can solve the problem and extract the motion field:

$$d = (A^T A)^{-1} A^T b$$

Thus, using the image notation:

$$\begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} \sum I_x^2 & \sum I_y I_x \\ \sum I_x I_y & \sum I_y^2 \end{bmatrix}^{-1} \begin{bmatrix} -\sum I_x I_t \\ -\sum I_y I_t \end{bmatrix}$$

In practice motion is large, the assumption that it is small fails, consequently the approach using Taylor expansions. For two reasons, the linearity does not hold, in order to solve it, we apply an iterative refinement, which consists in compute the displacement, apply it to the pixels, and compute it again till it converges. The other one is there are local minimum and it will fail into it. To solve it, we need to utilize a coarse to fine approach, the idea is to use multiresolution to compute optical flow, the basic is that in a low resolution image the motion between pixels is very small and we can compute optical flow.

So, in order to do so, we use image pyramids, this consists in downsample these images to specific resolution, then in top level, we compute the motion field using the previous stated method, then we upsample the motion field and the images, We apply a transformation to one image according to the motion field computed in the previous level and then compute the optical flow between that transformed image and the other image, we apply this algorithm in all the resolutions

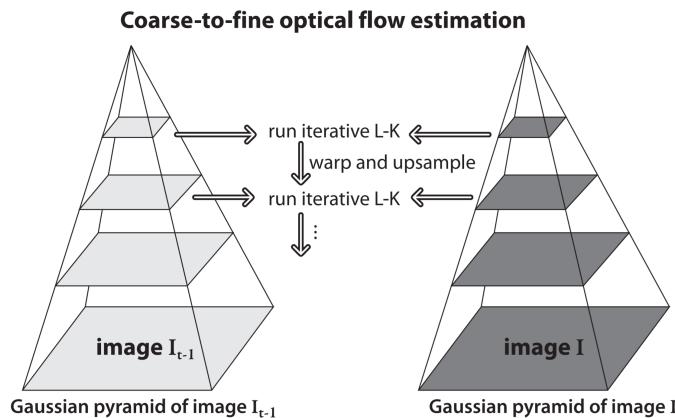


Figure 3.7: Optical flow with pyramids.

## 3.2 Person re-Identification

One of the problems of tracking is the reidentification of pedestrians. Add a module of reidentification helps to maintain the identity of the pedestrians.

Person identification is thoroughly studied in the field of Biometrics, it consists of knowing one biometric characteristic, and comparing with a claiming identity query. Specifically, the topic of pedestrian identification based on images has raised in the last years, this is due to the growth of surveillance applications. Also inside the topic of tracking there is a subfield called data association which studies this problem, it consists in matching tracks with further pedestrian detections.

Let's define mathematical the problem, consider  $G$  is a gallery composed of  $N$  images, denoted as  $(g_i)_{i=1}^N$ . They belong to  $N$  different identities  $1, 2, \dots, N$ . Given a query image  $q$ , its identity is determined by:

$$i^* = \arg \max_{i \in 1, 2, \dots, N} sim(q, g_i)$$

where  $i^*$  is the identity of probe  $q$ , and  $sim(., .)$  is some kind of similarity function. There are several categories of this similarity function [42]:

- **Hand-crafted systems.** This involves two components, an image descriptor and a distance metric algorithm. The most common image descriptors are those used in computer vision too, like colour [43], texture [44], SIFT [45], bag of word [46]. The general idea of distance metric learning is to keep all the vectors of the same class closer while pushing vectors of different classes further apart. The most commonly used formulation is based on the class of Mahalanobis distance function [47], [48]. Other works focus on learning discriminative subspaces [49].
- **Deep learning techniques.** Two types of CNN models are commonly employed in the community, the first is the classification model as used in image classification, the output is an identity label, and the second one is the siamese model using image pairs as input. The major drawback of the classification models is that they need a great quantity of training data by category, and most of the identifications datasets only provide a few examples for identity. So currently methods focus on siamese models.

The main differences between them, is that in hand-crafted methods, feature representation of the data and the metric are not learned jointly, instead, deep learning techniques jointly optimize the representation of the input data conditioned on the *similarity* measure being used [50].

### 3.2.1 Siamese networks

The first work with siamese architectures were developed by LeCun [51], [52] and they addressed the identification of signatures, besides the siamese networks are used in a variety of problems like: image recovery [53], feature descriptor [54], comparing patches [55], one shot learning [56], and learning visual similarity [57].

Siamese CNN topologies can be grouped under three main categories, depending on the point where the information from each input is combined:

- **Cost function.** Input patches are processed by two parallel branches featuring the same network structure and weights. Finally, the top layers of each branch are fed to a cost function.
- **In-network.** The top layers of the parallel branches processing the two different inputs are concatenated and some more layers added on top of that.
- **Joint data input.** The two input patches are stacked together forming a unified input to the CNN.

Graphically we can observe those differences in 3.8.

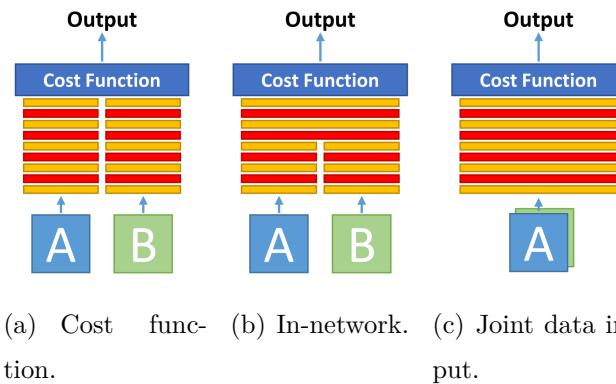


Figure 3.8: Siamese CNN topologies.

While the two first approaches have yield good results and historically were dominant, the best performance is obtained with the joint data input strategy. As pointed out by [55] and further corroborated by [58], and [59], jointly using information form both images from the first layer tends to deliver a better performance.

In the field of person re-identification, the community has used these architectures, and they also, have developed their own loss function, what is called *contrastive loss*, this loss is an extension of the Hinge loss of the SVM. This loss longs for getting close similar pairs and moving away according to one defined margin, dissimilar pairs. Although, the binary cross entropy is used by the community. Also, the community has focused in the developing of the datasets, increase the size and quality but there are not any landmark dataset.

There are several papers in the literature, one of the most famous is developed by Ahmed [?], they used In-network architecture although in order to join to the convolutional layers, they used *cross-input neighborhood differences* layer, this layer tried to increase the differences between the features of the inputs and obtain richer representation to the classification layer.

Another paper was published by Leal-Taixé [?], they are also the authors of the MOT challenge, their network used a cost function architecture besides they used as inputs the two images and their optical flow. They used the network as part of a data association algorithm.

# Chapter 4

## Software implementation

In this chapter, we explain the algorithm that we have designed for solving the visual people tracking problem and its software implementation.

### 4.1 System overview

The main contribution of this work is to develop a tracking algorithm that utilizes a neural network and do not miss the real time objective. To do so we use the tracking-by-detection framework, it combines people detection using a neural network, somehow, slow but very accurate, and link these detections with feature tracking, very quick but prone to drift.

The architecture of the system is summarized in the diagram 4.1, its inputs is a sequence of frames and its output is a CSV file. This file has the structure of the MOT's evaluation software requires. We divided the computing in two threads, the object detector thread and the main thread. The first is responsible of given a frame compute their pedestrian detections and send them to the tracking thread, and the second one, the tracking thread, given a detection it computes the tracking procedure, in addition, when a new detection appears it combines with the tracker, this is called data association, and finally it saves the results of the algorithm. We called data association to the process of combine a detection with the tracks, this is different to the data association concep in the tracking-by-detection nomenclature. In our algorithm the data association, like in the tracking-by-detection nomenclature, is the tracking module which links the detections.

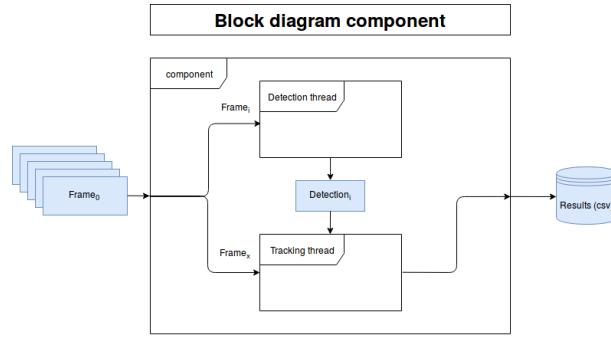


Figure 4.1: Block diagram of the component.

Temporally the system works as follows, when the algorithm starts, first of all it launches the object detector thread, it begins to compute the detections of the first frame. Meanwhile, the tracking thread remains idle, waiting for the detections. When the object detector finishes, it sends the detection to the tracking thread and at the same time starts to compute the detections of a predefined new frame. When the tracking thread has got a detection, it begins to compute the tracking between frames. Considering that the threads are not synchronized, the object detection thread goes ahead of the tracking thread it will have a detection which the tracking thread could not incorporate, thus, it saves the detections. When the tracking thread arrives to a frame that has a detection it has a mechanism to mix detection with tracks, what is called data association module, after running this module it will continue computing the tracking, when necessary it has a person re-identification module, to solve some possible identity incongruities. We can observe this temporal process in the next figure ??, when  $T$  represents a temporal step. When the object detector thread finishes computing all the detections it will *die* and the tracking thread computes their task previously stated till it does not have any frame to process.

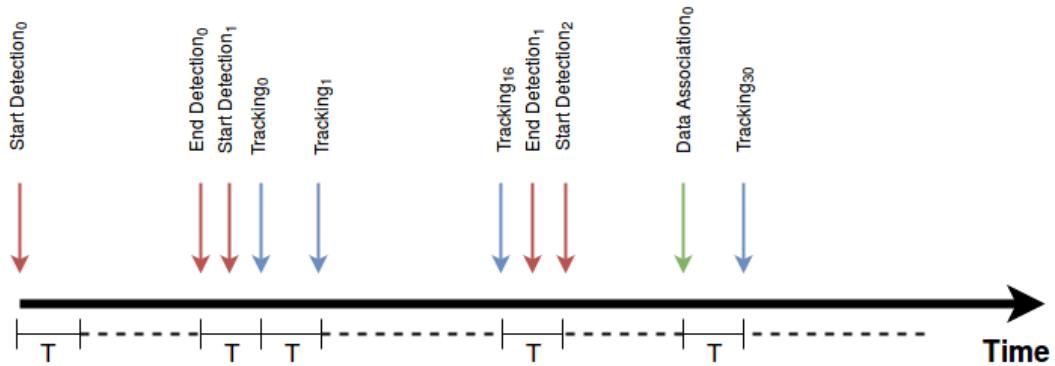
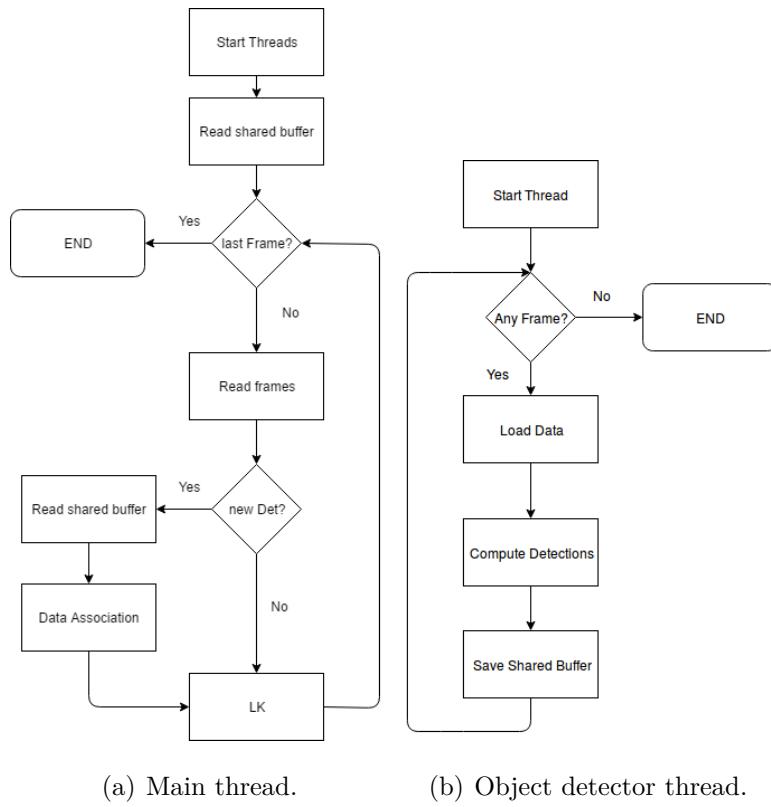


Figure 4.2: Timing of the component.

