



Júlio Miguel Braz da  
Costa Silva

**Recuperação e Identificação de Momentos em  
Imagens/Vídeos**

**Recovery and Identification of Moments in  
Images / Videos**

# **DOCUMENTO PROVISÓRIO**





Júlio Miguel Braz da  
Costa Silva

**Recuperação e Identificação de Momentos em  
Imagens/Vídeos**

**Recovery and Identification of Moments in  
Images / Videos**

Dissertação de Mestrado apresentada à Universidade de Aveiro, para  
obtenção do grau de Mestre em Engenharia Eletrónica e de Telecomuni-  
cações, sob orientação do Professor Doutor António Neves ...

**DOCUMENTO  
PROVISÓRIO**



**o júri / the jury**

presidente / president

**ABC**

Professor Catedrático da Universidade de Aveiro (por delegação da Reitora da Universidade de Aveiro)

vogais / examiners committee

**DEF**

Professor Catedrático da Universidade de Aveiro (orientador)

**GHI**

Professor associado da Universidade J (co-orientador)

**KLM**

Professor Catedrático da Universidade N



## agradecimentos / acknowledgements

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetuer id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum....

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris....



## Palavras-Chave

HEVC, ...

## Resumo

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.



**Keywords**

HEVC, ...

**Abstract**

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.



---

# Contents

---

<b>List of Figures</b>	<b>vi</b>
<b>List of Tables</b>	<b>vii</b>
<b>Acronyms</b>	<b>ix</b>
<b>Glossary</b>	<b>xi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Context . . . . .	1
1.2 Motivation . . . . .	2
1.3 Challenges . . . . .	2
1.4 Objectives . . . . .	2
1.5 Contributions . . . . .	2
1.6 Document Structure . . . . .	2
<b>Nomenclature</b>	<b>1</b>
<b>2 Image/Video Feature Extraction</b>	<b>3</b>
2.1 Fundamental Concepts . . . . .	4
2.1.1 Artificial Intelligence . . . . .	4
2.1.2 Machine Learning . . . . .	4
2.1.3 Deep Learning . . . . .	4
2.1.4 Computer Vision . . . . .	5
2.1.5 Image Annotation and Classification . . . . .	5
2.1.6 Object Detection, Segmentation and Recognition . . . . .	6
2.1.7 Features and Feature Space . . . . .	6
2.1.8 Object . . . . .	6
2.1.9 Image Description . . . . .	7
2.1.10 Datasets With Common Objects . . . . .	7
2.2 Computer Vision Libraries . . . . .	8
2.2.1 OpenCV . . . . .	8
2.2.2 VLFeat . . . . .	8
2.2.3 BOOFCV . . . . .	8
2.2.4 GluonCV . . . . .	9
2.3 Neural Networks . . . . .	9
2.3.1 Neural Network Training . . . . .	10
2.3.2 Types of Neural Networks architectures . . . . .	10
2.4 CNNs architectures For Image Classification . . . . .	12
2.4.1 SqueezeNet . . . . .	12

2.4.2	ResNet . . . . .	14
2.4.3	InceptionV3 . . . . .	15
2.4.4	DenseNet . . . . .	17
2.5	Regression based algorithms for Object Detection . . . . .	17
2.5.1	RetinaNet . . . . .	17
2.5.2	YOLOv3 . . . . .	18
2.5.3	TinyYoloV3 . . . . .	19
2.5.4	Single Shot MultiBox Detector (SSD) . . . . .	19
2.6	Classification Based Algorithms For Object Detection . . . . .	20
2.6.1	R-CNN Models Summary . . . . .	21
2.6.2	R-CNN . . . . .	21
2.6.3	Fast R-CNN . . . . .	21
2.6.4	Faster R-CNN . . . . .	22
2.6.5	Mask R-CNN . . . . .	22
2.7	State-Of-The-Art . . . . .	23
2.7.1	COCO Test-Dev . . . . .	24
2.7.2	ImageNet . . . . .	26
<b>3</b>	<b>Information Extraction From Text</b> . . . . .	<b>29</b>
3.1	Natural Language Processing . . . . .	30
3.1.1	Important NLP Terminologies . . . . .	30
3.1.2	Core Areas . . . . .	31
3.1.3	Application Areas . . . . .	31
3.2	Numerical Representation of Text . . . . .	31
3.2.1	Word Embeddings . . . . .	31
3.3	Static Word Embedding Models . . . . .	32
3.3.1	Word2Vec . . . . .	32
3.3.2	GloVe . . . . .	33
3.3.3	FastText . . . . .	34
3.4	Contextualized Word Embedding Models . . . . .	34
3.4.1	Context2vec . . . . .	34
3.4.2	ELMo . . . . .	35
3.5	Available NLP libraries . . . . .	36
3.5.1	SpaCy . . . . .	36
3.5.2	Natural Language ToolKit . . . . .	36
3.5.3	Stanford Core NLP . . . . .	37
3.5.4	Gensim . . . . .	37
3.5.5	Others . . . . .	37
<b>4</b>	<b>Initial Work</b> . . . . .	<b>39</b>
4.1	Image Recognition test runs . . . . .	39
4.1.1	Test Run Number 1 . . . . .	40
4.1.2	Test Run Number 2 . . . . .	41
4.1.3	Test Run Number 3 . . . . .	42
4.1.4	Results analysis . . . . .	43
4.2	Object Recognition test runs . . . . .	44
4.2.1	Test Run Number 1 . . . . .	44

4.2.2	Test Run Number 2 . . . . .	46
4.2.3	Test Run Number 3 . . . . .	48
4.2.4	Results analysis . . . . .	50
<b>5</b>	<b>ImageCLEF</b>	<b>51</b>
5.1	imageCLEF challenge . . . . .	51
5.2	Tasks . . . . .	51
5.2.1	LMRT sub-task . . . . .	51



---

## List of Figures

---

2.1	Feature extraction from an image. [5]	3
2.2	Generic picture of a family having a picnic.	7
2.3	Typical neural network architecture. [7]	9
2.4	Operations done by a neuron.	10
2.5	Example of a Feedforward Neural Network with one hidden layer (with 5 neurons) [29].	11
2.6	CNN architecture	12
2.7	SqueezeNet fire module. [31]	13
2.8	SqueezeNet architecture. [32]	14
2.9	Skipping connection example.[33]	14
2.10	ResNet Architecture.[34]	15
2.11	Inception Module. [35]	15
2.12	Mini-network replacing the $5 \times 5$ convolutions (Example of factorization). [36]	15
2.13	InceptionV3 architecture.[37]	16
2.14	A 5-layer dense block. Each layer takes all preceding feature-maps as input. [36]	17
2.15	RetinaNet architecture.[39]	18
2.16	Bounding Box Prediction : Predicted Box (Blue), Prior Box (Black Dotted).[41]	18
2.17	The network architecture of YOLO base model.[42]	19
2.18	SSD architecture. [45]	20
2.19	R-CNN model family summary. [46]	21
2.20	R-CNN architecture. [46]	21
2.21	Fast R-CNN architecture. [46]	22
2.22	Faster R-CNN architecture. [46]	22
2.23	Mask R-CNN is a Faster R-CNN model with image segmentation. [46]	22
2.24	Object Detection on COCO test-dev benchmark .[51]	23
2.25	Image Classification on ImageNet benchmark. [52]	23
2.26	CBNet Architecture for object detection.	25
2.27	ResNeXt architecture. [34]	25
2.28	Noisy Student Method. [66]	27
2.29	Comparison of different scaling methods: (a) is a baseline network example; (b)-(d) are conventional scaling that only increases one dimension of network width, depth, or resolution. (e) is the proposed compound scaling method that uniformly scales all three dimensions with a fixed ratio. [70]	27
2.30	EfficientNet-B0 architecture representation.	28
3.1	Natural Language Processing [73].	29
3.2	Classification of NLP. [74]	30
3.3	Example of text representation by one-hot vector.	32
3.4	One-hot encoding resulting vector. [81]	33

3.5	Word2Vec encoding resulting vector. [81] . . . . .	33
3.6	CBOW model and Skip-Gram model. [79] . . . . .	33
3.7	A 2D illustration of context2vec's embedded space and similarity metrics. Triangles and circles denote sentential context embeddings and target word embeddings, respectively [85]. . . . .	35
3.8	Closest target words to various sentential contexts, illustrating context2vec's sensitivity to long range dependencies, and both sides of the target word [85]. . . . .	35
3.9	Nearest neighbors to "play" using GLoVe and context embeddings from a biLM [86]. . . . .	36
4.1	First picture to be analysed. . . . .	40
4.2	Results obtained. . . . .	40
4.3	Second picture to be analysed. . . . .	41
4.4	Results obtained. . . . .	41
4.5	Third picture to be analysed. . . . .	42
4.6	Results obtained. . . . .	42
4.7	First picture to be analysed. . . . .	44
4.8	RetinaNet Detections. . . . .	44
4.9	RetinaNet Detections. . . . .	44
4.10	YOLO Detections. . . . .	45
4.11	YOLO Detections. . . . .	45
4.12	TinyYolo Detections. . . . .	45
4.13	TinyYolo Detections. . . . .	45
4.14	Second picture to be analysed. . . . .	46
4.15	RetinaNet Detections. . . . .	46
4.16	RetinaNet Detections. . . . .	46
4.17	YOLO Detections. . . . .	47
4.18	YOLO Detections. . . . .	47
4.19	TinyYolo Detections. . . . .	47
4.20	TinyYolo Detections. . . . .	47
4.21	Third picture to be analysed. . . . .	48
4.22	RetinaNet Detections. . . . .	48
4.23	RetinaNet Detections. . . . .	48
4.24	YOLO Detections. . . . .	49
4.25	YOLO Detections. . . . .	49
4.26	TinyYolo Detections. . . . .	49
4.27	TinyYolo Detections. . . . .	49

---

## List of Tables

---

2.1	COCO Test-Dev Benchmarks. . . . .	24
2.2	ImageNet Benchmarks. . . . .	26



---

## Acronyms

---

**AI** Artificial Intelligence.

**ANN** Artificial Neural Network.

**CV** Computer Vision.

**NLP** Natural Language Processing.

**ResNet** Residual Network.

**SSD** Single Shot Detection.

**YOLO** You Only Look Once.



---

## Glossary

---

**Pixel** Picture Element.



# CHAPTER 1

---

## Introduction

---

This chapter gives an introduction to the surrounding theme addressed in this thesis. In that sense, firstly the contextualization of the theme and the respective motivation will be presented. The different challenges, the objectives that are intended to be reach and the contributions given to the community are also described. Finally the document structure and organization is explained.

