

Università degli Studi di Padova

DIPARTIMENTO DI MATEMATICA "TULLIO LEVI-CIVITA"

CORSO DI LAUREA IN INFORMATICA



Object detection with YOLO v3

Tesi di laurea triennale

Supervisor

Prof. Alessandro Sperduti

Candidate

Francesca Lonedo

ACADEMIC YEAR 2017-2018

To my parents, who gave me the opportunity to graduate twice.

Abstract

The following is a description of the work done by the candidate Francesca Lonedo during her three-hundred-and-twenty hours long curricular internship at THRON S.p.A.

Object detection is a computer vision related technology that deals with identifying instances of semantic class objects, and locating their position in digital images and videos. It has application in various domains, such as self-driving cars, video surveillance and image retrieval.

My project's goal was to develop a proof-of-concept prototype application to detect class instances and their location in image data. There are different approaches to solve such task, but I was asked to follow a *deep learning*_g method that uses a *convolutional neural network*_g called *yolo v3*_g. Furthermore, it was required to use a production-ready framework for development, and I was given the freedom to experiment with both *Mxnet* and *Tensorflow* and to choose whichever I found more suitable for the project.

“There is nothing either good or bad, but thinking makes it so”

— William Shakespeare’s Hamlet

Thanks

I want to thank my parents, for their endless love and support. They accepted my weird interests and gave me the opportunity to follow my dreams and graduate twice.

I want to thank my partner, who re-introduced me to computer science. He taught me a lot of things and always spoils me.

I want to thank my friends, old and new, who kept me company on this journey.

I want to thank my tutor and everyone at THRON for making my internship a wonderful learning experience.

Padova, December 2018

Francesca Lonedo

Contents

1	Introduction	1
1.1	THRON	1
1.2	Birth of the project	1
1.3	Content organization	2
2	Deep learning	3
2.1	Overview	3
2.2	Deep neural networks	5
2.2.1	Convolutional neural networks	5
2.2.2	Object detection	5
2.2.3	YOLO v3	6
2.3	Training a network	7
2.3.1	Dataset	8
2.3.2	Training	9
2.3.3	Validation	9
2.3.4	Transfer learning	9
2.3.5	Hyperparameters	10
2.4	Object detection for THRON	10
3	Internship description	11
3.1	Project introduction	11
3.2	Risk analysis	11
3.3	Goals and requirements	12
3.4	Schedule	13
4	Analisi dei requisiti	15
4.1	Casi d'uso	15
4.2	Tracciamento dei requisiti	16
5	Design and development	19
5.1	Technology and tools	19
5.2	Design	21
5.2.1	Overview	21
5.2.2	Dataset creation	22
5.2.3	Training	22
5.2.4	Inference service	25
5.2.5	Standalone inference script	27
6	Verifica e validazione	29

7 Conclusioni	31
7.1 Consuntivo finale	31
7.2 Raggiungimento degli obiettivi	31
7.3 Conoscenze acquisite	31
7.4 Valutazione personale	31
A Appendice A	33
Bibliografia	37

List of Figures

2.1	Visual representation of the activation function that determines a neuron's output and the parameters that influence it.	4
2.2	Schematic representation of how a network updates the weight to decrease the loss and achieve better accuracy.	4
2.3	Visual representation of a convolutional layer.	6
2.4	Comparison on object detection performances on Nvidia Titan X GPU. A higher mAP value means a better accuracy.	7
2.5	Visual representation of <i>YOLO</i> network layers in its first version. . . .	8
4.1	Use Case - UC0: Scenario principale	15
5.1	Structure of a *. <i>lst</i> file containing dataset information.	23
5.2	UML describing the relationship between the server application (App class) and the implementation of its methods.	25
5.3	UML flow chart displaying the HTTP message exchange in the application.	28

List of Tables

4.1	Tabella del tracciamento dei requisiti funzionali	17
4.2	Tabella del tracciamento dei requisiti qualitativi	17
4.3	Tabella del tracciamento dei requisiti di vincolo	17

Chapter 1

Introduction

1.1 THRON

THRON S.p.A.¹ is an Italian company that develops a *marketing DAM*_g software. Founded in 2000 as New Vision by CEO Nicola Meneghello and CTO Dario De Agostini, it was one of the pioneers at delivering web applications in Italy.

In 2004 the software house launched 4ME, a cloud based service for content management which was the precursor of the current company-homonym product.

Nowadays THRON has four premises around the globe and innovation remains one of its core values.

In the following I will refer to the company as THRON and to the product as Thron to distinguish them.

1.2 Birth of the project

THRON's homonymous product is a *marketing DAM*, a software application that can be seen as a centralized archive to manage multimedia content; any type of file can be stored, but contextually images are predominant. A DAM provides various features such as advanced item search, property management and tools to share content on social media and web platforms.

The core feature of a DAM is its ability to categorize the indexed media, making it easy for the user to find and edit objects with specific properties or tagged in a certain way. Thron makes this process even simpler thanks to its intelligence, a system that automatically elaborates the inserted media and suggests how to tag and categorize it. Currently Thron's intelligence is powered by general purpose *artificial intelligence*_g tools, hence the ability to recognize only a limited set of classes/features. Since most of Thron's users are big brand-name companies, teaching the system how to detect custom classes relevant for each client would improve greatly the software's relevance on the market.

As stated above, the vast majority of the stored media is in an image format; this both determines the problem and its solution: what is needed in Thron is an object detection system, and given the considerable database already available the most straightforward

¹<https://www.thron.com>

way to create it is to follow a deep learning² method.

1.3 Content organization

Il secondo capitolo describe ...

Il terzo capitolo approfondisce ...

Il quarto capitolo approfondisce ...

Il quinto capitolo approfondisce ...

Il sesto capitolo approfondisce ...

Nel settimo capitolo describe ...

This works adheres to the following typographic conventions:

- * acronyms, abbreviations, ambiguous and field-specific terms are defined in the glossary at the end of the document;
- * the first occurrence of glossary terms is emphasized as follows : *glossary term_g*;
- * technical terms are emphasized as follows: *technical term*.

²I will exhaustively explain what deep learning is and why it's the best solution for this problem in chapter two.

Chapter 2

Deep learning

In this chapter I will dive deep (no pun intended) in the reasons behind the choice of a deep learning approach to detect custom-class objects in image data.

2.1 Overview

Deep learning is a part of the broader family of *machine learning*_g methods that focuses on learning features of interest by experiencing them through data samples. In contrast to *machine learning*, which still needs the problem's layout to be described by a rather complex human-coded algorithm, a *deep learning* model only needs a dataset to autonomously learn from. Thanks to the ever-growing availability of computational power and labeled data, and the relative simplicity in which a network can be trained, the deep learning approach to solve artificial intelligence problems is flourishing, bringing artificial intelligence to everyday life.

At the core of every *deep learning* model there is an artificial neural network, whose architecture is strongly inspired by the structure and functioning of a biological brain. Basically, an artificial neural network is composed by a variable number of neuron layers densely connected and stacked upon each other, hence the name *deep neural networks*. The information is elaborated by traveling from layer to layer until the top is reached.