In the figure 4.3 there is a flow chart of the algorithm and it summarizes as follows: The object detector thread reads images, processes the forward pass of the neural network and saves the detections in the shared buffer, it repeats this sequence until it has processed all the predefined list of frames. In another hand, the main tread, activates the object detector thread and waits till it gets the first detection, after this, it starts the tracking algorithm. It reads the images and computes the motion of all the regions of interest. However, at beginning of each cycle it checks whether it has newer detection to mix in. This is the main operating mode of the algorithm and it summarizes in the figure ???. In the next sections we will develope in detail these parts.



(a) Main thread.

(b) Object detector thread.

Figure 4.3: Flow chart of the system.

We represent each person with a bounding box, in this bounding box we extract some features and compute how they movement through the frames, based on the movement of those features we will infer the movement of the bounding box, therefore, the movement of the person.

With this approach we accomplish to join two opposite technologies, we can exploit their benefits and reduce their drawbacks. We can get an accurate detection every 30 seconds and in between, we link those detections with the tracking by features module. In this

way, we reduce the fragility of the tracking, it tends to miss feature points and cause a drift of the estimation, with a new detection we can correct that drift. Also, we compensate the slowness of the detection with the speed of the feature tracking.

## 4.2 System design

In this section we explain the parts of the algorithm in detail.

### 4.2.1 Object detector by neural networks

As we stated previously we compute the pedestrian detector based on a CNN. This types of systems are very accurate but slow. We are constraint by its execution time, it takes 0.92 seconds for compute each detection, this allows us to get a new detection after 30 frames.

We selected the Single shot multibox detector as object detector, in section 6.1.1 we explain why chose this detector. Finally, in 4.4 we can observe the result of this step.



Figure 4.4: Detections of the algorithm.

### 4.2.2 Feature tracking

The essence of the tracking thread it is the tracking module, called LK in the figure 4.3. It stands for Lucas-Kanade algorithm, due it is the method that we used in this work. This module for each detection computes its displacement by computing the displacement of the features inside this bounding box. It computes for each ROI at each step.

For extracting the features, we use the OpenCV routine `goodFeaturesToTrack()`, this function determines strong corners on an image, according the Shi-Tomasi method. His parameters are the following:

- `image`, input image
- `maxCorners`, maximum number of corners to return. If there are more corners than are found, the strongest of them is returned.
- `qualityLevel`, parameter characterizing the minimal accepted quality of image corners.
- `minDistance`, minimum possible Euclidean distance between the returned corners.
- `mask`, optional region of interest.
- `blockSize`, Size of an average block for computing a derivative covariation matrix over each pixel neighborhood.
- `useHarrisDetector`, Parameter indicating whether to use a Harris detector.
- `k`, Free parameter of the Harris detector.

We applied an equalization to the image to obtain more high contrast points, we can observe this in the figure 4.5 and for each ROIs it looks like 4.6.



Figure 4.5: Shi-Tomasi points on a person.



Figure 4.6: Detections with feature points.

Once we have all the detections with their points, we can compute its displacement, to do so, we used the OpenCV's routine `calcOpticalFlowPyrLK()`, this function implements a sparse iterative version of the Lucas-Kanade optical flow with pyramids. And his parameters are the following:

- `prevImg`, first image.
- `nextImg`, second image.
- `prevPts`, vector of 2D points for which the flow needs to be found.
- `nextPts`, output vector of 2D points containing the calculated new positions of input features in the second image.
- `status`, output status vector, it tells you whether the flow has been found.
- `err`, each element of the vector is set to an error for the corresponding feature.
- `winSize`, size of the search window at each pyramid level.
- `maxLevel`, number of pyramid levels.
- `criteria`, parameter, specifying the termination criteria of the iterative search algorithm.

For the example's person we can observe the matching between the points of consecutive frames. Once we have the correspondances, we compute the displacement as the median of all of them in each dimension.



Figure 4.7: Matched feature points.

In the figure 4.8 we can observe a representation of this displacement for each ROI.



Figure 4.8: Displacement of each ROI.

### Problemas

Next we erased this points, we compute the motion as the median of all the contributions displacement vectors in each dimension. Also, the scale change is computed as follows: for each pair of points, a ratio between current point distance and previous point distance is computed, bounding box scale change is defined as the median over these ratios.

At this point, we have an implementation of the tracking algorithm given a set of bounding boxes. Nevertheless, we notice running our algorithm on the dataset, that the bounding box could compute the motion of the assigned pedestrian including the points belonging to another pedestrian who appears also in the bounding box. This interaction will cause a wrong estimation of the movement of the target and eventually the pedestrian will not be embedded by the bounding box. We can observe this event in figure 4.9.



Figure 4.9: Tracking failure.

So, we need a mechanism to detect these failures, therefore we studied how the motion algorithm behaves in these situations. When it has got a trajectory without crossing with other pedestrian, the vertical and horizontal displacement behave like a damping sine wave (if it goes away of the camera) or amplified sine wave ( if it goes closer to the camera ). But when it has got an interference with another pedestrian, it has an steep change in that wave. We can measure that change as the differences between the current displacement and the previous one normalized by current displacement, if this value overtake a thresolhd, we consider that tracket as a lost tracket. We can observe this process in the next figure 4.10, it belongs to previous trajectory 4.9

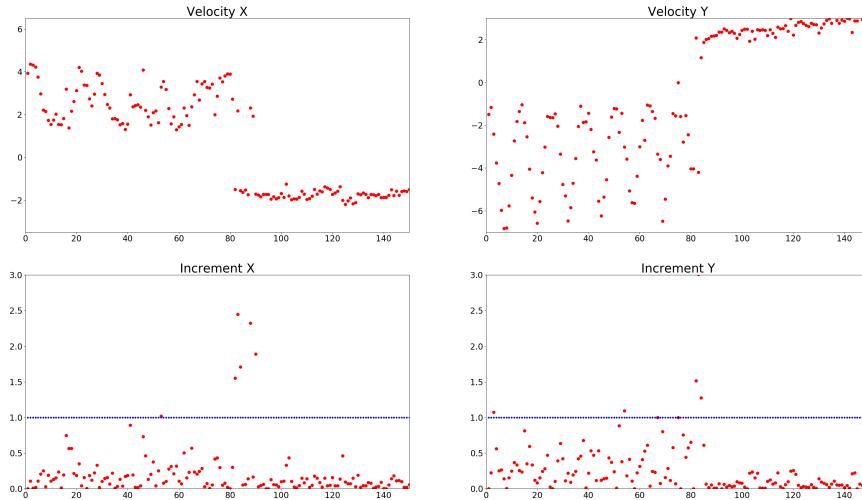


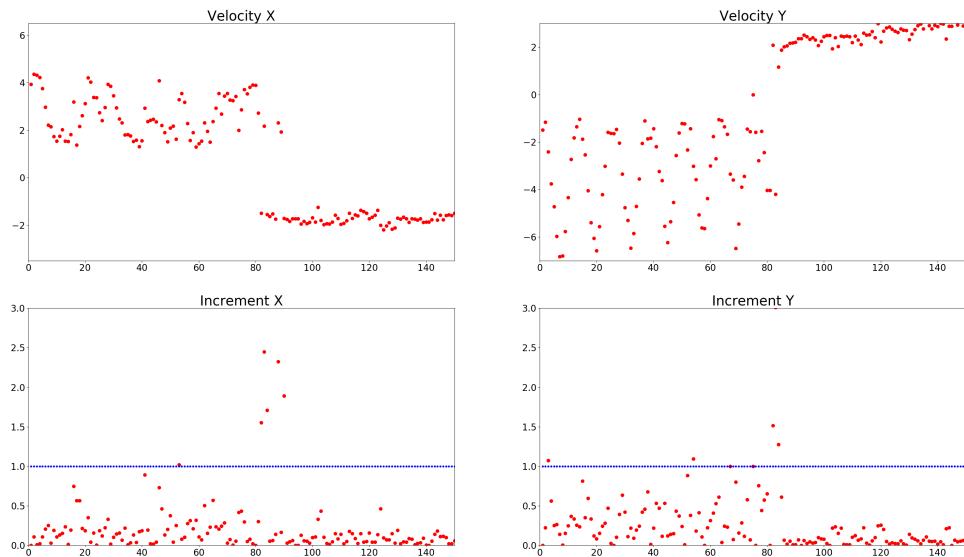
Figure 4.10: Tracking failure.

In contrast, when it does not cross with another pedestrian, the displacement does not get disrupt, then the normalized different with the previous displacement gets a low value,

we can observe this process in figure 4.11. We set a threshold to notice this interference and delete them from the current tracking execution, but we save the bounding box for following processings.



(a) Trajectory.



(b) Plots movement.

Figure 4.11: Wrong trajectory.

### 4.2.3 Data association

Once we computed the trajectories, in the next iteration we might have to add a detection, so we need a module to combine these trajectories with detections. Thus, for each pedestrian we distinguish three situations:

- **Situation 1**, the tracket has got a nearby detection, then the detection replace the

tracket bounding box. This is what is called spatio-temporal constraint.

- **Situation 2**, the tracket has not got a nearby detection, then the bounding box tracket continues.
- **Situation 3**, the pedestrian has not got a tracket but has got a detection. In this case we need to decided whether this pedestrian is new in the scene or it has been seen before ( it is a lost tracket ).

We can observe the procedure for situations number one and two in the figure 4.13, in green colour we can observe the detections and in blue colour the trackets. We defined nearby as the distance between the centres of the boundings boxes, this distance has to be lower than a threshold to be consider nearby.