### 1.1 | Context

The pervasive creation and consumption of content, especially visual content, is ingrained into our modern world. In the past, the main purpose given to pictures was to save moments of past events. People nowadays are constantly consuming visual media content. Pictures, images and photos have many different usages, they can either be used to save a moment or to document processes; we use them in engineering, in art, in science, in medicine, in entertainment and also in advertising. [1]

With the rapid development of Internet of things (IOT) this growth in consumption of visual media content has increased the usage of wearable and smart technologies making the subject of lifelogging more prevalent in the recent years. Lifelogging is the task of tracking and recording personal data created through the activities and behaviour of individuals during their day-to-day life in the form of images, video, biometric data, location and other data. The name given to this dataset is "lifelog data" and is rich in resources for contextual information retrieval. [2]

Some great examples of the usefulness of lifelogging is using it as memory extension for people who suffer from memory impairments such as Alzheimers, to find lost items during the day or even to understand human behaviour.

Most of the technical problems associated with creating, compressing, storing, transmitting, rendering and protecting image data are already solved. However we still face two main challenges which are the issues associated with image location and the continuous growth of image data (big data).

Locating images involves analysing them to determine their content, classifying them into related groupings, and searching for images. In order to solve these problems, the current technology relies heavily on the image description, usually called as "image metadata". This data can either be captured automatically at creation time or manually added afterwards.

In the present time the development in the area of content-based analysis (indexing and searching of visual media) is increasing, this is where most of the research in image management is concentrated. Automatic analysis of the content of images, which in turn would open the door to content-based indexing, classification and retrieval, is an inherently difficult problem

and therefore progress is slow. [1]

- 1.2 | Motivation
- 1.3 | Challenges
- 1.4 | Objectives
- 1.5 | Contributions
- 1.6 | Document Structure

# CHAPTER 2

## Image/Video Feature Extraction

The task of automatically recognizing and locating objects in images and videos is of extreme importance for computers to be able to understand and interact with their surroundings. Some major applications of this particular task are pedestrian face detection, surveillance, autonomous driving and text digitalization, where object detection is a crucial challenge. [3]

Feature extraction plays an important role in image classification and object detection systems which are two core components of computer vision. It is characterized by two important aspects, the mapping of image pixels into the feature space (explained in more details in section 2.1.7) and with the extraction of various attributes of an object. Only after extracting useful features from either images or videos, the computer is able to define what an object is or what a certain environment contains.[4]

Putting it simply, feature extraction is the first step to convert an image into text.

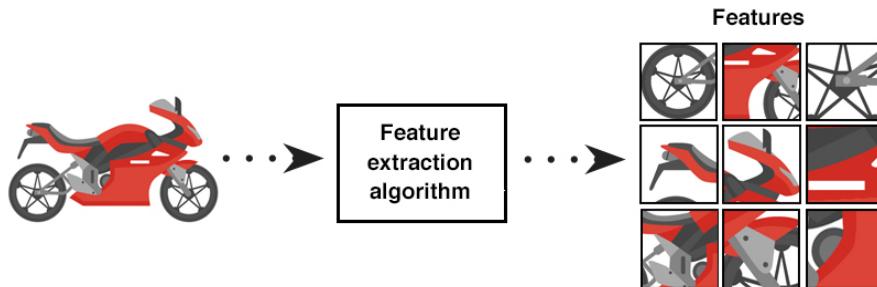


Figure 2.1: Feature extraction from an image. [5]

This chapter starts with fundamental concepts in order to understand how feature extraction in computer vision works. Section 3.5 gives a brief introduction of the most common computer vision libraries. Neural networks are introduced in section 2.3, where the most common architecture types are presented. All CNN architectures and regression based algorithms used during the development of this thesis are explained with detail in sections 2.4 and 2.5. Classification based algorithms, presented in 2.6, were not used in this work, but they are a required reading in order to understand the differences between them and the regression based algorithms. Finally, a state-of-the-art in object detection and image classification can be read in section 2.7.

## **2.1 | Fundamental Concepts**

### **2.1.1 | Artificial Intelligence**

Artificial Intelligence (AI) is the artificial simulation of human intelligence by a computer system in a way that it can perceive its environment, understand its behaviors and take action. Two important areas of AI are machine learning and deep learning. [6]

### **2.1.2 | Machine Learning**

Machine learning can be defined as a data analytics technique that allows computers to learn from experience. There are two types of machine learning techniques, which are supervised learning and unsupervised learning.

Normally, supervised machine learning is used to train a model to predict future outputs, this is done by inputting and outputting known data. Supervised learning uses two different techniques which are classification and regression. Classification techniques are used to classify input data into categories while regression techniques are used to predict continuous responses.

Unsupervised learning is mostly used to find hidden patterns or intrinsic structures in input data. The most common unsupervised learning technique is clustering which is used for data analysis exploration, in order to find hidden patterns or groupings in data. [7]

### **2.1.3 | Deep Learning**

Deep Learning is a subset of Machine Learning that is inspired by the structure and function of the human brain. In order to achieve this, deep learning resorts to artificial neural networks (ANNs).

The idea behind an ANN is that it tries to replicate the working of the human brain in the processing of data and creation of patterns, which is important for decision making. These ANNs are capable of learning unsupervised data that can either be unstructured, unlabeled or both.

Putting it as simple as possible, deep learning is a machine learning technique that teaches computers to learn by example, like a human would. [8]

Thanks to the new digital era, there has been an exponential increase in all forms of data, from every region of the planet. This data is defined as "big data" and comes from sources like social media, search engines, live streaming services and many others. Even though all of this information is easily accessible, it is unstructured. The problem with unstructured data is that the human brain cannot comprehend it efficiently enough to extract relevant information. However, using deep learning, all of this unstructured data can be usable.

A computer model learns how to perform classification tasks directly from data, being it text, images or sound. Current deep learning models are able to achieve such levels of accuracy that they can outperform humans.

In deep learning, models are trained with the usage of a large set of labeled data and neural network architectures that contain many layers. This is one of the disadvantages of deep learning, in order to improve the results of an ANN it requires to be trained with large amounts of labeled data.

Since deep learning deals with such great volumes of information, this introduces another disadvantage to deep learning, which is the extreme need of higher and higher computing power.

Some use cases for deep learning being used currently in the real world are: [8]

- Automated Driving : For the detection of pedestrians.
- Aerospace and Defense : To identify objects from satellites and identify safe or unsafe zones for troops.
- Medical Research : For the automatic detection of cancer cells.

However, the main purpose of deep learning, for this work, will be to apply it to the images obtained from lifelogging. In simple words, lifelogging is the process of tracking and record personal data created through our activities and behaviour [2]. More on lifelogging can be read in chapter 5.

#### **2.1.3.1 | How Deep Learning Works**

The term "deep" comes from the usage of an extensive quantity of hidden layers in the neural network. A normal neural network usually contains 2-3 hidden layers whereas a deep neural network can go up to 150 hidden layers or more.

As explained previously, deep learning models are trained by the usage of large sets of labeled data and neural network architectures that learn features directly from the data, without the need for manual feature extraction. This automated feature extraction makes deep learning models highly accurate for computer vision tasks such as object classification.

Deep Learning also offers "end-to-end learning", this means that a network can learn how to automatically classify raw data.

In addition, deep learning algorithms scale with data, whereas machine learning methods bottleneck at a certain level of performance when more examples and training data are added, which gives deep learning networks a key advantage since they improve as the size of the data increases. [8]

#### **2.1.4 | Computer Vision**

Computer vision is a field of artificial intelligence and computer science that aims at giving computers a visual understanding of the world [9] [10]. It is related with pattern recognition, which is a common way to train a computer so that it can understand visual information. Pattern Recognition is the training a computer goes through when it is fed with different labeled images and then subjected to different algorithms, allowing the computer to hunt for patterns in every element related to those specific labels [11].

#### **2.1.5 | Image Annotation and Classification**

Image classification is the process of associating an entire image with just one label. A simple example of image classification is labeling types of animals, cars or plants. [12]

Image annotation, one of the most important tasks in computer vision, is the process of manually annotating an image with labels. These labels are predetermined in order to give the computer vision model information about what is shown in the image, they are a combination of a bounding box in specific coordinates of the image and a description of the object inside of it. [13]

Feeding this kind of annotated image data to a computer model teaches it to recognize the visual characteristics of that specific label, this makes the model able to categorize new unannotated images of the same type of that label.