As stated, each layer is composed of computational units that work like *neurons*: each *neuron* receives an input, and using an activation function it produces an output that propagates to activate a *neuron* in the following layer. Actually each *neuron*'s output is influenced by weights and thresholds that need to be adjusted to achieve the best performance. When learning, the network indeed adjusts the weights and thresholds that activate each *neuron*'s function, so that it can get closer to the wanted result in an iterative manner. After each passage through the network, a *loss function*_g is calculated to determine how well it is performing; then a *gradient descent function*_g derives how this performance can be improved by decreasing the loss and accordingly *back-propagation*_g is performed to adjust the weights used by the neurons in their activation functions.

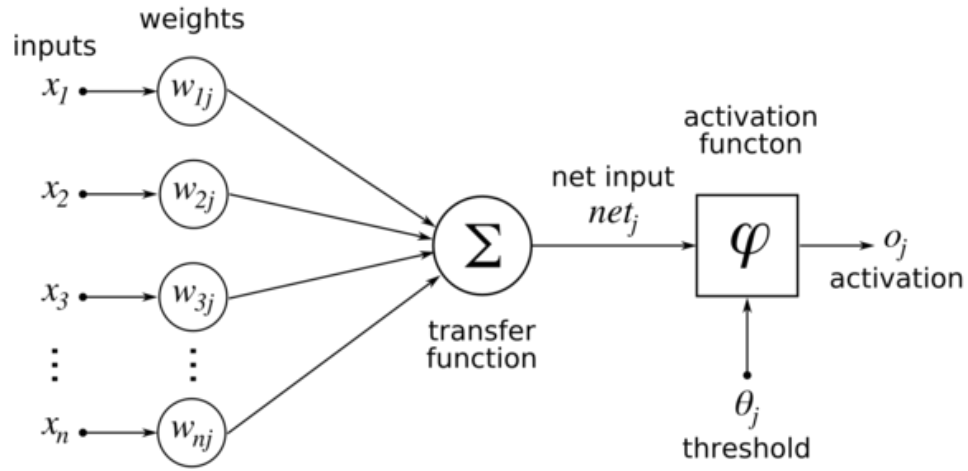


Figure 2.1: Visual representation of the activation function that determines a neuron's output and the parameters that influence it.

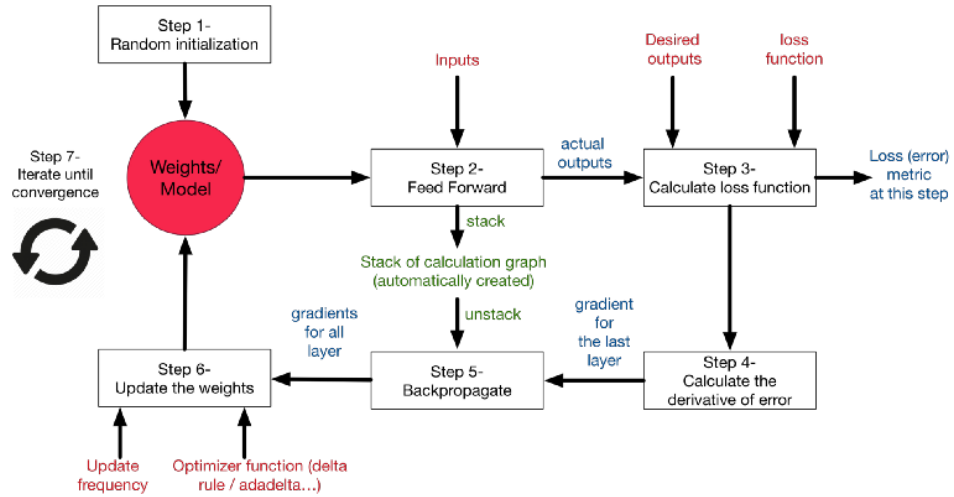


Figure 2.2: Schematic representation of how a network updates the weight to decrease the loss and achieve better accuracy.

2.2 Deep neural networks

As mentioned above, the *deep learning* approach to solve a problem is to train a *deep neural network* to recognize the wanted features so it can later make predictions about them. Basically when you train a *deep neural network* you feed it thousands of data samples, so it can learn from the wanted data representation by "seeing" it, just like you would show a child pictures of kittens to teach him what a kitten is. There are two different ways to train a network, each one with different goals and contexts of use:

- * supervised training: the network is given a labeled *dataset*_g to learn from; for each data sample the ground truth is provided, so the network knows what it is learning. This is the preferred approach whenever labeled data is available.
- * unsupervised training: the network is given an unlabeled *dataset* to learn from; no ground truth is given, so the network has to identify patterns and divide them into different categories on its own. This approach is preferred when the goal is to group data in categories or when data is too complex to be labeled.

There are different network architectures, each one dedicated to a specific type of problems. When solving *computer vision*_g task the best performances are held by a class of networks called *convolutional neural networks*.

2.2.1 Convolutional neural networks

convolutional neural networks (CNNs) are *deep neural networks* that consist of an input layer, an output layer, and a variable number of hidden layers with different purposes. What names this family of networks is in fact on particular type of layer, called *convolutional layer*; a single network typically contains various *convolutional layers* of different size, and each perform a *convolution* operation that filters its input stimuli before passing them to the next layer in order to reduce the parameters number, thus allowing the network to be deeper.

CNNs, among others, are mainly used to perform computer vision tasks; in particular they achieve good results in object detection tasks.

2.2.2 Object detection

Object detection is a computer vision problem that concerns the identification of class instance objects in an image (or video), and locating the actual position of said object in the picture. Commonly the spacial orientation of the detected object is framed by a rectangular *bounding box* that determines its height and width. Thus, object detection is far more powerful than mere image classification, not only because it "draws" a box where the object is located, but also because it can identify multiple object instances in a single image, while classification models have the limit of labeling only the one predominant object in the scene. The capability of labeling and locating multiple instances opens object detection models to a new set of application, such as video surveillance (e.g. moving subjects can be tracked) and instance counting for industrial purposes (e.g. counting boxes in a warehouse). Furthermore, object detection has proven to be more reliable than classification to scan images where the subject of interests occupies only a small part of the picture (e.g. a street sign in the corner). Commonly an *object detection* model is actually built on top a classifier that works