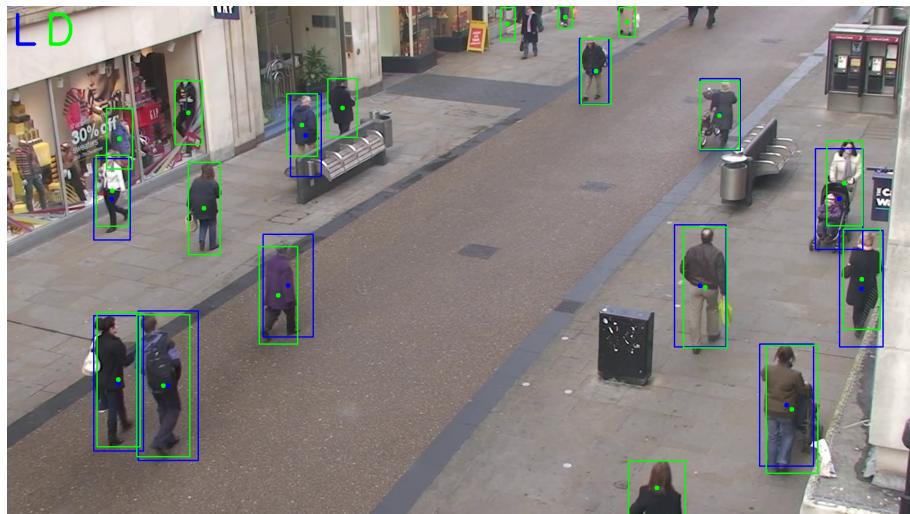


Figure 4.12: Spatio-temporal data association.

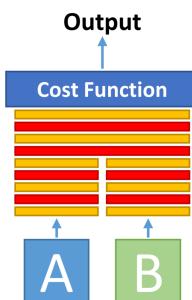


Figure 4.13: Spatio-temporal data association.

# Chapter 5

## Datasets and evaluation procedures

Next we explain the datasets and measures utilized in this thesis.

### 5.1 Datasets for object detection

This section describes the most common datasets used in object detection tasks. Throughout the history of computer vision research datasets have played a critical role. They not only provide a means to train and evaluate algorithms, they drive research in new and more challenging directions. In order to accomplish this, they provide:

- a collection of challenging images and high quality annotation.
- a standard evaluation methodology, so the performance of the algorithms can be compared.

In the next subsections, we will explain the main datasets for object detection task.

#### 5.1.1 Pascal Visual Objects Classes

The Pascal Visual Object Classes (VOC) challenge [62] is a benchmark in visual object category recognition and detection. Organised annually from 2005 to 2012, the challenge and its associated dataset has become accepted as one of the landmark benchmark for object detection. All the images are from the flickr consumer photographs website and annotated with the Mechanical turk tool.

The most popular editions of the challenge for object detection are those from years 2007 and 2012.

### 5.1.2 VOC07

The challenge of the year 2007 contains ten thousand images in the trainval and test sets, with almost twelve thousand objects. This was one the first datasets for object detection before the era deep learning. Also, it is very useful for researchers, due it has 2.5 mean object per image and it is very challenging. In the figure 5.1 we can observe the distribution of images and objects instances.

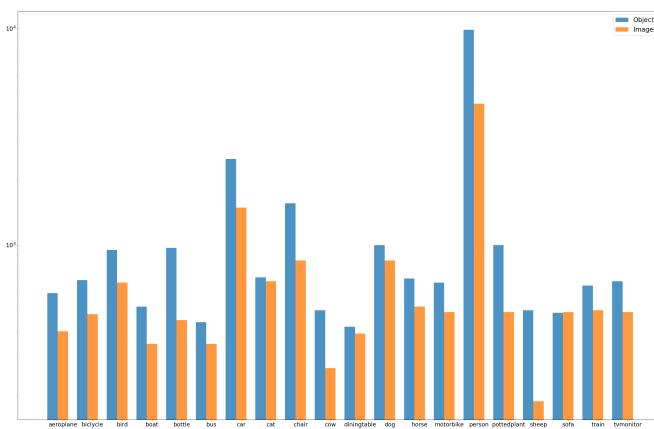


Figure 5.1: Distribution of VOC07 dataset.

### 5.1.3 VOC12

The 2012's edition is also, one of the most used dataset in object detection tasks. It increases the volume of images of the 2007 edition up to ten thousand images on trainval and test set and similar quantity of instances per image.

### 5.1.4 ImageNet

ImageNet project [63] with the challenge ImageNet Large Scale Visual Recognition Challenge [ILSVRC] was the first large-scale database, temporally developed to supply the deep learning techniques, eager of feed with tons of images. ImageNet aims to populate the majority of the 80000 synsets of WordNet with an average of 500-1000 clean and full resolution images. The collection was based on the query of that words on several image search engines and human refined on the Amazon Mechanical Turk platform.

In 2016, the project collects more than 10 million of annotated images with 1000 classes.

Although its main purpose is image classification, it has an object detection challenge with 200 categories with over a 1 million images with objects annotated.

### 5.1.5 COCO

The Microsoft Common Objects in Context also known as COCO dataset [64], is a dataset that address the three core research problems in scene understanding:

- detecting non-iconic views of objects, for many dataset most of the objects have an iconic representation, they appear unobstructed , near the center of the photo and with their canonical shape. So in this dataset, they included images to struggle the object recognition task, like objects in the background, partially occluded, amid clutter. Therefore, reflecting the composition of actual everyday scenes.
- contextual reasoning between objects, nowadays natural images contain multiple objects, and their identity can only be solved using context, due to small size or ambiguous appearance in the image, so in this dataset, images contain scenes rather isolated objects.
- the precise 2D localization of objects, also the detailed spatial understanding of object layout will be a core component of an image understanding system, so this dataset struggle to do so.

So, the three main tasks of this challenge are object classification, object detection and semantic scene labelling. This dataset contains 91 object categories, with 2.5 million labelled object instances in 328 thousand images, labeled with the Amazon Mechanical Turk tool.

### 5.1.6 Dataset comparison

The datasets from Pascal challenge are very useful to test object detection algorithms, their quantity is very handy ( a few thousands of images ) and contains a challenging quantity of objects per image, very interesting for the algorithms. But its little amount of images does not permit to train a network on this dataset, although it can be used to finetune the network.

The datasets for the ImageNet challenge are not used to much in object detection tasks, it contains a few instances per image, this did not encourage researcher to use it.

The COCO datasets, is the most recent one, is the one focus on object recognition, and the detection suppose a challenge due the objects are in common places and are very challenging to detect. And it is very interesting to due of the quantity of instances per image.

The COCO challenge contains 91 object categories with 82 of them having more than 5 thousand labeled instances. In total the dataset has 2.5 million labeled instances in 328 thousand images. In contrast to ImageNet dataset, COCO has fewer categories but more instances per category. Also, it has more instances per category than the VOC dataset. We can observe that difference in the chart 5.2. This fact aid in learning detailed object models capable to chance the variability and also its 2D location.

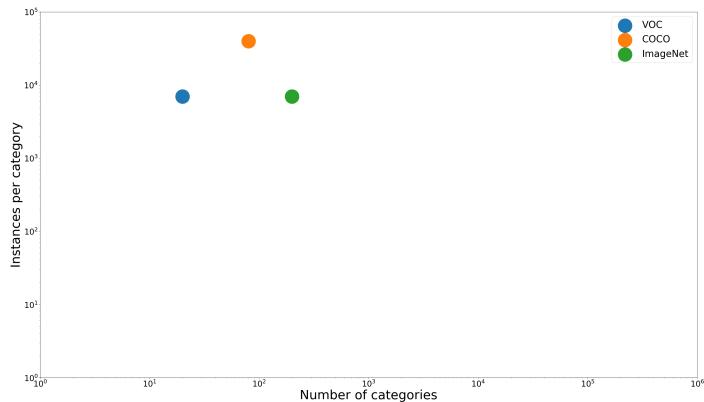


Figure 5.2: Distribution of pascal.

In addition, another prominent feature of the COCO over the other two, is the number of labelled instances per image which may aid in learning contextual information. This difference can be seen in ?? chart.

Moreover, the COCO dataset uses images from non-canonical point view, allowing to the algorithm to be robust to everyday views. This feature can be observed in the plot 5.3, in this plot we can observe differences views of the same category. And clearly the coco's images is the most uni-conic representation.



Figure 5.3: Distribution of pascal.

Finally, the table 5.1 summarizes the main statistics of the dataset stated previously.

	<b>VOC07</b>	<b>VOC12</b>	<b>ImageNet [ 2014 ]</b>	<b>Coco [ 2015 ]</b>
<i>trainval set</i>	5011	11540	476688	165482
<i>test set</i>	4952	10991	40152	81434
<i>Number of classes</i>	20	20	200	80
<i>Mean obj per image</i>	2.5	2.4	1.1	7.2
<i>Number person instances</i>	4690	8566	-	300000

Tabla 5.1: Datasets tables

### 5.1.7 Evaluation of object detection algorithms

In order to compare the performance of the different datasets, each challenges establish a clear measure. In this thesis, we used the interpolated average precision (AP), used in the Pascal VOC challenge (based on [65]).

For each class, the precision-recall curve is computed from a method's ranked output.

- Recall, is defined as the proportion of all positives examples ranked above a given rank.
- Precision is the proportion of all examples above the rank which are from the positive class.

The AP summarises the shape of the precision/recall curve, and is defined as the mean precision at a set of eleven equally spaced recall levels [0,0.1,...,1]:

$$AP = \frac{1}{11} \sum_{r \in (0, \dots, 1)} p_{interp}(r)$$

The precision at each recall level  $r$  is *interpolated* by taking the maximum precision measured for a method for which the corresponding recall exceeds  $r$ :

$$p_{interp}(r) = \max_{\hat{r}: \hat{r} > r} p(\hat{r})$$

The authors justified this measurement as a way to reduce the impact of the 'wiggles' in the precision/recall curve, caused by small variations in the ranking of examples. In the figure 5.4, we can observe this effect on the curve.

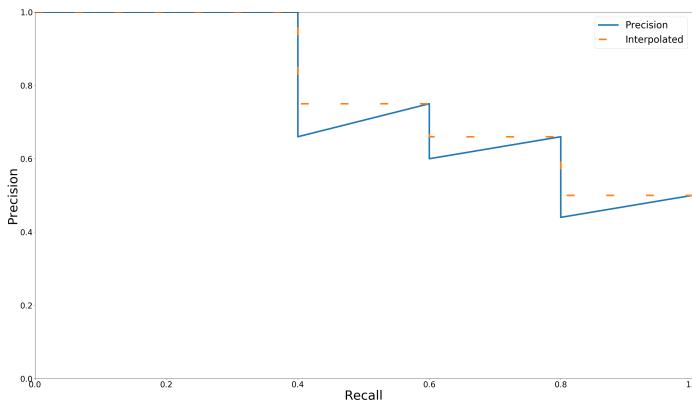


Figure 5.4: Comparison of interpolated and normal curve.

In addition, detections were assigned to ground truth objects and judged to be true/false positives by measuring bounding box overlap. To be considered a correct detection, the area of overlap  $a_0$  between the predicted bounding box  $B_p$  and ground truth bounding box  $B_{gt}$  must exceed 0.5 by the formula:

$$a_0 = \frac{\text{area}(B_p \cap B_{gt})}{\text{area}(B_p \cup B_{gt})}$$

where  $B_p \cap B_{gt}$  denotes the intersection of the predicted and ground truth bounding boxes and  $B_p \cup B_{gt}$  their union. The threshold of 50 % was set deliberately low to account for inaccuracies in bounding boxes in the ground truth data. Multiple detections of the same object in an image were considered false detections.