### **2.1.6 | Object Detection, Segmentation and Recognition**

Object detection is the name given to the process that combines image classification with object localization [14]. As previously explained, image classification is the prediction and assignment of a class label to an image, while object localization is the prediction and drawing of a bounding box around one or more objects in the image. In other words, object detection is the task that deals with the detection of objects of a certain class (e.g "flower", "table", "plane") in images, making it a natural extension of the classification problem.

The object detection task is considered to be a supervised learning problem, since the objective is to design an algorithm which can accurately locate and correctly classify as many instances of objects as possible, in a bounding box, while avoiding false detections in a given set of training images.

As an added challenge, many object detection applications require the problem to be solved in real time, which can be achieved. However, in order for a detector to be faster accuracy is usually reduced.

Finally, object segmentation is the task of grouping pixels from the same object into a single region and object recognition is the recognition of an object contained in a bounding box. [3]

### **2.1.7 | Features and Feature Space**

A feature is considered to be a measurable piece of data in the image which is unique to a specific object, it can be color, texture or shape. Usually these features are extracted from the image and used in order to represent an object. Color is the most straightforward visual feature for indexing and image retrieval, while shape representation is the most difficult. This is because a 3-D real world object is represented in a 2-D plane in an image, which means that one dimension of information is completely lost. Texture features are very important in pattern recognition and is an important cue in region based segmentation of images.

The similarity between images can be determined through features which are represented as a vector.

To sum things up, feature space is a collection of features related to some properties of the object, while a feature is an individual measurable characteristics of the object. [4]

### **2.1.8 | Object**

An object is used to identify specific items in an image or specific frames in a video. It is possible to label multiple objects in an image. An example of objects in an image of a car might be wheels, headlights, etc.

Usually an object is represented by a group of features in form of a feature vector that is used to recognize objects and classify them. [4]

In object detection, small objects are normally the ones that give worst results and lower performance when being detected. This happens because the information available to detect them is more compressed and hard to decode without some prior knowledge or context. [3]

### 2.1.9 | Image Description

Image description is the meaning of an image and humans can understand it with relative ease. However computers only see the digital representation of images, only detecting pixels, and therefore they are not able to recognize the semantic of the image. This problem makes the semantic gap the main challenge in computer vision [15]. This gap is defined by the lack of coincidence between the information extracted from visual data and the interpretation in a given situation. [3]

As an example, picture 2.2 shows an image of a family having a picnic. Feeding this image to a computer will output very different results from what a human would say.



Figure 2.2: Generic picture of a family having a picnic.

- Computer output: Tree, bottle, person, apple, cup.
- Human output: A family having a picnic in the park.

A computer is only able to output the objects detected but it is incapable of giving them any sort of meaning.

### 2.1.10 | Datasets With Common Objects

A dataset is a collection of images and videos that contain every day life objects that are manually labeled. State-of-the-art object detection models require deep learning neural networks, and in order for neural networks to be trained, they require training datasets, as previously explained.

A few examples of some available datasets are: MS COCO [16], ImageNet [17] , VisualGenome [18], OpenImages [19] and Pascal-VOC [20]

Some of these datasets propose challenges, where teams are able to compete in order to achieve state-of-the-art results. This subject is discussed in section 2.7

## 2.2 | Computer Vision Libraries

### 2.2.1 | OpenCV

OpenCV is an open source computer vision and machine learning software library originally developed by Intel in the year 2000 [21].

The library has more than 2500 optimized algorithms, which includes a comprehensive set of both classic and state-of-the-art computer vision and machine learning algorithms. These algorithms can be used to detect and recognize faces, identify objects, classify human actions in videos, track camera movements, track moving objects, extract 3D models of objects, produce 3D point clouds from stereo cameras, stitch images together to produce a high resolution image of an entire scene, find similar images from an image database, remove red eyes from images taken using flash, follow eye movements, recognize scenery and establish markers to overlay it with augmented reality, etc. [22]

OpenCV Supports the deep learning frameworks like Tensorflow, Torch/PyTorch, Caffe and it is the most standardized tooling for computer vision.

#### 2.2.1.1 | Tensorflow

Tensorflow is currently the most popular open source framework for numerical computation and large-scale machine learning introduced by google and was originally created for tasks with heavy numerical computations. [23] [24]

Tensorflow is written in c++ which enables extremely fast compile times, non the less, it can still be accessed by other languages, such as Python and also supports CPUs, GPUs and distributed processing.

The name given to tensorflow comes from the inputs, since it receives inputs as a multi-dimensional array, also known as tensors. The input (tensor) goes on one end and then it “flows” throughout a system of operations and comes out on the other end as output.

Tensorboard is a feature of tensorflow that allows the monitoring of what tensorflow is doing graphical and visually.

### 2.2.2 | VLFeat

The VLFeat open source library implements popular computer vision algorithms specializing in image understanding and local features extraction and matching. Algorithms include Fisher Vector, VLAD, SIFT, MSER, k-means, hierarchical k-means, agglomerative information bottleneck, SLIC superpixels, quick shift superpixels, large scale SVM training, and many others. It is written in C for efficiency and compatibility, with interfaces in MATLAB for ease of use, and detailed documentation throughout. It supports Windows, Mac OS X, and Linux. [25]

### 2.2.3 | BOOFCV

BoofCV is an open source library written from scratch for real-time computer vision. Its functionality covers a range of subjects, low-level image processing, camera calibration, feature detection/tracking, structure-from-motion, fiducial detection, and recognition.

BoofCV is organized into several packages: image processing, features, geometric vision, calibration, recognition, visualize, and IO. Image processing contains commonly used image

processing functions which operate directly on pixels. Features contains feature extraction algorithms for use in higher level operations.

Calibration has routines for determining the camera's intrinsic and extrinsic parameters. Recognition is for recognition and tracking complex visual objects. Geometric vision is composed of routines for processing extracted image features using 2D and 3D geometry. Visualize has routines for rendering and displaying extracted features. IO has input and output routines for different data structures [26].

#### 2.2.4 | GluonCV

GluonCV provides implementations of state-of-the-art (SOTA) deep learning algorithms in computer vision. It aims to help engineers, researchers, and students quickly prototype products, validate new ideas and learn computer vision. [27]

### 2.3 | Neural Networks

A neural network can be considered a computer program that operates identically to how a human brain would, in the sense that it is able to be teachable to do certain tasks like problem-solving. The appeal of a neural network is the ability to emulate the human brain in pattern recognition skills.

Neural networks are composed of many small cells called neurons. These neurons are grouped into several layers that form columns. The connection between columns are formed also through their neurons. Each neuron of each layer is connected to another neuron of another layer. A visual representation of a generic neural network architecture is shown in figure 2.3.

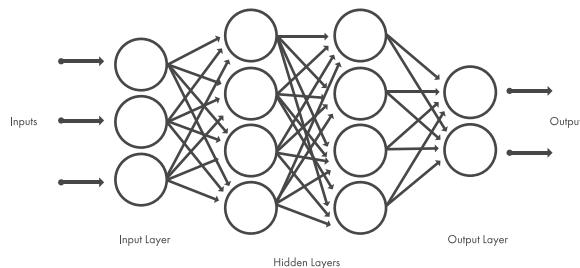


Figure 2.3: Typical neural network architecture. [7]

The connections between layers are called weighted connections and they are adjusted with a real-valued number attached to them. This number is important, because each neuron takes the value of the attached neuron (in their layer) and multiplies it by their connection weight. The bias value is an additional parameter in the neural network which is used to adjust the output along with weighted sum of the inputs to the neuron. The sum of the bias value with the weights is put through an activation function which mathematically transforms the value and assigns it to the connected neuron in the adjacent layer. This is propagated through the whole network. See figure 2.4 for a clear representation of this process.

To put it simply, a neural network can be compared to a filter that goes through all of the possibilities, so that the computer is able to come up with the correct answer.

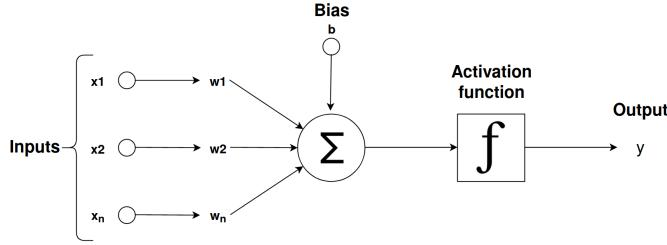


Figure 2.4: Operations done by a neuron.

Sometimes, an object might be too similar to another object which can make the network output a wrong answer. The solution to this problem is the usage of a back-propagation algorithm. This algorithm allows the network to adjust the connections back through the network, check if all the bias values are correct and all of the connections are weighted properly. [28]

### 2.3.1 | Neural Network Training

The best way to train a neural network from scratch is to design a network architecture that will learn through the feeding of a large dataset of labeled data. This allows it to learn the features and model. The problem with this is that depending on the learning rate of the network and the amount of data, these networks can take a lot of time to train (days, maybe weeks).