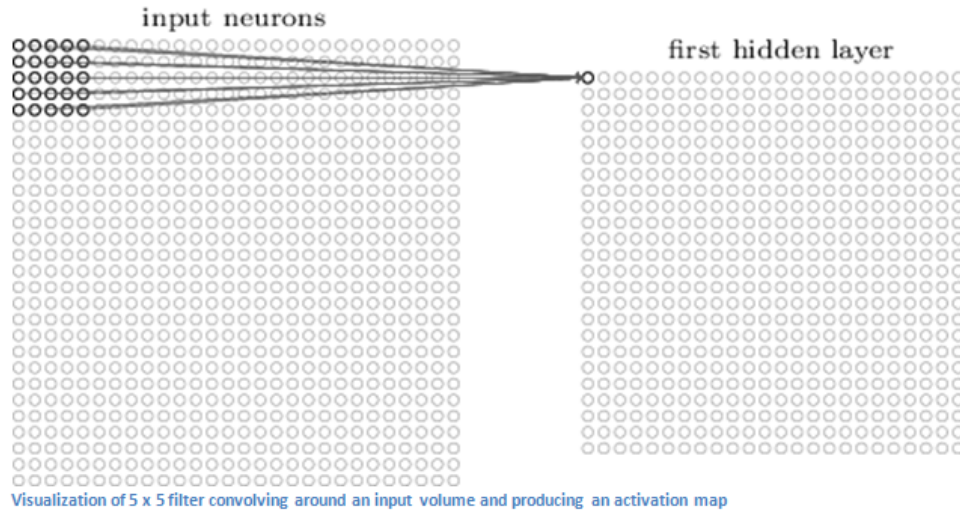


Figure 2.3: Visual representation of a convolutional layer.

as a *feature extractor*, but this exudes the topics of my project and I am introducing this notion only because it is noticeable in the naming convention of model networks, which are indeed important in my work.

There are three popular network architectures for object detection:

- * SSD (Single Shot Detection);
- * R-CNN (Region-based *convolutional neural network*, and its upgrades Fast R-CNN and Faster R-CNN);
- * YOLO (You Only Look Once).

While both *SSD* and *YOLO* follow the same approach of computing the image data through a single network, *R-CNN* splits the image in various regions trying to guess where objects of interest might be and processes each area separately, thus requiring more computational power. State-of-the-art models share comparable performances in accuracy, but due to the time overhead of processing a picture multiple times, *R-CNN* is a bit slower than *SSD* and *YOLO*.

2.2.3 YOLO v3

*YOLO*¹ is a *convolutional neural network* built on top of *Darknet* classifier, now on its third version *YOLO v3*. Both *YOLO* and *Darknet* are developed by Joseph Redmon on his C-written open source neural network framework *Darknet*².

What makes *YOLO* stand out compared to its competitor is its speed: at par of accuracy *YOLO* can perform inference in less than half the time, achieving the ability to track moving subjects in real time in 30fps videos using a gaming-tier GPU³.

For my project I was asked to use *YOLO* in its newest version, *YOLO v3*. As I

¹<https://pjreddie.com/darknet/yolo/>

²To avoid misunderstandings I will call *Darknet* the framework and *Darknet* the classifier.

³Real time inference is possible on Nvidia GeForce Titan X GPU.

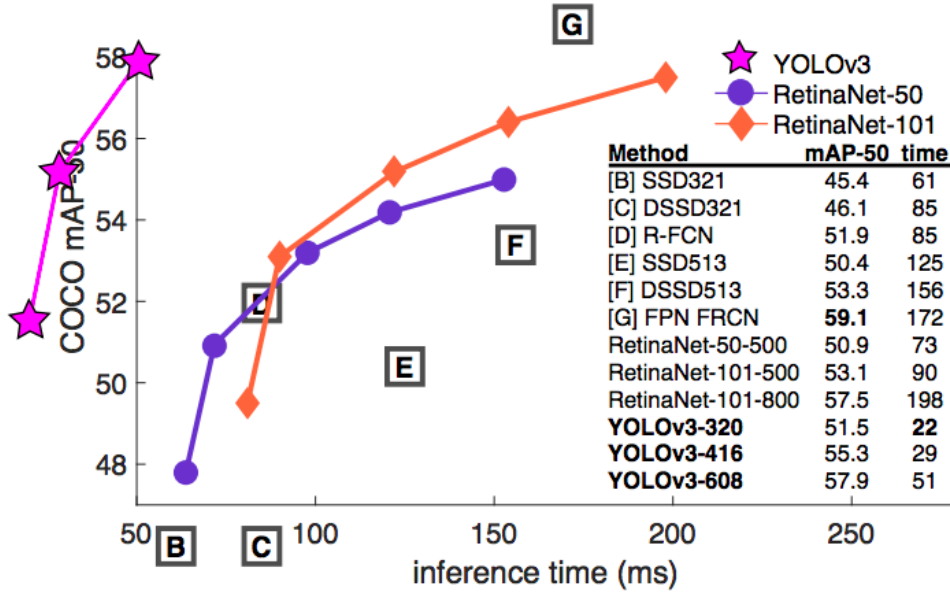


Figure 2.4: Comparison on object detection performances on Nvidia Titan X GPU. A higher mAP value means a better accuracy.

already said, *YOLO* was firstly implemented in Joseph Redmon’s Darknet framework, which is written in C so it is fast, but unfortunately it is suitable only for research purposes. Since I was asked to develop a production-ready application, I had to choose another framework to work on, and I will discuss about it in later chapters.

Structure

YOLO v3 is an *convolutional neural network* for *object detection* built on top of *Darknet53*, an image classifier that uses 53 convolutional layers. With *Darknet53* as its feature extractor, *YOLO* processes every image once, looking at it as a whole, hence encoding extra information about the context each class lives in. To do so it models detection as a regression problem, dividing the image into an $S \times S$ grid, and for each grid cell predicts B bounding boxes, confidence for those boxes, and C class probabilities; when a class-object’s center falls inside a grid cell, that grid cell is responsible of detecting that object.

2.3 Training a network

Training a *convolutional neural network* is a computation-heavy task, that involves the processing of thousands of pieces of data, which in the context of computer vision is in the format of images. Therefore best way to process thousands of images is to use a GPU; using a GPU (or even clusters of GPUs) it is possible to process the images in parallel among the *cores*, speeding up the process greatly. As a comparison, while training on CPU would take days even on a small dataset, on a GPU all the work would be done in just few hours.

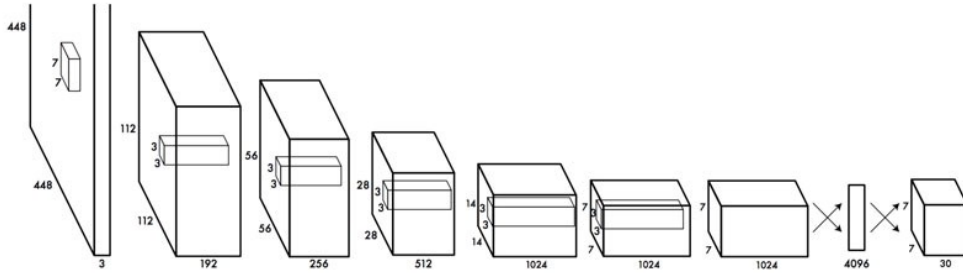


Figure 2.5: Visual representation of *YOLO* network layers in its first version.

2.3.1 Dataset

The first and most important step to train a network is to create the *dataset* it will learn from. Note that the quality of the dataset itself will have great influence on the final accuracy of the network. For instance you will want to feed your network examples of your objects of interest from every perspective and immerse in their context; also, when training a network to perform an *object detection* task, you want your *bounding boxes* to be as accurate and tight around the figures as possible.