Finally, we want to point out, even we don't take into account in our implementation, setting the threshold IoU to a value of 0.5 could cause misdetections of small objects [63],

they propose an adaptive setting of that threshold based on the size of the ground truth and so detect correctly small objects. In practice this change only affects 5.5% of objects in the detection validation set.

## 5.2 Datasets for multiple object tracking

Evaluating and comparing multi-target tracking methods is not trivial for numerous reasons.

- First, the perfect solution one aims to achieve is difficult to define clearly. Partially visible, occluded, or cropped targets, reflections, and objects that are very closely resemble targets; all impose intrinsic ambiguities, such that even humans may not agree on one particular ideal solution.
- Second, a number of different evaluation metrics with free parameters and ambiguous definitions often lead to inconsistent quantitative results across the literature.
- Finally, the lack of pre-defined test and training data makes it difficult to compare different methods fairly.

In contrast to other research areas in computer vision, still lacks large-scale benchmarks.

### 5.2.1 PETS

Targeted primarily at surveillance applications [66]. The 2009 version consisted of 3 subsets, S1 targeted at person count and density estimation, S2 targeted at people tracking, and S3 targeted at flow analysis and event recognition. In the 5.6 we can observe one image from this dataset.



Figure 5.5: Example of Pets.

Even for this widely used benchmark, we observe that tracking results are commonly obtained in an inconsistent fashion: involving using different subsets of available data, different detection inputs, inconsistent model training that is often prone to over-fitting, and varying evaluation scripts. Results are thus not easily comparable [67].

### 5.2.2 MOT challenge

The aim of the Multiple object tracking [MOT] is to standardize the use of multiple people trackings datasets, in order to do so, they solve the problems in this kind of dataset, explained above.

### 5.2.3 Evaluation of multiple people tracking algorithms

A critical point with any dataset is how to measure the performance of the algorithms. A large number of metrics for quantitative evaluation of multiple target tracking have been proposed. Choosing unique general evaluation is still ongoing.

On the one hand, it is desirable to summarize the performance into one single number to enable a direct comparison. On the other hand, one might not want to lose information about the individual errors made by the algorithms and provide several performance estimates, which precludes a clear ranking.

We will explain two sets of measures that have established themselves in the literature the CLEAR metrics [68], and a set of track quality measures [69].

As in the object detection metrics, we can classify each tracklet, whether it is a true positive, that describes an actual (annotated) target, whether the output is a false alarm

( or false positive, FP ). This decisions is typically made by the well-known thresholding measure of Intersection of the union [IoU]. Also a target that is mised by any tracker is a false negative.

Due to we are working with multiple object, we assume that each ground truth trajectory has one unique start and one unique end point, that is not fragmented. So we need to penalty re-identification. This is called, identity switch [IDSW], and it is counted as if a ground truth target  $i$  is matched to track  $j$  and the last known assignment was  $k = j$ . The next figure summarizes the stated measures ( the grey area indicate the the matching threshold ).

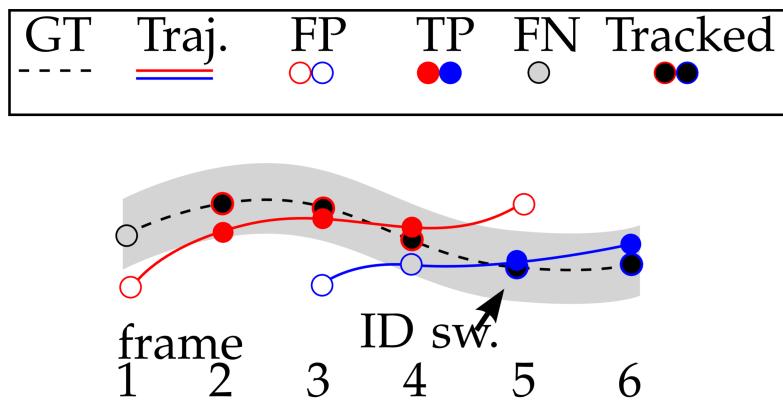


Figure 5.6: Example of measures.

Then after determining true matches and establishing the correspondances it is now possible to compute the metrics over all the sequences.

The multiple object tracking accuracy [MOTA] [68] is perhaps the most widely used figure to evaluate a tracker's perfomance. The main reason for this is its expressiveness as it combines as it combines three sources of errors defined above:

$$MOTA = 1 - \frac{\sum_t (FN_t + FP_t + IDSW_t)}{\sum_t GT_t}$$

where  $t$  is the frame index and GT is the number of ground truth objects. This measures gives an indication of the overall performance.

The multiple object tracking precision [MOTP] is the average dissimilarity between all true postives and their corresponding ground truth targets. For bounding box overlap, that is computed as

$$MOTP = \frac{\sum_{t,i} d_{t,i}}{\sum_t c_t}$$

where  $c_t$  denotes the number of matches in frame  $t$  and  $d_{t,i}$  is the bounding box overlap of target  $i$  with its assigned ground truth object. Thereby gives the average overlap between all correctly matched hypotheses. So, the MOTP is a measure of localization precision.

As we have stated above, another metrics are the track quality measures. Each ground truth trajectory can be classified as mostly tracked (MT), partially tracked (PT), and mostly lost (ML). This is done based on how much of the trajectory is recovered by the tracking algorithm. A target is mostly tracked if it is successfully tracked for at least 80% of its life span, without consider if there are an identity switch. If a track is only recovered for less than 20% of its total length, it is said to be mostly lost (ML). All other tracks are partially tracked. Finally antoher quality measure is track fragmentations (FM), it counts how many times a ground truth trajectory is resumed at a later point.

### 5.3 Datasets for pedestrian identification

A number of datasets for image-based re-Identification have been released, and some commonly used datasets are summarized in table 5.2.

Name	Date	Images	IDs	Cameras	Label	Evaluation
<i>VIPeR</i> [70]	2007	1264	632	2	hand	CMC
<i>iLIDS</i> [71]	2009	476	119	2	hand	CMC
<i>GRID</i> [72]	2009	1275	250	8	hand	CMC
<i>CAVIAR</i> [73]	2011	610	72	2	hand	CMC
<i>PRID2011</i> [74]	2011	1134	200	2	hand	CMC
<i>WARD</i> [75]	2012	4786	70	3	hand	CMC
<i>CUHK01</i> [76]	2012	3884	971	2	hand	CMC
<i>CUHK02</i> [77]	2013	7264	1816	10	hand	CMC
<i>CUHK03</i> [78]	2014	13164	1467	2	hand/DPM	CMC
<i>RAiD</i> [79]	2014	1264	43	4	hand	CMC
<i>PRiD 450S</i> [80]	2014	900	450	2	hand	CMC
<i>Market-1501</i> [81]	2015	32668	1501	6	hand/DPM	CMC/mAP

Tabla 5.2: Statistical comparision datasets.

Over recent years, progress can be observed, dataset's scale is increasing. Many of these datasets are relatively small in size, especially those of early days, but recent datasets, such as CUHK03 and Market-1501, are larger. Both have over 1000 ID's and over 10000 bounding boxes, and both datasets provide good amount of data for training deep learning models. In addition, the bounding boxes tend to be produced by pedestrian detectors, instead of being hand-drawn. Also, more cameras are used during collection, this helps to increase generalization. Although quantity of datasets, there are not a prominent dataset in the literature.

### 5.3.1 Evaluation for pedestrian identification

When evaluating identification algorithms, the cumulative matching characteristics (CMC) curve is usually used. CMC represents the probability that a query identity appears in different sized candidate lists.

Formally [82], for each probe  $p$  from  $P_G$  we sort the similarity scores against gallery  $G$ , and obtain the rank of the match. Identification performance is then stated as the fraction of probes whose gallery match it at rank  $r$  or lower. If the set of probes with a close match is:

$$C(r) = \left\{ p_j : \text{rank}(p_j) \leq r \right\} \quad \forall p_j \in P_G$$

where the rank is defined as before. We now define the Cumulative match characteristic (CMC) to be the identification rate as a function of  $r$ :

$$P_I(r) = \frac{|C(r)|}{|P_G|}$$

which we plot as the primary measure of identification performance. It gives an estimate of the rate at which probe images will be classified at rank  $r$  or better. One drawback of the characteristic is its dependence on gallery size,  $|G|$ .

# Chapter 6

## Experiments

In this chapter we explain how we chose the elements of our algorithm and the results of it on the datasets.

### 6.1 Validation experiments

#### 6.1.1 Object detector

For the choice of the detector we compare the detectors studied in the theoretical review on the *MOT16* dataset. In the figure ?? we can observe the ROC curve of different detectors.

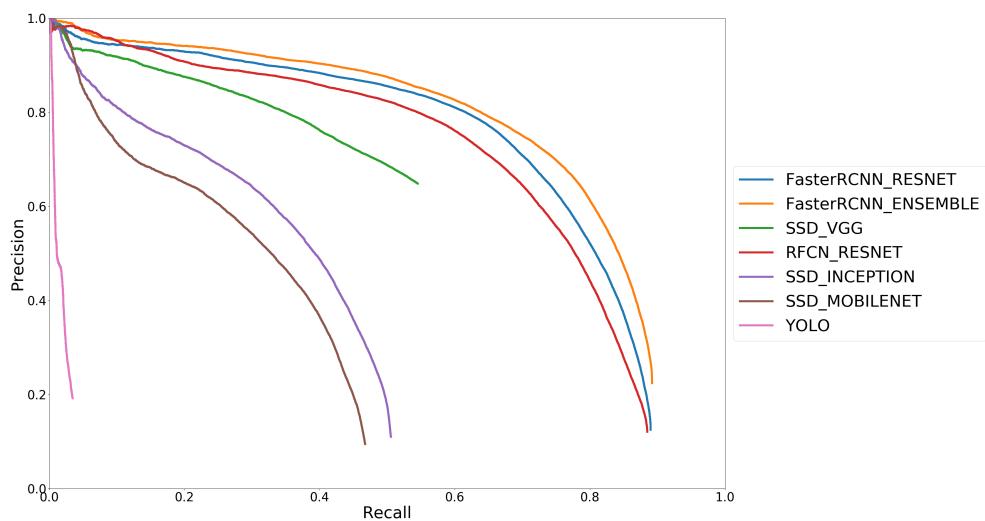


Figure 6.1: Barplot of the timming.

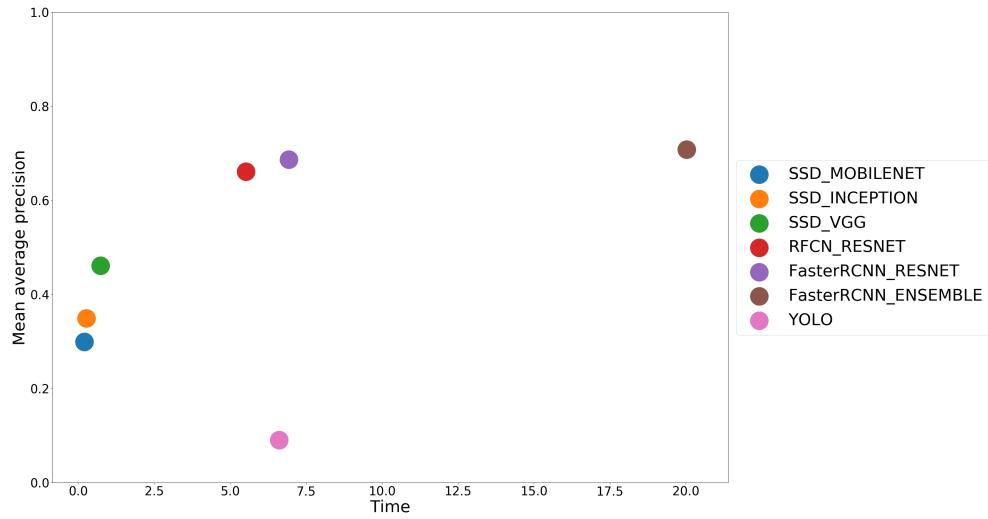


Figure 6.2: Barplot of the timming.