To solve the problem of time, deep learning applications can recur to the usage of transfer learning. Transfer learning is a process that involves the fine-tuning of a pretrained model. This works by using an existing network like GoogLeNet, and feed it new data of previously unknown classes to the network. After some tweaks to the network, it will be able to categorize only a specific object instead of many different ones. This not only allows the network to be more precise in categorizing that one specific object, but it will also save lot of computation time. [8]

### 2.3.2 | Types of Neural Networks architectures

#### 2.3.2.1 | Feedforward Neural Network

A Feedforward neural network has the most simple architecture, the data only travels in one single direction. It goes through the input node and exits at the output node. Since there is no back-propagation algorithm this neural network is not able to correct itself. [28]

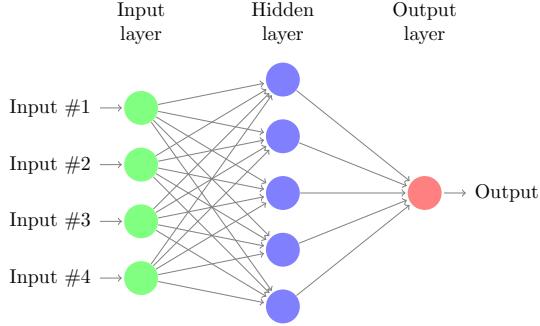


Figure 2.5: Example of a Feedforward Neural Network with one hidden layer (with 5 neurons) [29].

### 2.3.2.2 | Radial Basis Function Neural Network

This network is composed of two layers. In the first one features are combined with a radial basis function in the inner layer. The second one is the output, where these features are taken in consideration while computing the same output in the next function.

A radial basis function means that the distance of a point is considered with respect to the center. [28]

### 2.3.2.3 | Recurrent Neural Network (RNN)

Recurrent Neural Networks (RNNs) are designed to recognize sequential data characteristics and use patterns to predict the next likely scenario. In these kind of neural networks the signals are propagated in both directions as well as within the layers. They work on the principle of saving the output of a layer and feeding it to the input to help in the prediction of the outcome of the layer.

RNNs use the back-propagation algorithm which allows to make sure that the output is correct almost 100% of the time. [28]

### 2.3.2.4 | Convolutional Neural Network (CNN)

CNNs, also known as ConvNets, are a class of deep neural networks that employ the mathematically convolutional operation in at least one of its layers and have a deep feed-forward (not recurrent) architecture [2]. They share similarities with feedforward neural networks, since neurons also have weights and biases that are able to learn. In this network the input features are taken like a filter, which allows the network to have memory, since it can remember the images in parts and compute operations like conversion of the image from RGB or HSI to grayscale, allowing the detection of edges and images that can be classified into different categories. [28]

The only notable difference between CNNs and traditional ANNs is that CNNs are primarily used in the field of pattern recognition within images. This allows the encoding image-specific features into the architecture, making the network more suited for image-focused tasks, while further reducing the parameters required to set up the model. [30]

CNN convolves learned features with input data, and uses 2D convolutional layers which make this architecture one of the best to process 2D data, such as images. They also remove

the necessity of manual feature extraction. There is no need to identify features used to classify images since CNNs work by extracting them directly from images. This is important because relevant features are not pretrained, they are learned while the network trains on a dataset. [8]

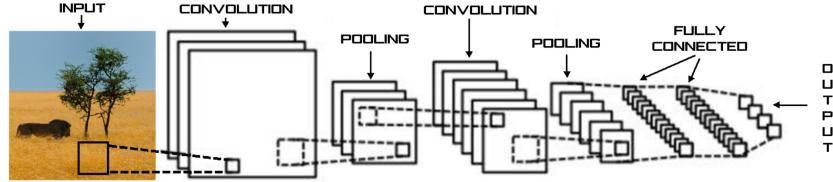


Figure 2.6: CNN architecture  
[2]

### Input Layers

The input layer is the first layer of a CNN and serves the purpose of resizing an image in order for it to pass onto further layers for feature extraction [2]. It also holds the pixel value of the images [30].

### Convolutional Layers

The convolutional layers extracts low-level features from an image, such as edges, color or gradient orientation, according to the applied filter or kernel. [2]

### Activation Functions

The activation function introduces nonlinearity in order for CNNs to learn functionalities. They serve as decision functions and help in learning complex patterns. Some examples of activation functions are sigmoid, softmax and ReLU. [2]

### Pooling Layers

The pooling layers serve the purpose of reducing the parameters required and computation in the network by controlling the overfitting. This is achieved by reducing the spatial size of the network. [2]

Overfitting happens when a model learns the detail and noise in the training data to an extent that it negatively impacts the performance of the model on new data. [30]

### Fully Connected Layers

The final layer in a CNN is usually a fully connect layer used for classification purposes. They take all features from the previous layer and compute class probabilities or scores. These features are then translated into a different class. [2]

## 2.4 | CNNs architectures For Image Classification

### 2.4.1 | SqueezeNet

SqueezeNet is a deep neural network for computer vision that is more efficient for distributed training, since it requires less parameters to be transferred. The main goal of SqueezeNet creation was to obtain a smaller neural network with fewer parameters that could

more easily fit into a computer memory, making it more easily transmitted over a computer network. This neural network was firstly implemented on top of the caffe deep learning software framework and later ported to the chainer deep learning software framework and Apache MXNET framework.

The basis of SqueezeNet consists on 3 ideas [31]:

- Replacing 3x3 filters with 1x1 filters and reduce the number of input channels. This improves computation speed and alleviates the computer resources required, since 1x1 filters have 9 times less parameters than 3x3 filters.
- Utilize 1x1 filters as a bottleneck layer to help reducing the computation required for the following 3x3 filters.
- Keeping a big feature map by down sampling late.

This neural network is built with fire modules, which are represented in figure 2.7.

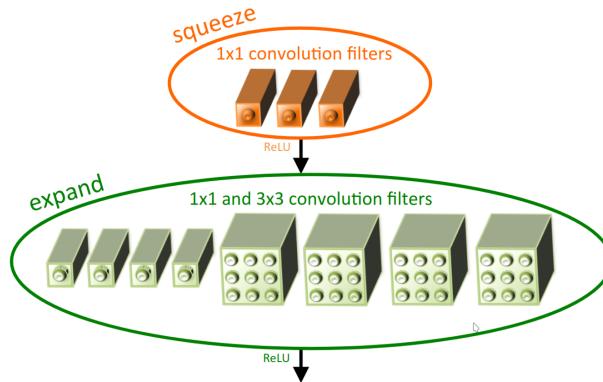


Figure 2.7: SqueezeNet fire module. [31]

The fire module contains both a squeeze layer and an expand layer. SqueezeNet stacks fire modules and pooling layers (this can be seen in figure 2.8). The squeeze layer and expand layer maintain the same feature map size, while the pooling layers reduce the depth to a smaller number, later increasing it. Reducing the depth means the expand layer has fewer computations to do, boosting the speed.

**Squeeze layer architecture:** Consists on 1x1 convolutions, it essentially combines all the channels of the input data into one (and thus reducing the number of input channels needed in the next layer).

**Expand layer architecture:** Consists on 1x1 convolutions mixed alongside 3x3 convolutions. The 1x1 convolutions combine the channels of the previous layers in various ways. The 3x3 convolutions detect structures in the image since 1x1 convolutions can't.

**SqueezeNet architecture:** SqueezeNet doesn't fully connect layers and it consists of 8 fire modules and a single convolution's layer as input and output. It uses Global Average Pooling, taking each channel from the previous convolution layer and builds an average over all values.

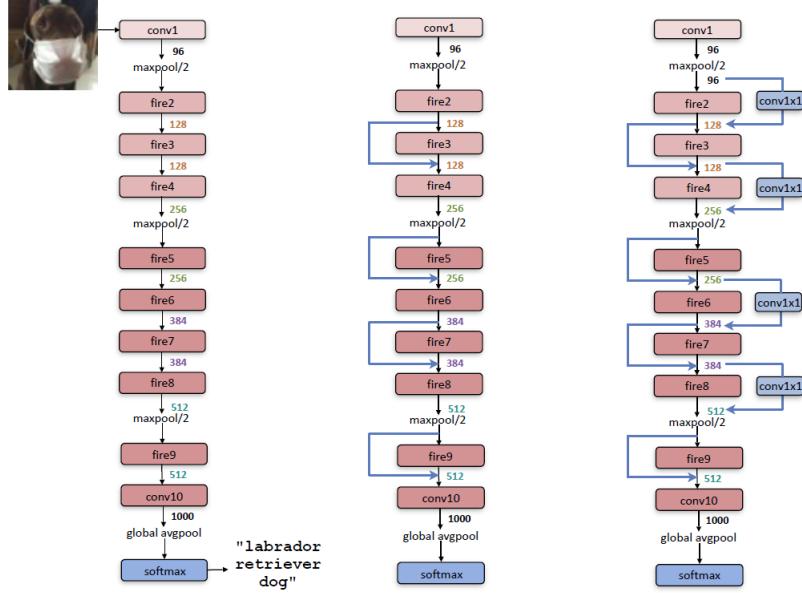


Figure 2.8: SqueezeNet architecture. [32]

#### 2.4.2 | ResNet

The core idea of ResNet (Residual Neural Network) is introducing skip connections (also called identity shortcut connection, represented in figure 2.9). The way this works is by adding the output of an earlier layer to a later layer in order to jump over some layers.

The vanishing of gradients problem makes deep neural networks hard to train, this happens because as the gradient is propagated back to earlier layers, repeated multiplications may turn the gradient too small, this results in a rapidly performance degradation.

Skipping over layers helps avoiding the vanishing of gradients problem and improves the accuracy of the neural network.

Having the skip connection allows the training of extremely deep neural networks, more than 150 layers, successfully and still being able to achieve a compelling performance. [33]

This architecture is represented in figure 2.10.