Clearly, creating a custom dataset is a demanding job, since piece of data must be manually labeled by a human. For my project I was given a custom dataset created by one of my colleagues, containing classes relevant for one of THRON's clients. Creating a *dataset* exudes the goals of my project, since on production it will be possibly performed by clients⁴.

When creating an *object detection dataset* there are there are two main qualities to determine its goodness:

- * *Cardinality*, which is the means of the number of labels per image;
- * *Density*, which is the cardinality divided by the number of classes.

A density too low has negative influence on the final accuracy of the network, because basically you are trying to teach it to recognize a high number of classes, but you're giving it only a few examples for each. On a lesser degree, a cardinality too high has a negative influence on the final accuracy as well because the examples you are giving are too cluttered with information. Furthermore, when creating your own *dataset*, you might want to insert a similar amount image samples/labels for each class.

After the full dataset is created, it is common practice to divide it into two separate *datasets*: a *training dataset* and a *validation dataset*, with the suggested ratio of 70% training and 30% validation. Note that these two *datasets* contain different samples to avoid *overfitting*_g.

⁴As explained in the introduction, THRON's product is a DAM, which in this context can basically seen as a database of categorized images; with possibly minor changes to data labeling (e.g. introduction of bounding boxes for the objects of interest) and database export, image data can be easily turned into a dataset for computer vision.

2.3.2 Training

Once the dataset is ready and formatted in the network's preferred format⁵, the training process can begin.

The training process is conceptually simple; as previously stated, teaching a machine to recognize kittens can be compared to showing pictures to a child and pointing at the kitten in them to make him understand what a kitten is. When doing so with a child you only need a few examples, when with a machine you'll probably need a few hundreds. The next step is to check whether the machine truly learned what a kitten is; again this can be roughly compared to checking whether the child truly learned what a kitten is by giving him a picture and asking him to point at the kitten, if there is one.

Now that we talked about kittens and children we can explore in detail how a machine learns the way a human would. The training process is a loop in which the network is "shown" the *dataset* it should learn from; since a *dataset* is composed by thousands of images, it can't be entirely loaded into memory and must be split in smaller units for processing, called *batches*; in practice maximum *batch* size is determined by your hardware specifics. Your network will be fed all of the training *dataset batch* after *batch*. Each loop of computation on the whole training dataset is called an *epoch*_g. To train your network you will want to perform at least a couple of hundreds epochs; this means that to learn the *dataset* your network will go through all of it hundreds of times.

2.3.3 Validation

Every fixed epoch number you will want to check your network's learning progress; to do so you will perform a validation loop, which means that you will feed the network your validation *dataset* and compare the network's prediction with the ground truth. There are various metrics to calculate a network's accuracy, and the one I used on my project is called *mean average precision*_g (mAP). This evaluation metric uses an *intersection over union*_g (IoU) threshold to determine how well the detected positives overlap with the ground truth; setting the IoU threshold higher or lower influences the wanted precision at locating objects⁶, and I set it to a value of 0.5.

2.3.4 Transfer learning

Training a network from scratch is a quite demanding operation in terms of computational time and power, because the weights for each class need to be calculated from a random initialization. Luckily, it is possible to train a network to recognize new custom classes starting from pre-trained weights for other (general) classes, and this process is called *transfer learning*.

Basically when you apply transfer learning on a pre-trained network⁷, you teach it to recognize new classes by adjusting its weights (which have a meaning and are not random anymore) and the structure of the two layers that determine the final output. To benefit from *transfer learning* you want to choose a base network with weights for

⁵Depending on the framework, the preferred format changes. To work with YOLO v3 on MXNet data can be provided in both RecordIO format or LST + Jpeg format.

⁶Note that ground truth bounding boxes are hand-drawn by a human so they already come with an error, thus setting the IoU threshold too high could prove to be useless and even counterproductive.

⁷All the major deep learning frameworks usually provide networks and weights pre-trained on popular datasets such as COCO, Pascal VOC or ImageNet.

classes that share some features with your custom ones for example if you want to be able to detect different species of animals, a good choice would be to *fine-tune* a network trained on a datasets that already contains cats and dogs.

2.3.5 Hyperparameters

The accuracy of a *convolutional neural network* is influenced by various parameters, called hyperparameters, that tweak how the weights are updated and how the results are evaluated. These values come in the form of "*magic numbers*", meaning that for every dataset each network has its own values to optimize the final performance and that they can be found only in a trial and error fashion⁸. Since this is rather time consuming, especially when training on CPU, I couldn't tune the hyperparameters for my project. However I will introduce briefly what the main hyperparameters do on a *YOLO v3* architecture network.

- * Learning rate. Speed at which the network learns. It is used when updating weights and a learning rate too high might cause overfitting.
- * Learning rate decay rate. As the network learns, its learning rate decreases; this value is important to avoid overfitting.
- * Momentum. Stochastic gradient descent momentum, used by the gradient descent function to calculate how to reduce the loss.

2.4 Object detection for THRON

As mentioned in Chapter 1, THRON already uses generic classification services in its *DAM*. Therefore, to provide a better service, it comes natural that the Company has interest in developing a system able to detect custom classes.

Since THRON already has a vast database of labeled data, following a deep learning approach is the easiest option, even more so now that framework providing almost out-of-the-box are increasing prominence.

At first it was still unclear whether the Company needed just a classifier, but with my proof-of-concept the benefits of an *object detection* model became clear. We particularly chose *YOLO v3* under the suggestion of a data scientist working in my team, since he already developed a research-purpose prototype application using *YOLO v2*. Due to the technology used, his application couldn't be released in a production environment, so it was my job to develop a production-ready application using *YOLO*.

⁸To determinate the best values it is common practice to re-train the network various times and compare the results.

Chapter 3

Internship description

My project at THRON S.p.A. was to develop a production-ready application for object detection using YOLO v3. I was given the responsibility to choose the framework to use, and I coded both the training and inference parts of the application.

3.1 Project introduction

The goal of my project was to develop a production-ready *object detection* application that could be integrated in Thron in the future. Since THRON strongly relies on web services like those provided by [AWS](#)_g, so it was my job to keep compatibility.

3.2 Risk analysis

Every project comes with its risks. Since deep learning is a relatively new field, my project was at high risk from the beginning.