Finally, the conclusions of the characteristics are the following:

- **Faster-RCNN**, it scores 0.81 average precision on the dataset, although we could not used due a software incompatibilities, it requires a Nvidia GPU to be executed.
- **R-FCN**, the code is not publicly available.
- **DPM**, it scores 0.69 average precision on the dataset, although we could not used due a software incompatibilities, it requires a Nvidia GPU to be executed.
- **YOLO**, it scores 0.09 average precision on the dataset, it takes 6 seconds per image [83].
- **PVANET**, the code is not publicly available.
- **SSD**, it scores 0.73 average precision on the dataset, it takes 0.9 seconds per image. We used the tensorflow version of the detector [84] instead of the original code [85].

According to these results we chose the SSD detector, as object detector on this thesis.

### 6.1.2 Tracking module

The tracking module is inspired by the well-known tracking algorithm *MedianFlow* by Zalal et al[60] with his correspondent implementation in Python[61]. Next we explain the tracking module in details.

According to [1], our method belong to tracking using matching, considering that the algorithm performs a matching of the representation of the target model built from the previous frames. We do not build an additional extended model of the appearance and merge it with the matching, instead we will apply an object detector to correct the drift of the matching.

The tracking using matching is based on the optical flow, explained in ??, it computes the new position through gradient descent in several frames. We will assume that the motion is pure translational. As we can observe in the sequences of frames 6.3 of the dataset, the pedestrians move in translation way in the image plane, so this assumption is achieved.



Figure 6.3: Sequence of translational movement.

In contrast to the previous figure, we can observe the next figure 6.4 where the assumption of translation motion is not fulfilled ( this sequences does not belong to the used dataset, only showed to contrast the previous idea) and a translational assumption will failed.



Figure 6.4: Sequence of no translational movement.

#### 6.1.2.1 Feature extraction

The strength of the further processing depends on the quality and quantity of this features, so in other to improve both, we apply some preprocessing to the image. We tried several preprocessings techniques like sharpening, image contrast, median filter, and equalization. In 6.3, we can observe the relation between number of points extracted and time consumption of those techniques.

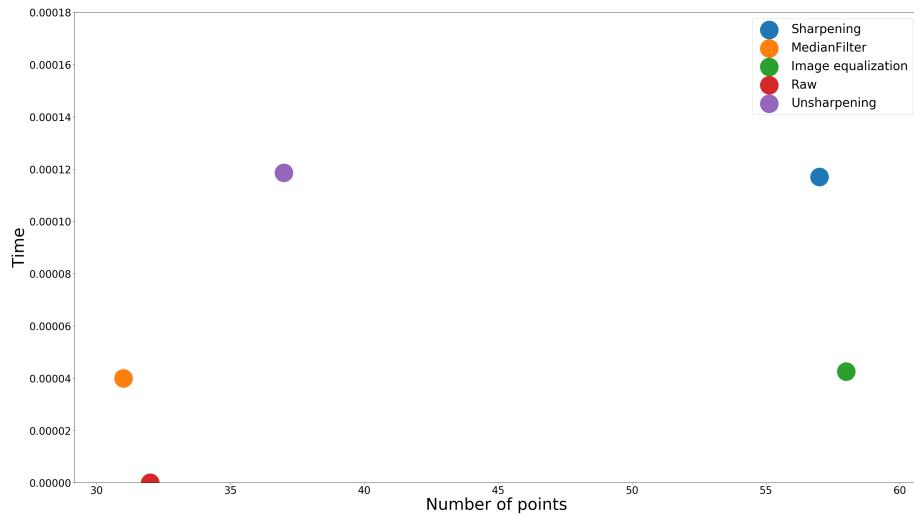


Figure 6.5: Barplot of the timming.

We realized that the best preprocessing in terms of speed and number of points, is to equalize the image. The computation is really simple, it only consists in equalize an histogram and apply that transformation to the image. It increases over 55% the number points in comparision to not apply it to the raw image. In the figure 6.6 we can observe the different number of features in the raw and in the equalized image.



(a) Raw image. (b) Equalized image.

Figure 6.6: Different preprocessings.

### 6.1.2.2 Matching module

Forward backward

We used the Lucas-Kanade method to matching the points. Also, we implement the same method used in [60], the proposed method is based on so called forward-backward consistency assumption. That assumption consists in that correct tracking should be independent of the direction of time-flow. Algorithmically, the assumption is exploited as follows. First, a tracker produces a trajectory by tracking the point *forward* in time. Second, the point location in the last frame initializes a validation trajectory. The validation trajectory is obtained by *backward* tracking from the last frame to the first one. Third, the two trajectories are compared and if they differ significantly, the forward trajectory is considered as incorrect. 6.7 illustrates the method when tracking a point between two images. Point number 1 is visible in both images and the tracker is able to localize it correctly. Tracking this point forward or backward results in identical trajectories. On the other hand, point number 2 is not visible in the right image and the tracker localizes a different point. Tracking this point backward ends in a different location than the original one. We also implemented the forward method.

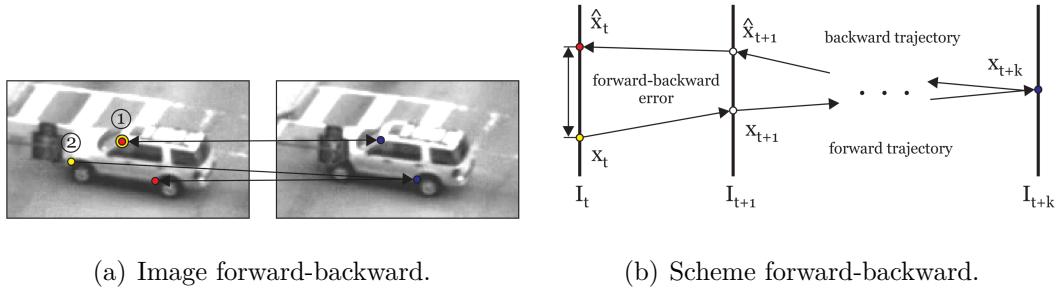


Figure 6.7: Illustration forward backward error.

### 6.1.3 Siamese networks

To maintain the identity of the pedestrian we need a method to compare miss tracks with no associated detections. So, we decided to solve it with deep learning techniques, the scheme that we considered are the following:

- **Siamese network: Cost function**, this is based on the idea of deep learning as feature extractor and top layers as classifier. Two branches that share parameters process the images and classify it.

- **Siamese network: In-network**, this is a mix of the previous models, where the information of the convolutional layers merges at some point before the classifier.
- **Siamese network: Joint data input**, according to the literature this architecture gives the best results compared with the other topologies. The input of the network is a concatenation of the two images and the network processes together.
- **Feature extractor with cosine distance**, we used well-known architectures for image classification to extract features from the images and then compare those features with the cosine distance.
- **Famous network fine-tuned**, we extract features for each image with a well-known architecture and merge it with a fully connected layer.

We can observe these architectures in the figure 6.8.

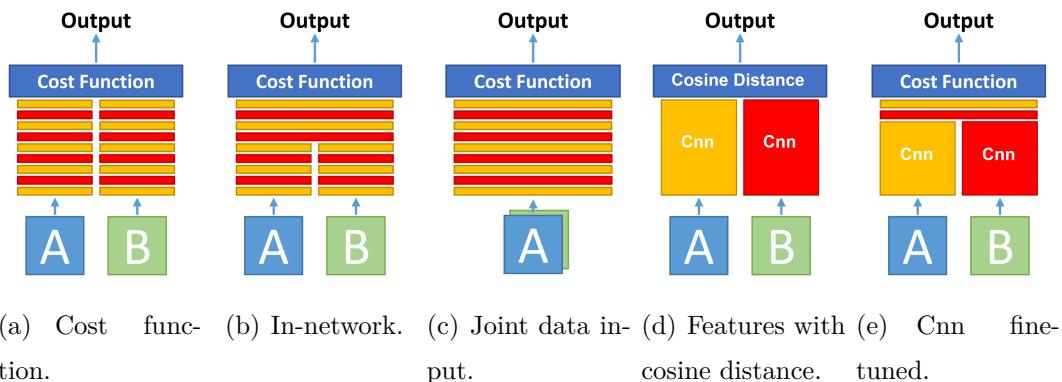


Figure 6.8: Siamese CNN topologies.

The main characteristics of the trained networks are the following:

- **Loss**, we used the binary cross entropy as a loss, we tried with the contrastive divergence but it did not converge.
- **Optimizer**, As optimizer we used Adam, even though it has a mechanism to decrease the learning rate, we add an exponential decay, it speeds up the convergence
- **Activation**, we used ReLu. Currently, there are other activation functions, but ReLu has been established as the reference.

- **Inizialization**, To initialize the wrights we used He. inizialization, in adition we initialize the biases with the value of 0.1, in this way we avoid the dead neurons in the firsts iterations.
- **Batch normalization**, We tested batch normalization, but it adds to much computation time and we discarded it.
- **Regularization**, we use Dropout in the fully connected layer to avoid overfitting.
- **Final layers**, In the junction between the convolutional layers and the fully connected historically, a flatten mechanism of the tensor has been used, but it increases dramatically the number of neurons in the fully connected layer, and it shows problems to converge. From the publication of InceptionV3, it appears with a global average pooling layer, it computes the spatial average of each layer of the tensor, reducing the number of parameters. Also we used the spatial pyramid pooling layers, it consists in a multiresolution max pooling. We can observe those differences in refsiameseData2.

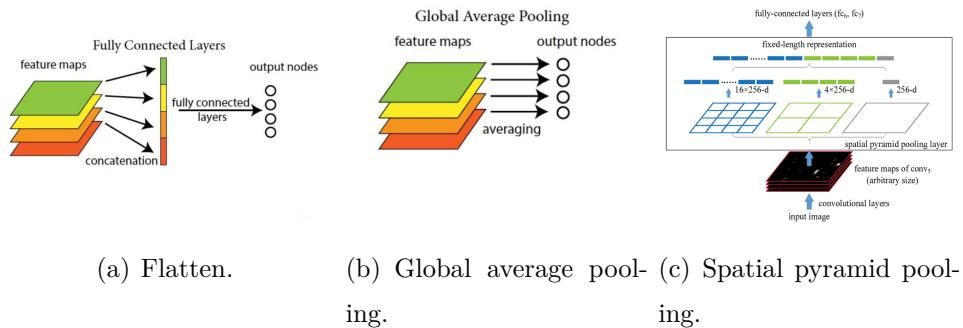


Figure 6.9: Final layers.