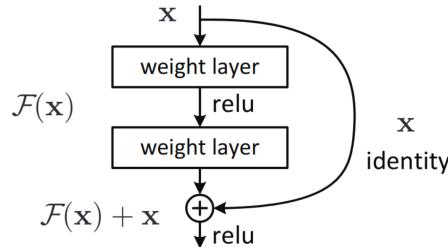


Figure 2.9: Skipping connection example.[33]

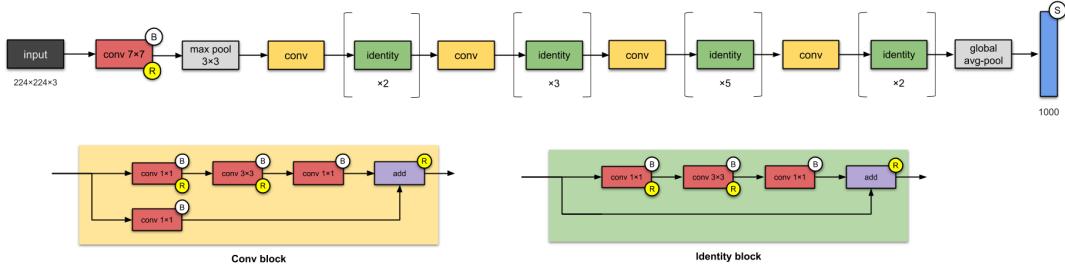


Figure 2.10: ResNet Architecture.[34]

### 2.4.3 | InceptionV3

Initially named GoogLeNet, the Inception-v1 architecture was proposed by researchers of Google company and was the winner of the ILSVRC 2014 competition, making it historically significant in Convolutional Neural Networks.

This network, trained on the imageNet dataset, introduced inception modules (shown in figure 2.11) that allowed for a more efficient computation and deeper network.

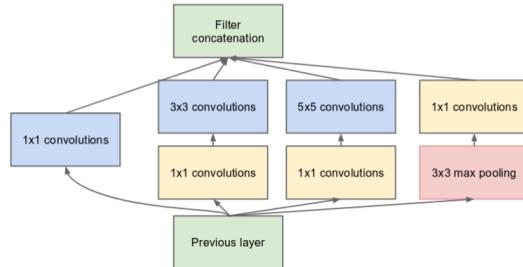


Figure 2.11: Inception Module. [35]

The Inception architecture (Inception-v1) was improved by the introduction of batch normalization (Inception-v2). [2]

InceptionV3, is 48 layers deep and able to classify images into 1000 different categories. The improvement over its predecessors is the adding of factorization ideas (figure 2.12 shows an example of this).

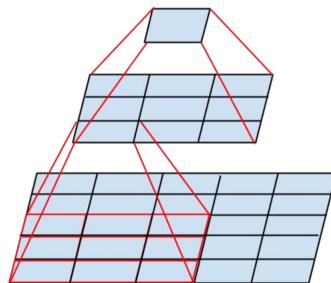


Figure 2.12: Mini-network replacing the  $5 \times 5$  convolutions (Example of factorization). [36]

This third interaction aims at factorizing convolutions, reducing the number of connections/parameters required while maintaining network efficiency. As an example, using a layer of  $5 \times 5$  filter requires  $5 \times 5 = 25$  parameters, this layer can be replaced by two  $3 \times 3$  layers which reduce the number of parameters required by 28%, since  $2 \times (3 \times 3) = 18$  parameters. Reducing the number of parameters required reduces the computational resources required and also prevents overfitting. This enables the network to go deeper. [37]

The inceptionv3 architecture can be seen in figure below.

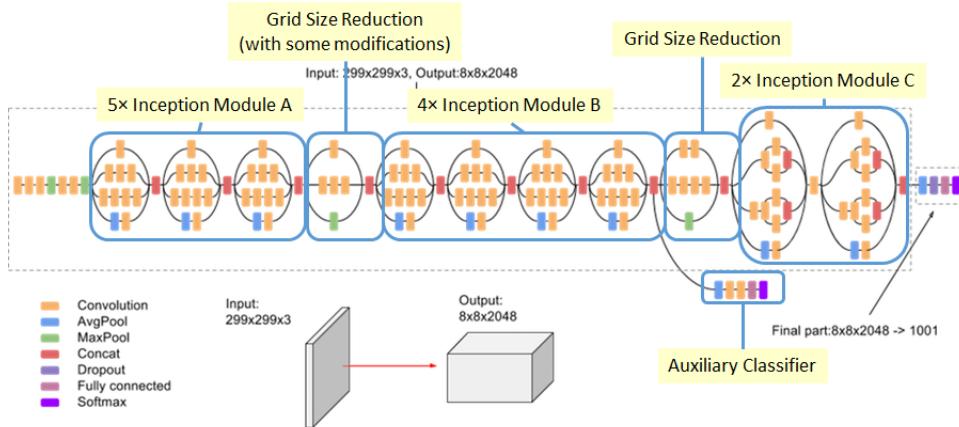


Figure 2.13: InceptionV3 architecture.[37]

Even though InceptionV4 [38] is already available it was not used in this work. The improvements over its predecessor are as follow:

1. Converting Inception modules to Residual Inception blocks.
2. Adding more Inception modules.
3. Adding a new type of Inception module (Inception-A) after the Stem module.

#### 2.4.4 | DenseNet

Densely Connected Convolutional Networks aim at expanding the depth of deep convolutional networks by connecting each layer to every other layer, in a feed forward fashion (this can be seen in figure 2.14), this reduces the number of parameters required, and alleviates the problem of the vanishing-gradients, while improving feature propagation (ensuring maximum information and gradient flow) and feature reuse which allows the learning of more compact and accurate models. This kind of neural network simplifies the connectivity pattern between layers introduced in other architectures (such as ResNets). [36]

The improved flow of information and gradients makes DenseNets easier to train, since each layer has direct access to the gradients from the loss function and the original input signal, leading to an implicit deep supervision.

DenseNets scale naturally to hundreds of layers, while exhibiting no optimization difficulties.

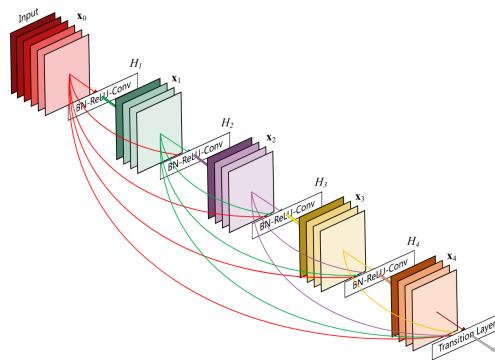


Figure 2.14: A 5-layer dense block. Each layer takes all preceding feature-maps as input. [36]

### 2.5 | Regression based algorithms for Object Detection

Regression based algorithms (also called single stage detectors) work differently than classification based algorithms. Instead of selecting multiple interesting parts of an image, they predict classes and bounding boxes for the entire image in one single run of the algorithm.

These algorithms are extremely fast but are not so accurate as classification based algorithms. [39] RetinaNet, YOLO and SSD are a few examples of object detection algorithms of this type.

#### 2.5.1 | RetinaNet

RetinaNet is a one-stage object detector presented at the 2017 International Conference on Computer Vision by the Facebook AI Research.

In order to improve performance a loss function was implemented, called Focal Loss, allowing the network to focus more on difficult samples. With the loss function, alongside a one-stage network architecture, RetinaNet is able to achieve state-of-the-art performance in terms of accuracy and running time.

This neural network is essential composed of one backbone network and two subnetworks. The backbone network is called Feature Pyramid Net [40], built on top of ResNet, and has

the purpose of computing convolutional feature maps of an image. Both subnetworks serve different purposes, one is for object classification using the backbone network output and the other subnetwork is responsible for performing the bounding box regression using the backbone network output.[39]

In figure 2.15 its observable the Feature Pyramid Network (FPN) on top of the convolutional neural network ResNet as a backbone network (a) to generate a rich convolutional feature pyramid (b). The class subnet (c) is for classifying anchor boxes, and the box subnet (d) is for regressing from anchor boxes to ground-truth object boxes.

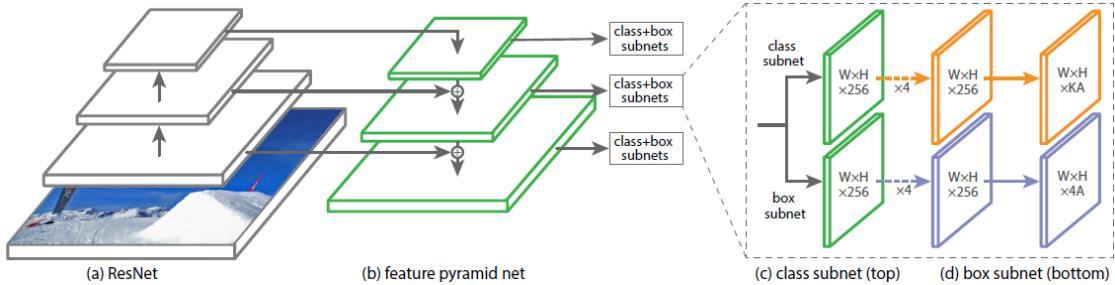


Figure 2.15: RetinaNet architecture.[39]

## 2.5.2 | YOLOv3

YOLOV3 (You Only Look Once Version 3) is a state-of-the-art, real-time object detection that's in the third iteration of the original YOLO, it's extremely fast and accurate (on par with the accuracy of focal loss from RetinaNet, but 4 times faster). YOLO allows the user to tradeoff between speed and accuracy simply by changing the size of the model.