Decription	Solution	Occurrence
Stable releases unavailable The technology I used is fairly new, so it is not surprising that the main toolkit I used, <i>gluon cv_g</i> was actually in pre-release version 0.3.0. This lead to some inconveniences such as bugs and inconsistencies between build and documented source code.	No real solution was identified. When necessary I submitted issues to the code developers and patched bugs myself when I could.	High
AWS support Since the technology I used was fairly new, it relied on the newest version of my framework of choice <i>MXNet</i> 1.3.0. During my internship, AWS provided out-of-the-box support only to <i>MXNet</i> 1.2.1 and prior, which isn't compatible with <i>Gluon CV</i> 0.3.0.	No real solution was identified. I tried other AWS services, but they are not as straightforward as the one I was supposed to use and proved to be rather troublesome.	High
Hardware Deep learning is a demanding technology that requires a lot of computational power. Specifically, it is suggested to execute training on a <i>text_g</i> GPU since it is much faster, but I didn't have one. Training on CPU, on the other hand, is indeed possible, but it requires fairly new hardware to run at all, which, during my first week, I didn't have.	Before I was assigned a newer CPU I got to work on a remote AWS machine via SSH. This machine provided a GPU so I could test example algorithms. When given a newer CPU I had to train my network on that, which required a couple of days and thus prevented me from tuning the hyperparameters.	High
Python Deep learning frameworks are written almost exclusively in Python. With no prior experience, writing a Python application could prove to be troublesome.	None of my coworkers was a Python expert, but they always helped me to find good libraries and frameworks to use.	Medium

3.3 Goals and requirements

Requirements are categorized with the following criteria:

- * MUST, the functionality is absolutely required;
- * SHOULD, the functionality is recommended;
- * MAY, the functionality is optional.

Due to the nature of the project, whose purpose was to explore a new technology and determine the possibility of deploying an application in a production environment,

requirements and goals changed along the way. In fact, some requirements proved impossible to be satisfied because of the unavailability of tools and services supporting them yet.

* Required:

- MUST 01: development of a training function;
- MUST 02: integration with Thron APIs;

* Recommended:

- SHOULD 01: development of a prototype user interface;

* optional:

- MAY 01: development of unit tests;
- MAY 02: extraction of metrics to evaluate the algorithm's performance;
- MAY 03: training process automation;
- MAY 04: integration of the prototype to product.

3.4 Schedule

Schedule is organized to follow the initial requirements; however, as already explained, these changed along the way influencing the schedule as well.

Assigned hours	Activity
8	Functional requirements analysis
8	Study of existing documentation
16	Training on the current image feature extraction system
16	Architecture design
16	UI design
160	Back end development
80	Front end development
16	Documentation

Chapter 4

Analisi dei requisiti

Breve introduzione al capitolo

4.1 Casi d'uso

Per lo studio dei casi di utilizzo del prodotto sono stati creati dei diagrammi. I diagrammi dei casi d'uso (in inglese *Use Case Diagram*) sono diagrammi di tipo [Unified Modeling Language \(UML\)](#) dedicati alla descrizione delle funzioni o servizi offerti da un sistema, così come sono percepiti e utilizzati dagli attori che interagiscono col sistema stesso. Essendo il progetto finalizzato alla creazione di un tool per l'automazione di un processo, le interazioni da parte dell'utilizzatore devono essere ovviamente ridotte allo stretto necessario. Per questo motivo i diagrammi d'uso risultano semplici e in numero ridotto.

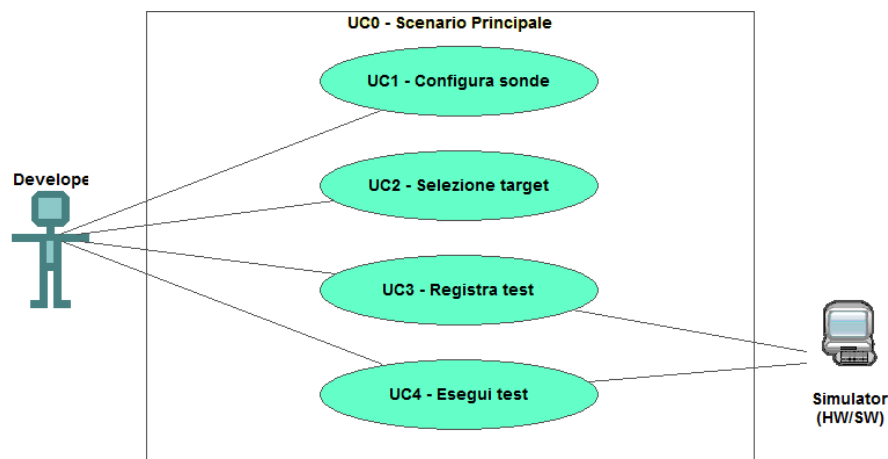


Figure 4.1: Use Case - UC0: Scenario principale

UC0: Scenario principale

Attori Principali: Sviluppatore applicativi.

Precondizioni: Lo sviluppatore è entrato nel plug-in di simulazione all'interno dell'IDE.

Descrizione: La finestra di simulazione mette a disposizione i comandi per configurare, registrare o eseguire un test.

Postcondizioni: Il sistema è pronto per permettere una nuova interazione.

4.2 Tracciamento dei requisiti

Da un'attenta analisi dei requisiti e degli use case effettuata sul progetto è stata stilata la tabella che traccia i requisiti in rapporto agli use case.

Sono stati individuati diversi tipi di requisiti e si è quindi fatto utilizzo di un codice identificativo per distinguerli.

Il codice dei requisiti è così strutturato $R(F/Q/V)(N/D/O)$ dove:

R = requisito

F = funzionale

Q = qualitativo

V = di vincolo

N = obbligatorio (necessario)

D = desiderabile

Z = opzionale

Nelle tabelle 4.1, 4.2 e 4.3 sono riassunti i requisiti e il loro tracciamento con gli use case delineati in fase di analisi.

Table 4.1: Tabella del tracciamento dei requisiti funzionali

Requisito	Descrizione	Use Case
RFN-1	L'interfaccia permette di configurare il tipo di sonde del test	UC1

Table 4.2: Tabella del tracciamento dei requisiti qualitativi

Requisito	Descrizione	Use Case
RQD-1	Le prestazioni del simulatore hardware deve garantire la giusta esecuzione dei test e non la generazione di falsi negativi	-

Table 4.3: Tabella del tracciamento dei requisiti di vincolo

Requisito	Descrizione	Use Case
RVO-1	La libreria per l'esecuzione dei test automatici deve essere riutilizzabile	-

Chapter 5

Design and development

In this chapter I will explain my technology choices and how these influenced the development of the final application.

5.1 Technology and tools

In the following I will introduce the technologies I used and the reasons behind their choice.

Deep learning framework

There are many popular deep learning frameworks out there, but I was asked to choose between *MXNet* and *TensorFlow* after experimenting with both for bit. Keep in mind that I was required to use *YOLO v3* as my *convolutional neural network* and it influenced greatly my final choice. Furthermore, due to Python being predominant in deep learning, I had little freedom in choosing a programming language¹

MXNet and Gluon CV

MXNet is an open source deep learning framework developed by Apache. It supports various languages, but the most extensive toolkit is written in Python; in particular *MXNet* Python provides *Gluon*, a high-level API that allows to easily build and train *deep neural networks*. *Gluon* is also used by *Gluon CV*, a *computer vision* toolkit that provides out-of-the-box implementation of several *neural networks*, pre-trained models and even complete *training scripts*. In particular *YOLO v3* is well supported and many useful tutorials are available.

Both *MXNet* and *Gluon CV* can be installed from [pypi](#) with [pip](#).

There are minor compatibility issues since *Gluon CV* 0.3.0 requires *MXNet* 1.3.0 and later.

¹Actually I did try Scala APIs on *MXNet*, but I had little success with them due to the restricted set of functionality provided when compared to Python.