- **Output**, We did not use softmax as ouput, we only used one neuron with sigmoid activation, in this way the output is constraintet between 0 and 1.

We developed our models in a VGG way, stacking several convolutional layers and finishing it with a fully connected layers. We started with a few convolutional layers and added more till we reach an overfitting condition, this conditional will manifest when adding more layers the score on the test set decline. We started with 3 and we end up with 7 convolutional layer as best performance.

For the dataset, it does not exist a prominent dataset in the field, so we decided to use the MOT16 as dataset to adapt the domain. In order to do so, we extract the detections with their identities, and then for each identity we selected all possible random pairs and for the negative set we selected several random identities. The negative dataset is much bigger than the positive dataset, so we limited it to have a balanced dataset. The problem with the MOT16 dataset, is that the ground truth was built with the detections of a classifier and there is not a human intervention, resulting in a messy ground truth. We inspected the dataset and around the 70% of the dataset was wrong, there are a lot of occlusion in detection resulting in erroneous pairs, pairs that are not matching with the same identity. Then, we decided to discard the MOT16 dataset and use the TownCenter dataset [86] from the University of Oxford, which has got a manual ground truth. We have got 29824 positive and negative pairs, then a dataset of 59648 image pairs. We split the dataset between training and validation set, 80% and 20% respectively. For testing we selected a set of identities of the MOT16 dataset. To regularize and enlarge our dataset we applied some data augmentation techniques to our dataset like we observe in figure 6.10. For each pair we added one transformation, so we double our dataset. We tried to apply all the transformation for each images but the dataset was too noisy and the network did not converge.

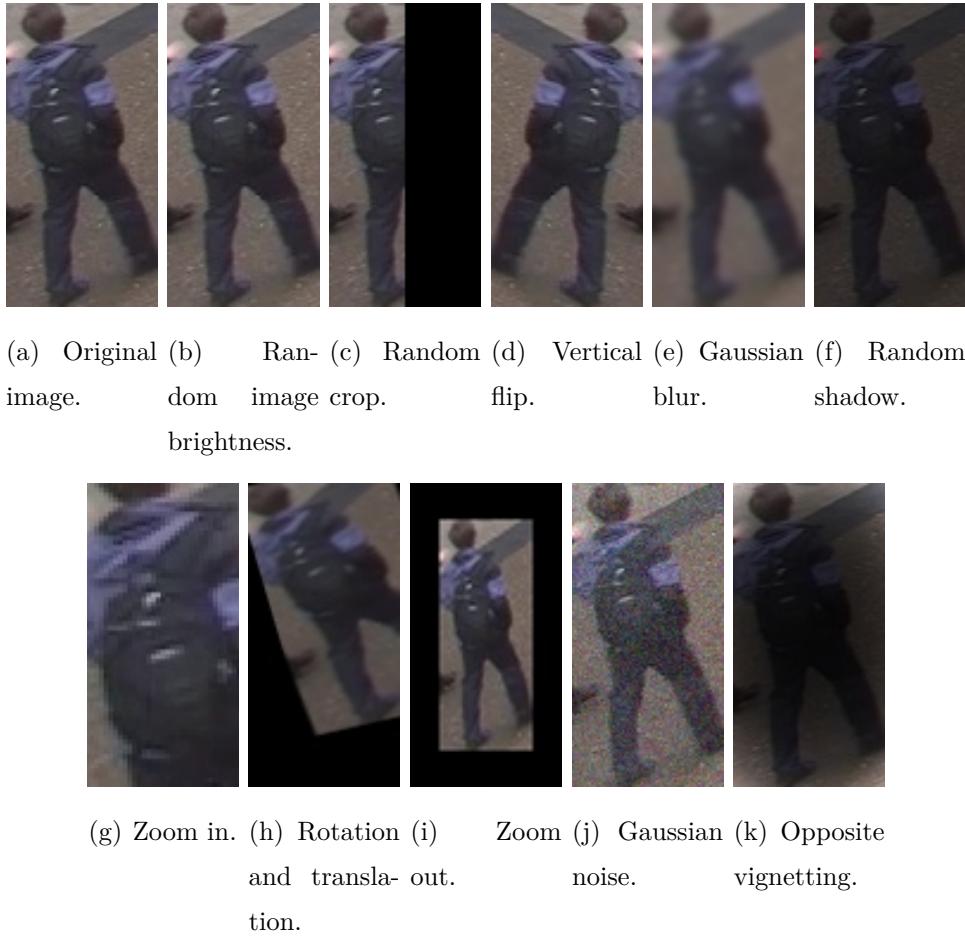


Figure 6.10: Data augmentation.

We trained all the models and obtained a graphs like the figure 6.11, we observe that the network converge, it decreases the loss and increases the accuracy, also we observe that tests plots are a little bit noisy, but we considered that the regularization techniques are enough. We notice that the joint data input outperforms the other siamese configurations, so we increased the number of layers of that architecture, *conv I*, refers to it, with  $I$  the number of layers.

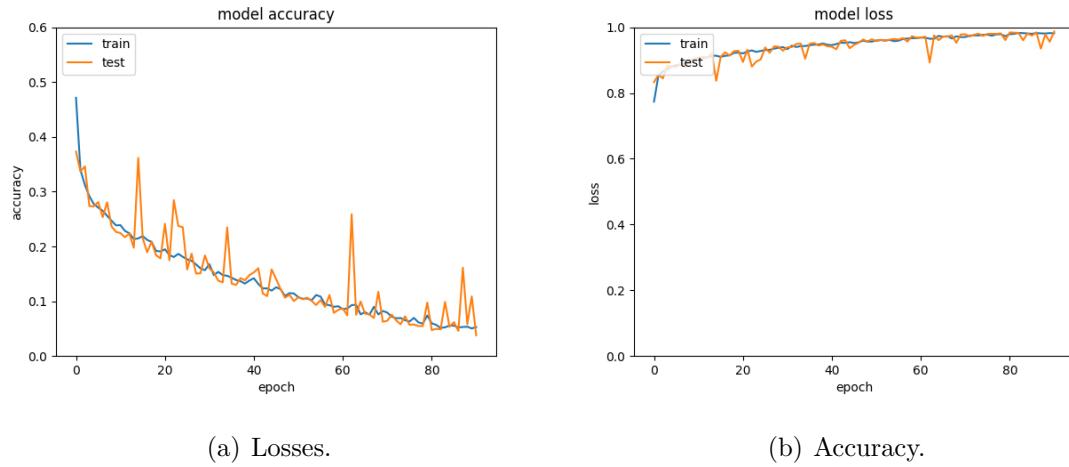


Figure 6.11: Results training.

Finally, we can observe the comparision using the CMC measure in the figure 6.12, we notice that the siamese network with joint data input with less layers than bigger models like Inception performs better, this remarks the idea of training jointly the feature extractor and the classifier and the need of task specific networks. Also, the siamese network with the configuration joint data input, outperforms the other siamese networks. Among the siamese joint data input, the performance increases till the 7 convolutional layer architecture, then the 8 convolutional layers drops.

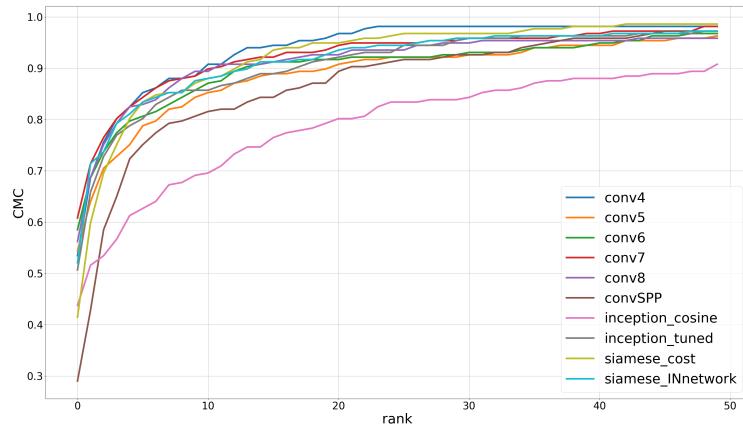


Figure 6.12: CMC plot.

Also in the next plot, we observe the performance against the time consumption the siamese network with the joint data input with 7 convolutional layers gets the best balance

between performance and execution time.

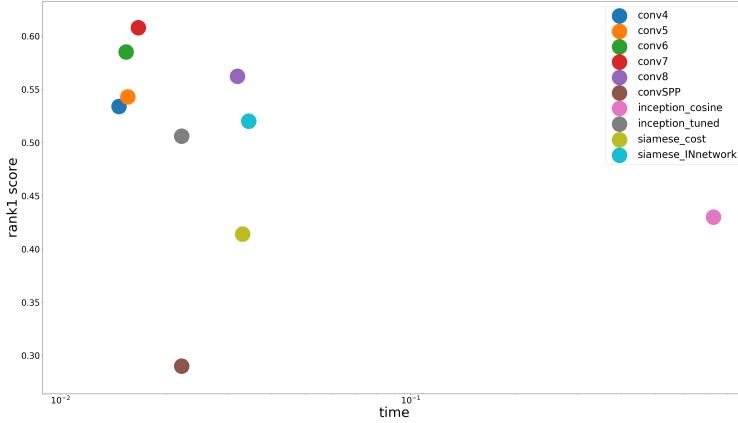


Figure 6.13: Performance-timing comparision.

## 6.2 Algorithm evaluation

In this section, we will analyse the performance and the timings of our solution.

### 6.2.1 Performance

Using the code provided by the MOT16 challenge organization, we evaluate our solution, we tested the version without re-identification module,  $v1$ , the version with optical flow forward-backward  $v2$ , and the version with re-identification module and without forward-backward method, also we compare with the state of the art, an anonymous submission to the MOT16's website. The results are showed in table 6.1.

	Rcll	Prcn	FAR	GT	MT	PT	ML	FP	FN	IDs	FM	MOTA	MOTP	MOTAL	FPS
$v1$	21.4	49.5	5.17	517	12	180	325	13339	78998	827	1053	9.8	69.1	6.5	18.2
$v2$	21.3	50.4	5.12	517	11	181	325	11212	75738	827	1056	9.7	67.3	6.1	9.0
$v3$	19.6	61.8	2.52	517	3	127	387	13373	78999	618	936	10.8	70.3	7.5	15.85
SOTA	-	-	-	517	92	219	206	5333	86795	391	-	49.3	79.0	-	0.8

Tabla 6.1: Results algorithm global.

According to this results, we decided to do not use the forward-backward method, due for its high computation demands and its limited improvement. Also, we observe the

usage of a mechanism of re-identification increment the MOTA performance, so the overall performance of the detector, with the cost of increase the running time. This improvement come from the reduction in around 24% the identity switching (ID's) parameter.

Finally we can observe the results in all the sequences 6.2, we get a MOTA of 10.8 and a frame rate of 15.85 FPS.