Compared to other classification networks that perform predictions multiple times for various regions in an image, YOLO architecture is more like a fully convolutional neural network do to the fact that it takes an image as input and passes it only once through the FCNN. The network divides the image into regions and predicts bounding boxes (weighted by predicted probabilities) and probabilities to each region, outputting a vector of bounding boxes and classes predictions.

YOLO works by dividing an image in an  $S \times S$  grid and assuming  $B$  bounding boxes per grid. Each of the bounding box predicts 4 coordinates, object and class probabilities. [3]

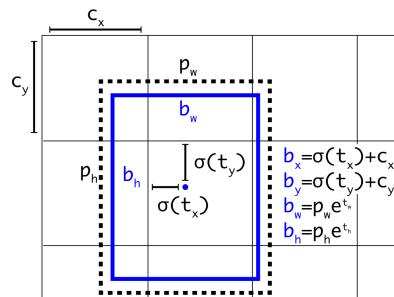


Figure 2.16: Bounding Box Prediction : Predicted Box (Blue), Prior Box (Black Dotted).[41]

YOLO image predictions are informed by global context in the image since it can look at the entire image at the test time. This gives it several advantages over classifier-based systems. In addition, this algorithm also uses an open source neural network called Darknet-53 for feature extraction, this neural network is written in C and CUDA and it supports CPU and GPU computation. [41]

The full architecture of YOLOv3 is represented in figure 2.17.

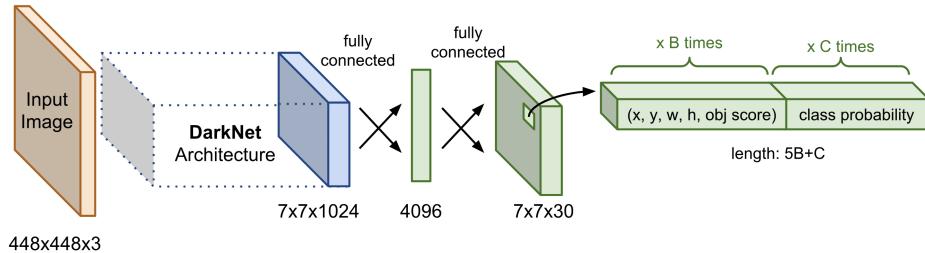


Figure 2.17: The network architecture of YOLO base model.[42]

### 2.5.3 | TinyYoloV3

TinyYOLOv3 is a smaller model of YOLOv3 that requires less computational resources since it doesn't occupy a large amount of memory , making it able to run in a smartphone. This model has a smaller number of convolutional layers, which improves the detection for small targets, therefore, it's a model best suited for constrained environments. In its architecture this network is composed of 7 convolutional layers and 6 pooling layers and can detect 80 different object categories. For complex scenes TinyYOLO is not accurate enough, however it is one of the fastest algorithms available. [43]

### 2.5.4 | Single Shot MultiBox Detector (SSD)

SSD is a method for detecting objects in images using a single deep neural network. This Multibox detector discretizes the output space of bounding boxes into a set of default boxes over different aspect ratios and scales per feature map location.

The base network of SSD is a VGG-16 network [44] followed by multibox convolutional layers. VGG-16 has the purpose of extracting the features for high quality image classification. The additional convolutional layers have the purpose of detecting objects, they are located at the end of the base network and decrease in size progressively, which helps with the detection of objects at multiple scales. The deep layers cover larger receptive fields and are helpful for larger objection detection, while the initial convolutional layers cover smaller receptive fields and are used for smaller objects detection. [45]

The added auxiliary structure can be summarized in the following key points:

- **Multi-scale feature maps for detection.** These layers decrease in size progressively and allow predictions of detections at multiple scales.
- **Convolutional predictors for detection.** Each added feature layer can produce a fixed set of detection predictions using a set of convolutional filters.

- **Default boxes and aspect ratios.** They associate a set of default bounding boxes with each feature map cell, for multiple feature maps at the top of the network. The default boxes tile the feature map in a convolutional manner, so that the position of each box relative to its corresponding cell is fixed.

The SSD architecture is represented in figure 2.18.

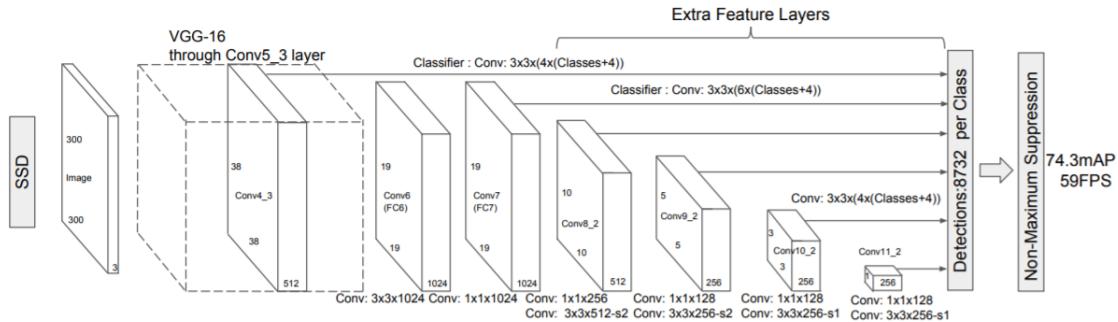


Figure 2.18: SSD architecture. [45]

The prediction of bounding boxes is done by multiple feature maps of different sizes that represent multiple scales. During prediction time, the network generates scores for the presence of each object category in each default box and produces adjustments to the box to better match the object shape. The network also combines predictions from multiple feature maps with different resolutions to naturally handle objects of various sizes.

SSD is simple relative to methods that require object proposals because it completely eliminates proposal generation and subsequent pixel or feature resampling stages and encapsulates all computation in a single network. This makes SSD easy to train.

The core of SSD is predicting category scores and box offsets for a fixed set of default bounding boxes using small convolutional filters applied to feature maps.

To achieve high detection accuracy, SSD produces predictions of different scales from feature maps of different scales, and explicitly separate predictions by aspect ratio.

These design features lead to simple end-to-end training and high accuracy, even on low resolution input images, further improving the speed vs accuracy trade-off.

This approach is based on a feed-forward convolutional network that produces a fixed-size collection of bounding boxes and scores for the presence of object class instances in those boxes, followed by a non-maximum suppression step to produce the final detections. [45]

## 2.6 | Classification Based Algorithms For Object Detection

Classification based algorithms work in two stages. Firstly, they select interesting regions from the image and secondly, they classify those regions using convolutional neural networks. The problem with this approach is that it can be extremely slow since a prediction is run for every selected region, however this approach is extremely accurate. [39]

RCNN, Fast-RCNN and Faster-RCNN are some types of classification based algorithms.

## 2.6.1 | R-CNN Models Summary

In the picture below a compact summary of all of the R-CNN models.

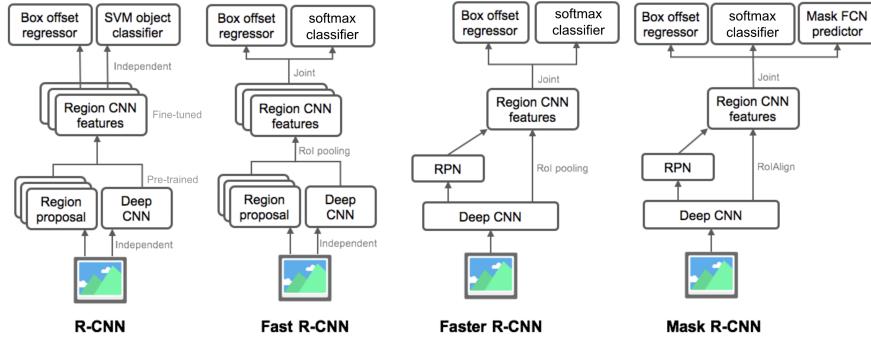


Figure 2.19: R-CNN model family summary. [46]

## 2.6.2 | R-CNN

The principal idea behind Region-based Convolutional Networks (R-CNN) can be split into two steps. In the first step the network identifies a number of regions of interest (bounding-box object region candidate) using a selective search method [46], which is a common algorithm to provide region proposals that can potentially contain objects [47].

In the second step it extracts CNN features from each region independently for the classification.

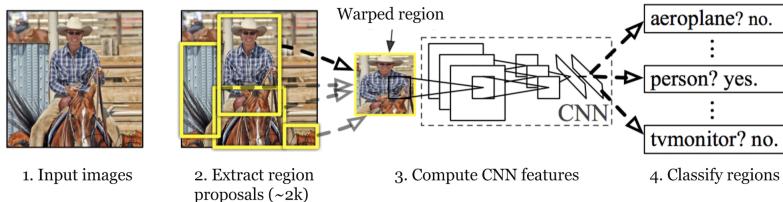


Figure 2.20: R-CNN architecture. [46]

## 2.6.3 | Fast R-CNN

The idea behind Fast-RCNN [48] is, as the name implies, to make R-CNN faster. In order to achieve this, the training procedure was improved by unifying three independent models into one jointly trained framework and increasing shared computation results.

In this new improved network, the CNN feature vectors are not extracted independently for each region proposal, instead this model aggregates them into one CNN forward pass over the entire image and the region proposals share the feature matrix.

This feature matrix is then branched to be used for learning the object classifier and the bounding-box regression.