TensorFlow

TensorFlow is an open source machine learning framework developed by Google Brain. It supports various language, but as common in this field the most extensive one is written in Python. *TensorFlow* is compatible with *Keras*, a high-level API to build and train *deep neural networks*; sadly *Keras* does not provide any network model for *computer vision*.

When looking to *YOLO v3* support I explored *TensorFlow*'s *GitHub* repository, where they host some research projects developed and maintained by independent researchers and not officially supported or distributed in release branches. Among the *computer vision* models on there I did not find *YOLO v3*, but I took my time to experiment with a *SSD* model just to see if *TensorFlow* was worth anyway.

TensorFlow can be installed from *PyPI* with *Pip*, but the research toolkit cannot.

My choice

Since *MXNet*'s *Gluon CV* provides a comprehensive toolkit with models, scripts and utilities to train and infer *YOLO v3 convolutional neural networks*, it easily became my framework of choice.

As I said I tried *TensorFlow* with *SSD*, but it was a low-level implementation that relied a lot on a [jupyter notebook_g](#) and it was not suitable for the goals of my project.

IntelliJ IDEA

I used [intellij idea_g](#) IDE to write my code and I installed *PEP Python* linter plugin to check my code style.

GitLab

I was given access to a [gitlab_g](#) repository to version control my work and I used *IntelliJ IDEA*'s built-in *Git* interface to access it.

Docker

[docker_g](#) is a software application to create containers, virtual environments with their own operative system and tools installed. Containers can communicate with each other and to the outer world through ports, and can mount hard drives when consistent memory is needed.

I used it in various steps of my project's development; at first I created a *Linux Docker* container to set up a *Python 2.7* and *MXNet* environment to run my code. I also created a *MongoDB* instance on another container.

At some point I even dockerized my application to run it on *AWS*.

OpenAPI + Swagger

Thron works with a service architecture, so I was asked to follow it when developing my inference web application. I wrote a [yaml_g openapis_g](#) using [swagger_g](#) to declare my application's interface, which also served as documentation.

Flask + Connexion

My server application for inference is written in *Python*, so I set up a *flask_g* server to host it. The *Flask* needed to adhere to the operations and responses I declared on the *OpenAPI*; matching everything manually would be a tedious job, but thankfully *Zalando*² developed a framework called *Connexion*³ that, provided the *OpenAPI* automatically handles the HTTPs requests and maps them to your Python functions.

Postman

I used *postman_g* to send *HTTP* request to my web server and check its behavior.

MongoDB + Robo 3T

I set up a *mongodb_g* instance in a *Docker* container to model a queue I needed in my inference application, and I used the *Pymongo*⁴ toolkit to work with it from my server application.

To easily access and manage my database, I used *robo 3t_g* client.

AWS SageMaker

I experimented with *AWS sagemaker_g* to test whether a *YOLO v3* training script could easily take advantage of *AWS*'s computational power in a future production environment.

5.2 Design

In this section I will cover the design choices I made while developing my application.

5.2.1 Overview

As stated before, *THRON'S DAM* categorizes images inserted by the clients to make their access easier. Details about how these images are stored in a non-relational database are irrelevant for my project, and I actually do not know them. However, how a database export is structured and what it contains do is important and had major influence in some implementation choices.

A database export is a *json lines_g* where each line represent a valid *Json* file containing information about an image in the database; the covered fields are the following:

- * ID: image unique identifier;
- * Client ID: client identifier⁵; the goal is to construct a custom dataset for detection for each client's needs so this parameter is vital;
- * Url: location where the image is stored, likely an *Amazon S3*⁶ cloud storage location.

²E-commerce platform.

³Zalando Connexion repository: <https://github.com/zalando/connexion>

⁴Pymongo website: <https://api.mongodb.com/python/current/>

⁵Thron has several clients and often uses client ID based sharding techniques to run services in order to guarantee various levels of fairness.

⁶AWS cloud storage service; Amazon S3 website: <https://aws.amazon.com/s3>

- * Data: image properties relevant in the dataset context such as *bounding boxes* coordinates and classes.

5.2.2 Dataset creation

Format choice

To train the network model for my prototype I was given a dataset created by one of my co-workers; his format of choice⁷, however, it was not compatible with my framework so I had to convert it with a script.

Gluon CV provides modules to load a custom dataset from two different formats:

- * Record file format; this format is composed by two different files, a **.rec* binary file, containing the encoded images and their labels, and a **.idx* file containing the indexes to provide random access to the binary;
- * Pure text file format; this format is composed by a **.lst* file⁸. Each line of this file contains the path to an image on disk and the labels describing its bounding boxes.

The record file format allows faster training, but I chose to use the **.lst* file format instead because it allows you to work directly from raw images without prior computation. There are two main reason behind my choice:

1. The production training process should run on *AWS SageMaker*, which reads input data from an *Amazon S3*; Thron already uses *S3* so uploading the training on a dedicated bucket on there would be the easiest approach;
2. In production environment the datasets should be created from *Json lines file* database exports, so simply converting the *Json lines* content describing the bounding boxes for each image to an **.lst* would be the easiest and fastest approach.

Dataset creation approach

The dataset-creation module I designed takes a *Json lines* database export file as its input and creates a **.lst* for each *client ID* found in the *Json lines*. The module also provides a function to split the dataset into *training dataset* and *validation dataset* after creation and saves an additional text file containing statistics about the dataset⁹. I chose to keep the creation of the full dataset and splitting it separated, as it is common practice to split the dataset right before training.

To track the class names with their numerical encoding used in the **.lst* dataset a **.txt* calles *synset.txt* is used. This file must be created manually.

5.2.3 Training

Model training for my *YOLO v3* network can be executed as a stand-alone *Python* script. This file receives as an input the *training* and *validation* datasets in **.lst* format and the *synset.txt* file.

⁷Each raw image in the dataset was labeled by an XML file containing the image size, bounding boxed and classes.

⁸It is a text file readable as a tab-separated **.csv file*.

⁹These statistics include number of images, number of classes, cardinality and density.