	Rcll	Prcn	FAR	GT	MT	PT	ML	FP	FN	IDs	FM	MOTA	MOTP	MOTAL	FPS
02	12.9	51.4	3.63	54	0	13	41	2181	15526	113	146	0.1	67.1	0.7	9.02
04	28.5	71.2	5.23	83	0	41	42	5495	33980	185	290	16.6	71.1	17.0	12.3
05	30.9	42.4	3.41	125	3	43	79	28571	4713	83	109	-12.2	67.8	-11.1	17.94
09	38.7	68.6	1.78	25	1	19	5	932	3225	62	71	19.7	71.2	20.9	10.52
10	5.4	62.4	0.62	54	0	4	50	404	11647	81	133	1.5	68.4	2.1	14.23
11	19.7	65.6	1.05	69	0	16	53	948	7366	72	121	8.6	71.4	9.4	17.49
13	6.2	34.9	1.75	107	0	9	98	1315	10743	32	59	-5.6	67.1	-5.4	20.5
Global	19.6	61.8	2.52	517	3	127	387	13373	88754	618	936	10.8	70.3	7.5	15.85

Tabla 6.2: Results algorithm by sequences.

The algorithm gets the best performance on sequences with a fixed camera from an elevated view point like sequences 4 and a low angle recording with a high frame rate and close targets like sequences 9, *and* 11, but it fails in sequences with low frame rate like sequences 5 and when the targets are away from the camera like 13.

The algorithm gets the best performance on sequences with a fixed camera from an elevated view point and a low angle recording with a high frame rate and close targets like sequences 4, 9, *and* 11, but it fails in sequences with low frame rate like sequences 5 and when the targets are away from the camera like 13.

### 6.2.2 Timing

As we stated above the mean frame rate of the algorithm with the person re-identification mechanism is 15.86. In figure 6.14 we can observe a barplot of per frame time consumption of our algorithm. We notice the peaks each 30 frames, these belong to the execution of the siamese network and it depends on how many detections without assignment there are.

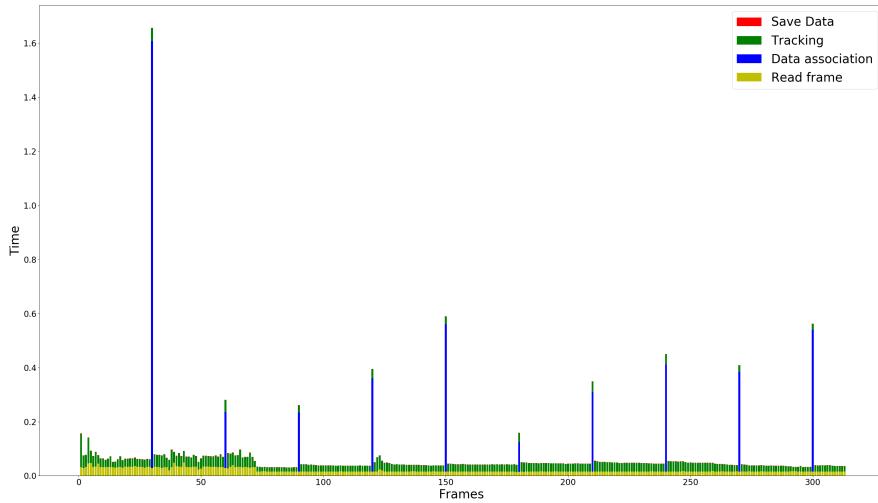


Figure 6.14: Barplot of the timming.

Getting a zoom in of the graph 6.16, we can observe when the object detector execution finishes around the frame 75, then the execution time of the algorithm decresases, the time of reading the frame remain constant, and the tracking gets a peak after the detection and afterwards decreases this is due it erase trackets with the lost mechanism.

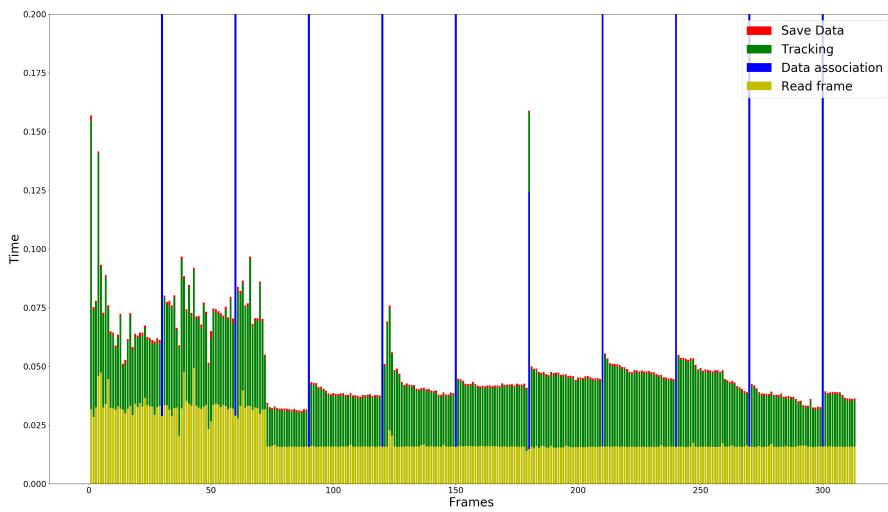


Figure 6.15: Zoom in of the barplot.

To optimize our code we studied how to reduce the execution time of the tracking module, plotting the execution time against some dependent variables, in this case,

the number of points and the size of the ROI, as we can observe in figure 6.16. We notice that the execution time is high correlated with the size of the pedestrian's ROI and not with the number of points. The main responsible is the OpenCV's routine `calcOpticalFlowPyrLK()`, but we could not modify this parameter, and reducing the number of points would have a remarkable importance in the execution time.

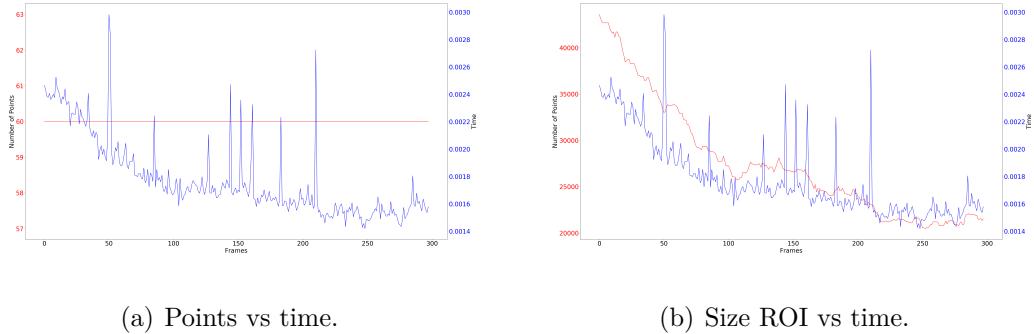


Figure 6.16: Time versus points and size of the ROI.

# Chapter 7

## Conclusions

In this thesis we studied deep learning techniques and their application in a tracking algorithm. We included an object detector and a person re-identification module based on those techniques and we implemented a tracking algorithm.

This tracking algorithm, takes the detections of the object detector and computes the optical flow on the sequences. When it misses some pedestrians it can preserve their identity with a person re-identification module.

Finally we evaluated our algorithm in a well-known challenge, MOT16, and analyses its performance and timing capabilties on it. The algorithms performs reasonably well in sequences of high frame rate and resolution, but in the opposite sequences, low frame rate and resolution the performance drop dramatically. This is happens beacause our algorithm rely on points with high texture, with a low resolution there are less points with this characteristics. Also, we have problems with people who wear low texture or are away from the camera, as we can observe in figure 7.1. The low frame rate, penalize the matching capabilites between frames, it produces wrong matches between images, therefore, wrong displacements.



(a) High texture person.      (b) Low texture person.      (c) Far away person.

Figure 7.1: Differences texture examples.

## 7.1 Future work

This is a first entrance on trackings algorithm, we have a resanoable results. We can add some details to improve our results.

- Port to C++. We used a sripting programming languagge, if we switch to a compiled programing language we would increase the time perfomance.
- GPU implementation. Computing displacement for each tracket could be computed in a parallel way. They are independent from each other and are a low demanding.
- Propbabilistic framework. Include bayesian filter techniques to increase perfomance.
- Siamese architectures. Study new siamese architecture to increase the accuracy of this modeule, like inception stem of InceptionV3 or include the optical flow information in the neural network.
- Data association. Use much confidence techniques to associate the detections, the current methods relay ond probablistic graphical models.

# Bibliography

- [1] A. W. M. Smeulders, D. M. Chu, R. Cucchiara, S. Calderara, A. Dehghan, and M. Shah. Visual tracking: An experimental survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7):1442–1468, July 2014.
- [2] Kai Briechele and Uwe D. Hanebeck. Template matching using fast normalized cross correlation, 2001.
- [3] Simon Baker, Ralph Gross, and Iain Matthews. Lucas-kanade 20 years on: A unifying framework: Part 3. *International Journal of Computer Vision*, 56:221–255, 2002.
- [4] H. T. Nguyen and A. W. M. Smeulders. Fast occluded object tracking by a robust appearance filter. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(8):1099–1104, Aug 2004.
- [5] Dorin Comaniciu, Visvanathan Ramesh, and Peter Meer. Real-time tracking of non-rigid objects using mean shift, 2000.
- [6] S. Hare, S. Golodetz, A. Saffari, V. Vineet, M. M. Cheng, S. L. Hicks, and P. H. S. Torr. Struck: Structured output tracking with kernels. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(10):2096–2109, Oct 2016.
- [7] Li Zhang, Yuan Li, and Ramakant Nevatia. Global data association for multi-object tracking using network flows. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8. IEEE, 2008.
- [8] Anton Milan, Stefan Roth, and Konrad Schindler. Continuous energy minimization for multitarget tracking. *IEEE transactions on pattern analysis and machine intelligence*, 36(1):58–72, 2014.
- [9] Zdenek Kalal, Krystian Mikolajczyk, and Jiri Matas. Tracking-learning-detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1409–1422, 2012.

## BIBLIOGRAFÍA

---

- [10] Federico Pernici and Alberto Del Bimbo. Object tracking by oversampling local features. *IEEE transactions on pattern analysis and machine intelligence*, 36(12):2538–2551, 2014.
- [11] Warren S. McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4):115–133, 1943.
- [12] B Widrow. An adaptive ‘adaline’neuron using chemical ‘memistors’, 1553–1552, 1960.
- [13] Marvin Minsky and Seymour Papert. *Perceptrons: An Introduction to Computational Geometry*. MIT Press, Cambridge, MA, USA, 1969.
- [14] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Parallel distributed processing: Explorations in the microstructure of cognition, vol. 1. chapter Learning Internal Representations by Error Propagation, pages 318–362. MIT Press, Cambridge, MA, USA, 1986.
- [15] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural Computation*, 1(4):541–551, Dec 1989.
- [16] Kunihiko Fukushima. Neocognitron: A hierarchical neural network capable of visual pattern recognition. *Neural Networks*, 1(2):119 – 130, 1988.
- [17] David H. Ackley, Geoffrey E. Hinton, and Terrence J. Sejnowski. A learning algorithm for boltzmann machines. *Cognitive Science*, 9(1):147 – 169, 1985.
- [18] Gerald Tesauro. Temporal difference learning and td-gammon. *Commun. ACM*, 38(3):58–68, March 1995.
- [19] Alexander Waibel, Toshiyuki Hanazawa, Geofrey Hinton, Kiyohiro Shikano, and Kevin J. Lang. Readings in speech recognition. chapter Phoneme Recognition Using Time-delay Neural Networks, pages 393–404. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1990.
- [20] Geoffrey E. Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural Comput.*, 18(7):1527–1554, July 2006.