In short, computation sharing improves the speed of R-CNN. [46]

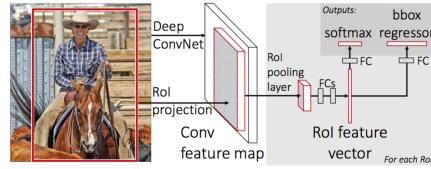


Figure 2.21: Fast R-CNN architecture. [46]

#### 2.6.4 | Faster R-CNN

The Faster R-CNN [49] improves upon the previous considered solutions since it integrates the region proposal algorithm directly into the CNN model. It can be seen as a single, unified model composed of a region proposal network and fast R-CNN with shared convolutional feature layers. [46]

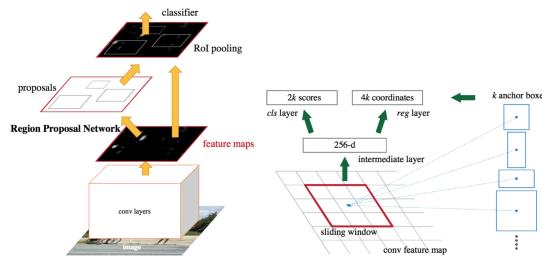


Figure 2.22: Faster R-CNN architecture. [46]

#### 2.6.5 | Mask R-CNN

The final model of the R-CNN family, mask R-CNN [50], extends faster R-CNN to pixel-level image segmentation by decoupling the classification and the pixel-level mask prediction tasks. It adds a third branch for predicting an object mask in parallel with existing branches for classification and localization, based of the Faster R-CNN framework. This new mask branch predicts a segmentation mask in a pixel-to-pixel manner.

Mask R-CNN improves the region of interest pooling layer because pixel-level segmentation requires much more fine-grained alignment than bounding boxes. This allows the region of interest to more precisely map regions of the original image. [46]

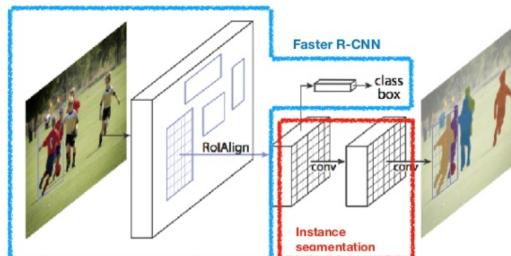


Figure 2.23: Mask R-CNN is a Faster R-CNN model with image segmentation. [46]

## 2.7 | State-Of-The-Art

Image classification and object detection are both subjects that are constantly innovating and improving upon previous results, every month new papers are published with new and more efficient networks.

In the figures 2.24 and 2.25 it is shown not only the current best methods for both image classification and object detection but also the development of the state-of-the-art throughout the years.

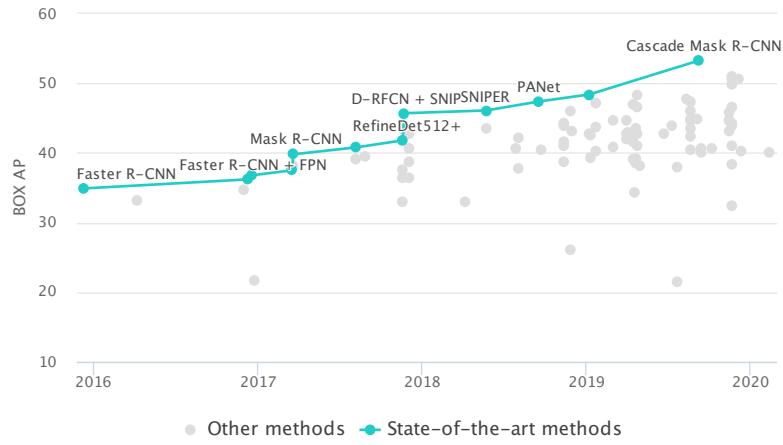


Figure 2.24: Object Detection on COCO test-dev benchmark .[51]

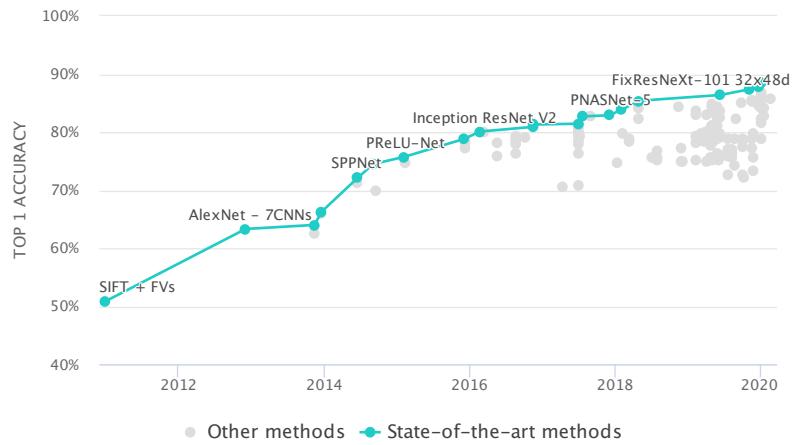


Figure 2.25: Image Classification on ImageNet benchmark. [52]

Due to the fact that object detection is a subject of great innovation, there is an extreme amount of papers that try to compete for the best results coming out every few months. So, in order to do a review of the state of the art the tables, 2.1 and 2.2 show the benchmarks for both imageNet and COCO test-dev. These tables were obtained from [2] and are based on the analysis of [52] and [51] which is a website dedicated to the current state-of-the-art for object detection and image classification.

### 2.7.1 | COCO Test-Dev

The COCO benchmark [16] is a dataset that places object recognition in the context of scene understanding. The evaluation metric used is the average precision (AP). Table 2.1 shows the current best architectures and their respective score for the COCO test-dev dataset.

Table 2.1: COCO Test-Dev Benchmarks.

Method	Backbone	AP (%)
Liu et al.(2019) [53]	ResNeXt-152	53.3
Tan et al. (2019) [54]	EfficientNet	51.0
Zhang et al. (2019) [55]	ResNeXt-101	50.7
Girshick et al. (2018) [56]	ResNeXt-152	50.2
Li et al. (2019) [57]	ResNet-101	48.4
Zhang et al. (2019) [55]	ResNet-101	46.3
Mahajan et al. (2018) [58]	ResNeXt	45.2
Zhao et al. (2019) [59]	VGG16	44.2
Cai et al. (2018) [60]	ResNet-101	42.8
Wang et al. (2019) [61]	ResNet-50	39.8
Lin et al. (2017) [39]	ResNet-101	39.1
Shrivastava et al. (2016) [62]	Inception-ResNet-v2	36.8
Kim et al. (2018) [63]	VGG-16	35.2

Liu et al. [53] achieved the best score in the COCO Test-Dev in 2019. They proposed better detection performance by creating a more powerful backbone network from previously existing backbones like ResNet [33] and ResNetXt [64]. They implemented a strategy for assembling multiple identical backbones (called Assistant Backbones and Lead Backbones) linked by composite connections between the adjacent backbones in order to form a more powerful backbone which was given the name of Composite Backbone Network (CBNet).

In typical CNN based detectors, the backbone network (the baseline of a network architecture) is used for basic feature extraction.

CBNet feeds the output features of the previous backbone as an input feature to the succeeding backbone through composite connections. At the final stage, the Lead Backbone outputs features for object detection.

This architecture was able to achieve the best result in the COCO Test-Dev with a 53.3% AP with single model by integrating a CBNet using triple ResNeXt-152 [64] backbones into the Cascade Mask R-CNN baseline.

Figure 2.26 presents the architecture for CBNet.

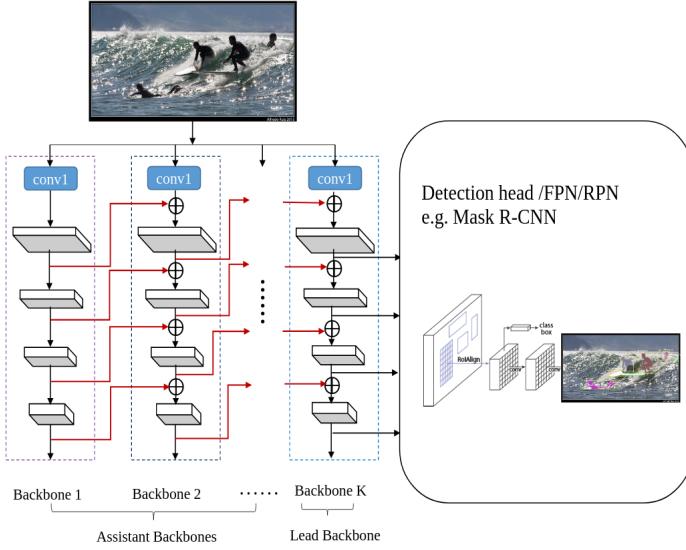


Figure 2.26: CBNet Architecture for object detection.

### 2.7.1.1 | ResNeXt

ResNeXt, also known as Aggregated Residual Transform Network was created by facebook researchers and it is a simple highly modularized network architecture for image classification.

The network is constructed by repeating a building block that aggregates a set of transformations with the same topology. The simple design results in a homogeneous, multi-branch architecture that has only a few hyper-parameters to set. This strategy creates a new dimension, which was given the name of "cardinality" (size of the set of transformations).