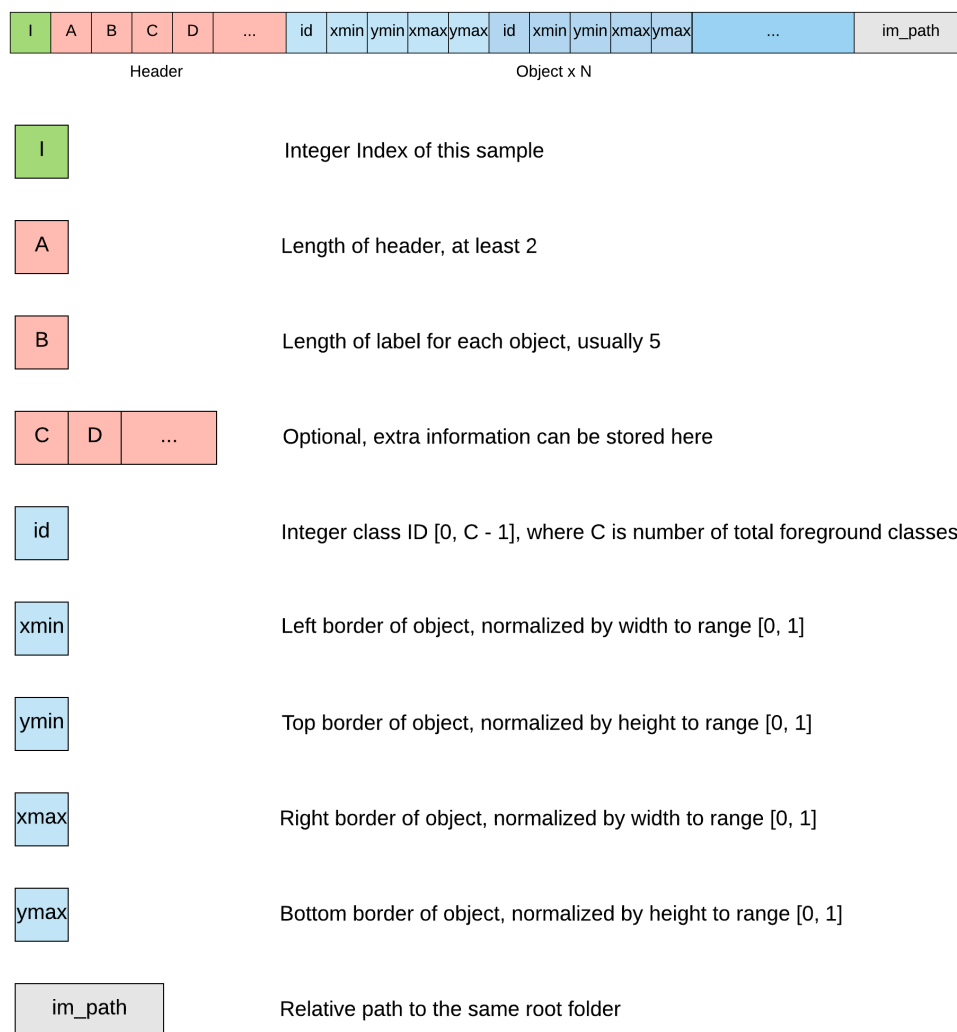


Figure 5.1: Structure of a **.lst* file containing dataset information.

My script actually performs *transfer learning* instead of training from scratch due to the better performances achievable this way; the base network is a YOLO v3 Darknet53¹⁰ pre-trained on the *VOC* dataset; the *VOC* weights are provided by the *Gluon CV* toolkit and should cover classes general enough to prove useful in most cases.

Hyperparameters

My scripts accepts arguments to tune the training *hyperparameters*. Follows the hyperparameters list and their description:

- * Data shape: input data shape; the shortest side of each image will be resized to match this value. Accepted values are 320, 416 and 608 (32 multiples) and it defaults to 416;
- * Batch size: number of images in each computing batch. It must be tuned accordingly to your hardware as a batch size too big would trigger a bus error. It defaults to 8 (which works for CPU), but with GPU training you might be able to use a bigger size such as 32;
- * Workers: number of workers used for training. You must tune it accordingly to your hardware as a number too high might result in a thread error;
- * Gpus: numbers of the GPUs to use as a comma separated list of integers; to train on CPU leave this argument empty;
- * Epochs: number of epochs to perform;
- * Learning rate: speed at which your model learns the dataset content. A high learning rate makes you model learn quicker, but might result in overfitting;
- * Learning rate decay rate: speed at which your learning rate decays;
- * Learning rate decay epoch: epoch interval where your learning rate should decrease;
- * Momentum: stochastic gradient descent momentum;
- * Syncbn: argument to synchronize devices when training in parallel.

Model export

There are two ways to export a trained (or partially trained) network:

- * Saving weights: it is possible to set the network weights as a **.PARAMS* binary file every set epoch; this type of file is mostly used as a training checkpoint to resume training from later.
- * Network export: it is possible to set an export epoch where both the network architecture with adjusted top layers and its weights are exported. The network itself is saved as a **.json* file describing its layers while the weights are saved as a **.PARAMS* similar to the checkpoint file. These files are coupled and must be used together to load the trained network; this export approach is used when you want to save a model to use for inference, as it can be loaded from different languages other than *Python*.

¹⁰YOLO v3 convolutional neural network built on top of a Darknet classifier composed by 53 layers

Resume training

Since training demands a lot of computational power, it is useful to be able to stop a training process in order to resume it later. The following arguments allows you to load weights from a training checkpoint to resume training:

- * Resume: path to the *.PARAMS binary file containing the training checkpoint;
- * Start epoch: epoch where training should be resumed.

It is important to load a checkpoint *.PARAMS binary file to resume training as the model export *.PARAMS binary file cannot be loaded without its paired *.json.

5.2.4 Inference service

While training usually requires GPU power to speed up the process, inference can be reasonably done on CPU, so I decided to keep my training script and my inference script separated.

RESTful Architecture

My inference request server uses *restful api*_g architecture that I designed using the API first approach; specifically, I wrote my *OpenAPIs* using the *Swagger*. To implement my APIs I took advantage of Zalando's *Connexion*, a framework based on *Flask* that tracks the operations declared in the OpenAPI to the application, validates the HTTP request exchange by comparing the actual Json content to what is declared in the API and manages error messages¹¹.

To increase decoupling I kept my server interface separated from its implementation; in fact operations are declared in a *JobOpsMongo* module called by the server interface via a *JobService* module. I followed this decoupling approach for both *MongoDB* related operations and *MXNet* inference related operations, but the latter aren't directly called by the inference request server but are used by a queue consumer I did not implement and that potentially resides on AWS.



Figure 5.2: UML describing the relationship between the server application (App class) and the implementation of its methods.

Basically the server's job is to accept and manage inference requests by inserting them into a *MongoDB* queue; from there a consumer will dequeue them and run the actual inference.

HTTP requests work with *Json* objects which, despite *Connexion*'s automatic conversion into *Python* dictionaries, still require tedious syntax to access their fields. On

¹¹Since I used *Connexion* instead of a pure *Flask* server application I couldn't override the standard HTTP error messages

the other hand *MongoDM* indexes its content in *Json* format. Since working with different formats of data would unnecessarily increase the application's complexity, I found it useful to model *Python* class objects to reflect what was declared in the *OpenAPI* and use them to work within the logic modules. My approach while designing these classes was to roughly follow the structure I declared in the *OpenAPI*, adding further information that could be useful inside the application, and to provide object creation, serialization¹² and deserialization¹³ methods to quickly convert between *Python* objects, *Python* dictionaries and *Jsons* as needed.

Asynchronous computation

While inference isn't too heavy on the computational side of the task, performing a prediction requires the trained network model to be loaded into memory and, since each model occupies over 250MBs of space, it is obvious that constantly switching from model to model would make the performance fall. The easiest solution to this problem is to perform inference by sharding the requests by client ID, so each batch of images runs on the same network model. A further advantage of client ID sharding is that it allows guarantee different levels of fairness accordingly to each client's service contract. Since computation is obviously performed asynchronously, a client needs a way to retrieve its inference once it's completed. There are two main options to manage this: *callbacks_g* and *polling_g*. I chose the latter because it comes more natural when working in a *RESTful* environment, as after having inserted the new job in the queue the server can simply send a HTTP 202 Accepted response containing the address to check for the requested resource.