## BIBLIOGRAFÍA

---

- [21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
- [22] M. Pieras. Wiki jderobot: Marcos pieras. <http://jderobot.org/Marcospieras-tfm> [2017-06-11], 2017.
- [23] M. Pieras. Github tfm: Marcos pieras. <https://github.com/RoboticsURJC-students/2015-TFM-Marcos-Pieras> [2017-06-11], 2017.
- [24] David Held, Sebastian Thrun, and Silvio Savarese. Learning to track at 100 FPS with deep regression networks. *CoRR*, abs/1604.01802, 2016.
- [25] Hyeonseob Nam and Bohyung Han. Learning multi-domain convolutional neural networks for visual tracking. *CoRR*, abs/1510.07945, 2015.
- [26] Martin Danelljan, Gustav Hager, Fahad Shahbaz Khan, and Michael Felsberg. Learning spatially regularized correlation filters for visual tracking. In *The IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [27] Ran Tao, Efstratios Gavves, and Arnold W. M. Smeulders. Siamese instance search for tracking. *CoRR*, abs/1605.05863, 2016.
- [28] Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. Tracking the untrackable: Learning to track multiple cues with long-term dependencies. *arXiv preprint arXiv:1701.01909*, 2017.
- [29] Jonathan Huang, Vivek Rathod, Chen Sun, Menglong Zhu, Anoop Korattikara, Alireza Fathi, Ian Fischer, Zbigniew Wojna, Yang Song, Sergio Guadarrama, and Kevin Murphy. Speed/accuracy trade-offs for modern convolutional object detectors. *CoRR*, abs/1611.10012, 2016.
- [30] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 91–99. Curran Associates, Inc., 2015.

## BIBLIOGRAFÍA

---

- [31] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.
- [32] Ross Girshick. Fast r-cnn. In *The IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [33] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.
- [34] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [35] Kye-Hyeon Kim, Yeongjae Cheon, Sanghoon Hong, Byung-Seok Roh, and Minje Park. PVANET: deep but lightweight neural networks for real-time object detection. *CoRR*, abs/1608.08021, 2016.
- [36] Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. *CoRR*, abs/1506.02640, 2015.
- [37] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott E. Reed, Cheng-Yang Fu, and Alexander C. Berg. SSD: single shot multibox detector. *CoRR*, abs/1512.02325, 2015.
- [38] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-FCN: object detection via region-based fully convolutional networks. *CoRR*, abs/1605.06409, 2016.
- [39] Chris Harris and Mike Stephens. A combined corner and edge detector. In *Alvey vision conference*, volume 15, pages 10–5244. Citeseer, 1988.
- [40] Jianbo Shi and Carlo Tomasi. Good features to track. In *Computer Vision and Pattern Recognition, 1994. Proceedings CVPR'94., 1994 IEEE Computer Society Conference on*, pages 593–600. IEEE, 1994.
- [41] Bruce D. Lucas and Takeo Kanade. An iterative image registration technique with an application to stereo vision. In *Proceedings of the 7th International Joint Conference*

- on Artificial Intelligence - Volume 2*, IJCAI'81, pages 674–679, San Francisco, CA, USA, 1981. Morgan Kaufmann Publishers Inc.
- [42] Liang Zheng, Yi Yang, and Alexander G. Hauptmann. Person re-identification: Past, present and future. *CoRR*, abs/1610.02984, 2016.
- [43] Douglas Gray and Hai Tao. *Viewpoint Invariant Pedestrian Recognition with an Ensemble of Localized Features*, pages 262–275. Springer Berlin Heidelberg, Berlin, Heidelberg, 2008.
- [44] Alexis Mignon. Pcca: A new approach for distance learning from sparse pairwise constraints. In *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, CVPR '12, pages 2666–2672, Washington, DC, USA, 2012. IEEE Computer Society.
- [45] R. Zhao, W. Ouyang, and X. Wang. Learning mid-level filters for person re-identification. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 144–151, June 2014.
- [46] Joost van de Weijer, Cordelia Schmid, Jakob Verbeek, and Diane Larlus. Learning color names for real-world applications. *Trans. Img. Proc.*, 18(7):1512–1523, July 2009.
- [47] M. Köstinger, M. Hirzer, P. Wohlhart, P. M. Roth, and H. Bischof. Large scale metric learning from equivalence constraints. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2288–2295, June 2012.
- [48] Kilian Q. Weinberger and Lawrence K. Saul. Distance metric learning for large margin nearest neighbor classification. *J. Mach. Learn. Res.*, 10:207–244, June 2009.
- [49] Shengcai Liao, Yang Hu, and Stan Z. Li. Joint dimension reduction and metric learning for person re-identification. *CoRR*, abs/1406.4216, 2014.
- [50] M. Chatterjee, Y. Luo. Similarity learning with cnn. [http://slazebni.cs.illinois.edu/spring17/lec09\\_similarity.pdf](http://slazebni.cs.illinois.edu/spring17/lec09_similarity.pdf) [2016-12-15], 2016.
- [51] Jane Bromley, Isabelle Guyon, Yann Lecun, Eduard Säckinger, and Roopak Shah. Signature verification using a "siamese" time delay neural network. In *In NIPS Proc*, 1994.

## BIBLIOGRAFÍA

---

- [52] S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with application to face verification. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 1, pages 539–546 vol. 1, June 2005.
- [53] Nam N. Vo and James Hays. Localizing and orienting street views using overhead imagery. *CoRR*, abs/1608.00161, 2016.
- [54] Edgar Simo-Serra, Eduard Trulls, Luis Ferraz, Iasonas Kokkinos, Pascal Fua, and Francesc Moreno-Noguer. Discriminative learning of deep convolutional feature point descriptors. In *The IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [55] Sergey Zagoruyko and Nikos Komodakis. Learning to compare image patches via convolutional neural networks. *CoRR*, abs/1504.03641, 2015.
- [56] Gregory Koch, Richard Zemel, and Ruslan Salakhutdinov. Siamese neural networks for one-shot image recognition. 2015.
- [57] Sean Bell and Kavita Bala. Learning visual similarity for product design with convolutional neural networks. *ACM Trans. on Graphics (SIGGRAPH)*, 34(4), 2015.
- [58] Philipp Fischer, Alexey Dosovitskiy, Eddy Ilg, Philip Häusser, Caner Hazirbas, Vladimir Golkov, Patrick van der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional networks. *CoRR*, abs/1504.06852, 2015.
- [59] Laura Leal-Taixé, Cristian Canton-Ferrer, and Konrad Schindler. Learning by tracking: Siamese CNN for robust target association. *CoRR*, abs/1604.07866, 2016.
- [60] Zdenek Kalal, Krystian Mikolajczyk, and Jiri Matas. Forward-backward error: Automatic detection of tracking failures. In *Pattern recognition (ICPR), 2010 20th international conference on*, pages 2756–2759. IEEE, 2010.
- [61] Jay Rambhia. Median flow tracker in python. <https://github.com/jayrambhia/MFTracker> [2015-08-15], 2015.

## BIBLIOGRAFÍA

---

- [62] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *International Journal of Computer Vision*, 111(1):98–136, January 2015.
- [63] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- [64] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014.
- [65] Gerard Salton and Michael J. McGill. *Introduction to Modern Information Retrieval*. McGraw-Hill, Inc., New York, NY, USA, 1986.
- [66] J. Ferryman and A. Ellis. Pets2010: Dataset and challenge. In *Proceedings of the 2010 7th IEEE International Conference on Advanced Video and Signal Based Surveillance, AVSS ’10*, pages 143–150, Washington, DC, USA, 2010. IEEE Computer Society.
- [67] Laura Leal-Taixé, Anton Milan, Ian D. Reid, Stefan Roth, and Konrad Schindler. Motchallenge 2015: Towards a benchmark for multi-target tracking. *CoRR*, abs/1504.01942, 2015.
- [68] Rainer Stiefelhagen, Keni Bernardin, Rachel Bowers, John Garofolo, Djamel Mostefa, and Padmanabhan Soundararajan. *The CLEAR 2006 Evaluation*, pages 1–44. Springer Berlin Heidelberg, Berlin, Heidelberg, 2007.
- [69] Bo Wu and R. Nevatia. Tracking of multiple, partially occluded humans based on static body part detection. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, volume 1, pages 951–958, June 2006.
- [70] UC Santa Cruz. Viper: Viewpoint invariant pedestrian recognition. <https://vision.soe.ucsc.edu/node/178> [2016-12-15], 2007.
- [71] T. Wang, S. Gong, X. Zhu, and S. Wang. Person re-identification by discriminative selection in video ranking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(12):2501–2514, Dec 2016.

## BIBLIOGRAFÍA

---

- [72] University of London. Qmul underground re-identification. [http://personal.ie.cuhk.edu.hk/~ccloy/downloads\\_qmul\\_underground\\_reid.html](http://personal.ie.cuhk.edu.hk/~ccloy/downloads_qmul_underground_reid.html) [2016-12-15], 2009.
- [73] Doug Gray, Shane Brennan, and Hai Tao. Evaluating appearance models for recognition, reacquisition, and tracking. In *In IEEE International Workshop on Performance Evaluation for Tracking and Surveillance, Rio de Janeiro*, 2007.
- [74] Martin Hirzer, Csaba Beleznai, Peter M. Roth, and Horst Bischof. Person re-identification by descriptive and discriminative classification. In *Proc. Scandinavian Conference on Image Analysis (SCIA)*, 2011.
- [75] Abir Das, Anirban Chakraborty, and Amit K. Roy-Chowdhury. *Consistent Re-identification in a Camera Network*, pages 330–345. Springer International Publishing, Cham, 2014.
- [76] Wei Li, Rui Zhao, and Xiaogang Wang. *Human Reidentification with Transferred Metric Learning*, pages 31–44. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013.
- [77] W. Li and X. Wang. Locally aligned feature transforms across views. In *2013 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3594–3601, June 2013.
- [78] 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*, pages 152–159, June 2014.
- [79] N. Martinel and C. Micheloni. Re-identify people in wide area camera network. In *2012 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, pages 31–36, June 2012.
- [80] Peter M. Roth, Martin Hirzer, Martin Köstinger, Csaba Beleznai, and Horst Bischof. *Mahalanobis Distance Learning for Person Re-identification*, pages 247–267. Springer London, London, 2014.
- [81] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *Computer Vision, IEEE International Conference on*, 2015.

## BIBLIOGRAFÍA

---

- [82] Patrick Grother, Ross J. Micheals, and P. Jonathon Phillips. *Face Recognition Vendor Test 2002 Performance Metrics*, pages 937–945. Springer Berlin Heidelberg, Berlin, Heidelberg, 2003.
- [83] J. Redmon. You only look once. <https://github.com/pjreddie/darknet> [2017-06-11], 2017.
- [84] Paul B. Single shot detector tensorflow. <https://github.com/balancap/SSD-Tensorflow> [2017-02-21], 2017.
- [85] Wei Lee. Single shot detector caffe. <https://github.com/weiliu89> [2017-02-21], 2017.
- [86] University of Oxford. Town centre dataset. [http://www.robots.ox.ac.uk/ActiveVision/Research/Projects/2009bbenfold\\_headpose/project.html#datasets](http://www.robots.ox.ac.uk/ActiveVision/Research/Projects/2009bbenfold_headpose/project.html#datasets) [2017-02-21], 2011.

# **Chapter 8**

## **Annex**