Queue

I modeled the queue on *MongoDB*, as it provides an easy way to sort the elements in its collections by insertion time. Furthermore, contents in *MongoDB* are represented as *Json* objects, so it proved to be convenient to work with the objects I declared in the *OpenAPI*.

The jobs queue is accessed by both my inference request server and a consumer, but I did not implement the latter.

Inference work flow

When a client needs an image to be inferred, it sends the server a HTTP Post request containing a *Json* stating the URL to the image and a client ID; upon receiving the request the server creates a new inference request job and inserts it in a *MongoDB* queue. If the insertion was successful, the server sends the client a HTTP 202 Accepted response, containing the URL where to check for the requested prediction; since computation runs asynchronously, the client might have to check several times for its requested inference result before it's actually available.

When the prediction is ready, the client can retrieve it from the given URL; inference contains the detected object classes, their confidence score¹⁴ and their *bounding boxes*,

¹²From *Python* object to *Json*

¹³From *Json* to *Python* objects.

¹⁴How much the network is confident the detected object belongs to the stated class; a higher score means higher confidence.

expressed in percentage on the image's size.
Predictions are executed by I consumer I did not implement.

Unit test

I wrote unit tests that cover both the methods provided by the *JobOpsMongo* module and the *PredictionOpsMxnet* modules. While the *JobOpsMongo* module was actually used by the application and I tried its functionality by sending HTTP requests via *Postman* as well, the *PredictionOpsMxnet* module has only been used by its unit tests. In fact, the *PredictionOpsMxnet* module provides the inference methods that the consumer should use, but I did not implement it because it should potentially run on *AWS* and therefore requires specific coding.

5.2.5 Standalone inference script

Before designing and developing the inference request server application, I wrote a *Python* standalone script to perform inference on an image. This script takes the image path and the network model to use as arguments, and then run inference. The result is composed by the classes of the detected objects, their *bounding boxes* and their confidence score; furthermore the inference script saves a copy of the input image where it draws *To edit images I used Python PILLOW library. PILLOW website: <https://pillow.readthedocs.io/en/5.3.x/>* the *bounding boxes* with their class and confidence score.

I used this inference script mainly when testing whether the network training was successful; since I used unlabeled data¹⁵ for testing, visually checking whether the detection was accurate proved to be faster than automatizing the process, as classes and their *bounding boxes* must be manually determined by a human anyway.

¹⁵Using dataset images for testing could give false performance, especially if overfitting has occurred.

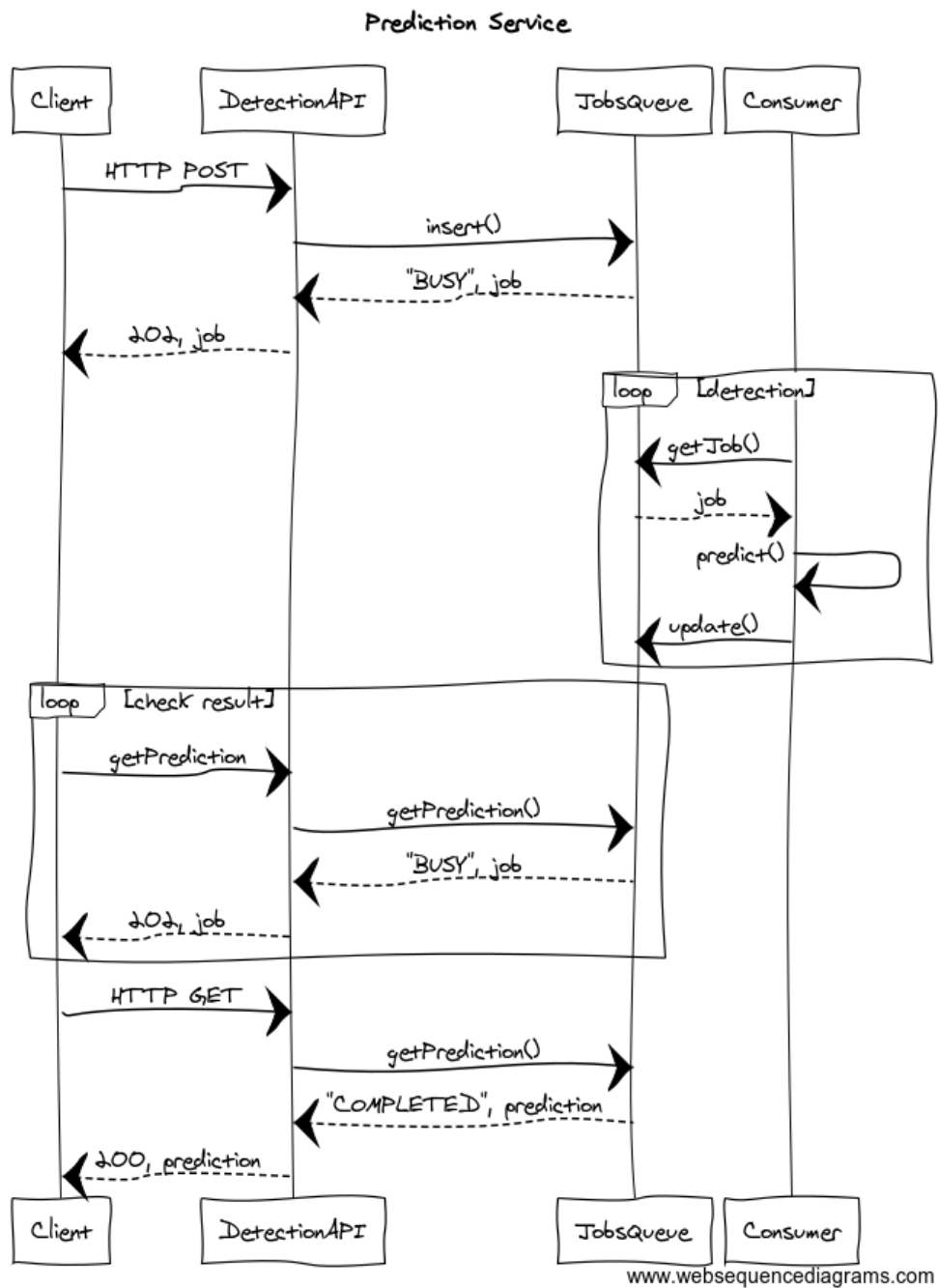


Figure 5.3: UML flow chart displaying the HTTP message exchange in the application.

Chapter 6

Verifica e validazione

Chapter 7

Conclusioni

7.1 Consuntivo finale

7.2 Raggiungimento degli obiettivi

7.3 Conoscenze acquisite

7.4 Valutazione personale

Appendix A

Appendice A

Citazione

Autore della citazione

Bibliografia