This architecture is an improvement over the Inception architectures, being more simple in design and adding more branches (towers) within modules. [64]

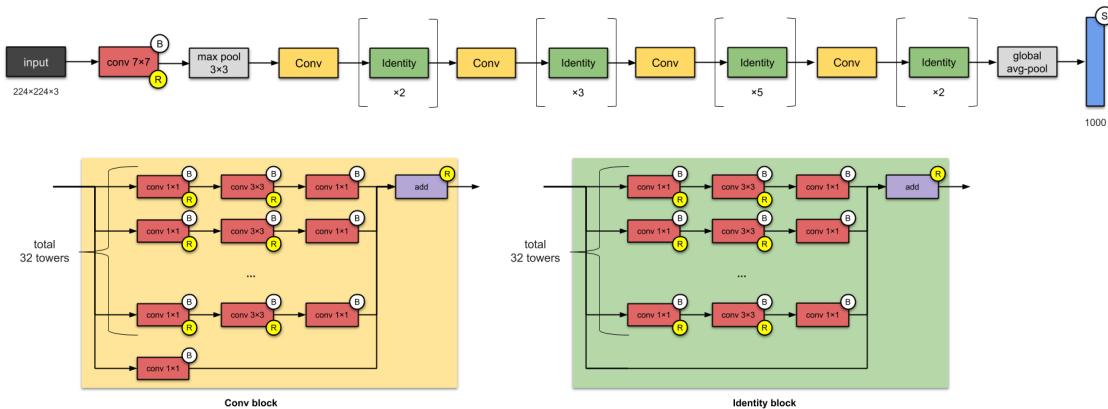


Figure 2.27: ResNeXt architecture. [34]

## 2.7.2 | ImageNet

The imageNet Large Scale Visual Recognition challenge [65] is a benchmark for object category classification and detection. The evaluation metrics used are top-1 and top-5 accuracy.

Table 2.2: ImageNet Benchmarks.

Method	Backbone	Top-1 Acc (%)
Xie et al. (2019) [66]	EfficientNet	88.4
Kolesnikov et al. (2019) [67]	ResNet-152	87.8
Touvron et al. (2019) [68]	ResNeXt-101	86.4
Xie et al. (2019) [69]	EfficientNet	85.5
Mahajan et al. (2018) [58]	ResNeXt	85.4
Tan et al. (2019) [70]	EfficientNet	84.4
Touvron et al. (2019) [68]	ResNet-50	82.5
Szegedy et al. (2017) [38]	Inception-resnet-v2	80.1
Szegedy et al. (2017) [38]	Inception-v4	80.0
Simonyan et al. (2014) [44]	VGG-16	74.4

Xie et al. [66] stated that current state-of-the-art vision models are still trained with supervised learning, which implies the necessity of large corpus of labeled images in order to work properly. The fact that current models are only shown labeled images causes an obvious limitations in the improvement of accuracy and robustness of current state-of-the-art models, this can be improved with the usage of the large available quantities of unlabeled images available.

Having this in mind, they decided to use unlabeled images to improve the state-of-the-art ImageNet accuracy and show that accuracy has an outsized impact on robustness. For this purpose, they used a much larger corpus of unlabeled images, where a large fraction of images did not belong to ImageNet training set distribution.

Using a self-training framework the model was trained with 3 main steps which consist in:

1. Training of a teacher model on labeled images.
2. Usage of the teacher to generate pseudo labels on unlabeled images.
3. Train a student model on the combination of labeled images.

The algorithm was iterated a few times by treating the student as a teacher to relabel the unlabeled data and training a new student.

An important discovery was made during the training of the algorithm. For the method to work well at scale the student model should be noised during its training while the teacher should not be noised during the generation of pseudo labels. This way, the pseudo labels are as accurate as possible and the noised student is forced to learn harder from the pseudo labels.

To induce noise in the model it was used RandAugment data, dropout and stochastic depth during the training. Figure 2.28 shows a brief view of how the method works.

This is where the name of the method "Noisy Student" comes from, since the student is noised to learn beyond the teacher's knowledge.

With this method they were able to show that it is possible to use unlabeled images to significantly advance both accuracy and robustness of state-of-the-art imageNet models.

The presented model uses EfficientNet (this architecture is explained in more detail in ??) as a backbone trained on images from imageNet dataset and was able to obtain the best results in the ImageNet benchmark dataset by achieving an accuracy of 88.4%.

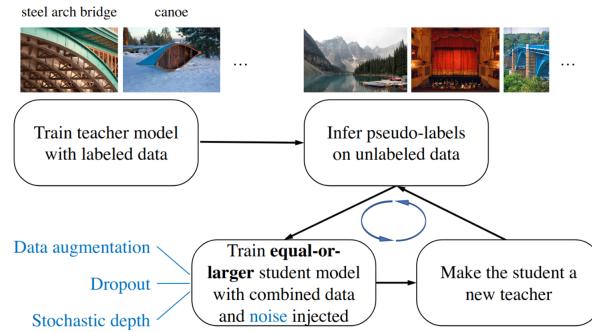


Figure 2.28: Noisy Student Method. [66]

### cnnarchitectures

Researchers at Google decided to study the impact of scaling up CNNs, in order to achieve better accuracy and efficiency. EfficientNet-B0 was developed based on a simple idea, scaling each of the dimensions of the network (width, depth and resolution) with a constant ratio, improves the overall performance [70].

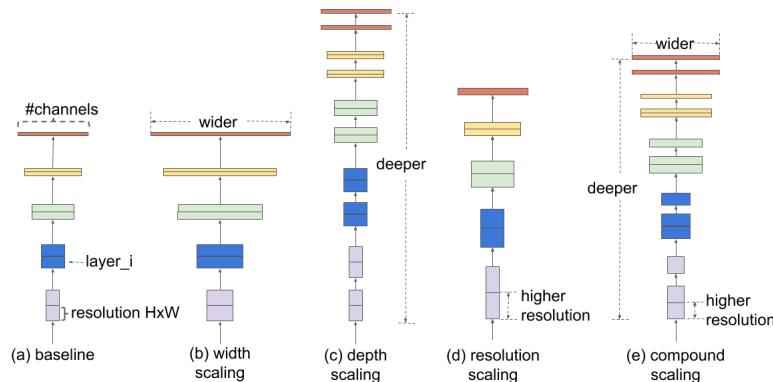


Figure 2.29: Comparison of different scaling methods: (a) is a baseline network example; (b)-(d) are conventional scaling that only increases one dimension of network width, depth, or resolution. (e) is the proposed compound scaling method that uniformly scales all three dimensions with a fixed ratio. [70]

The baseline network architecture, EfficientNet-B0, uses mobile inverted bottleneck convolution (MBConv), similar to MobileNetV2 [71] and MnasNet [72]. Figure shows the baseline network architecture EfficientNet-B0.

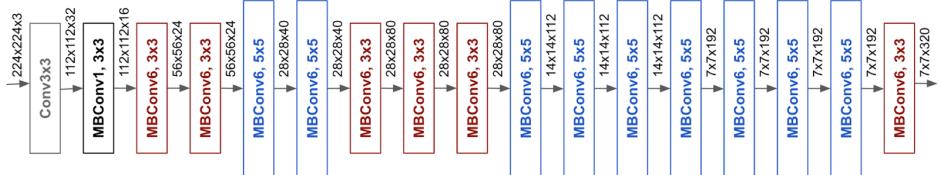


Figure 2.30: EfficientNet-B0 architecture representation.

# CHAPTER 3

## Information Extraction From Text

As previously discussed in chapter 2 the problem of big data has become more relevant in the recent years. There is too much data to be analysed and the human brain is incapable of processing such large quantities of information. Therefore, information extraction becomes more prevalent as data increases.

Information Extraction (IE) is the task of automatically extracting pre-specified information from textual sources, a trivial example is the usage of IE to analyse large quantities of documents only retrieving the relevant information. In most cases the task of extracting information is concerned with the processing of human language texts by means of natural language processing (NLP).

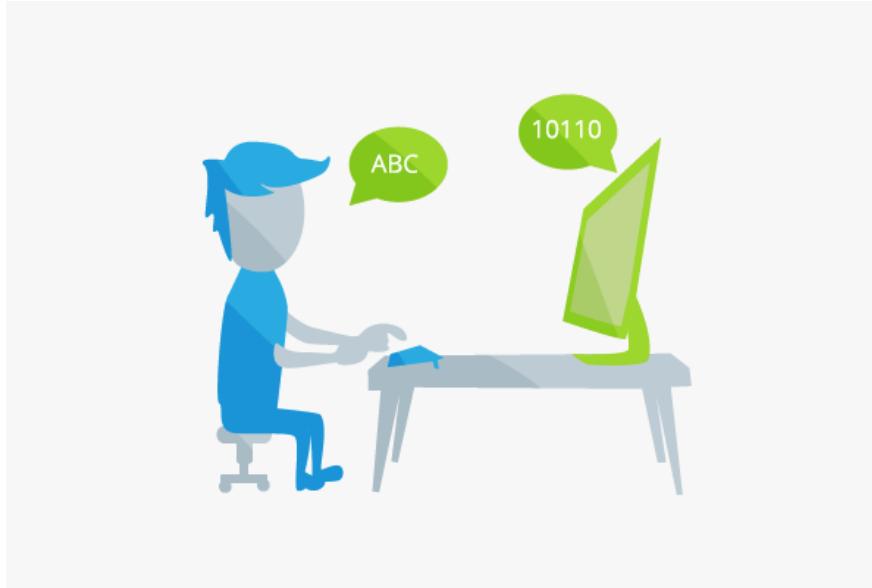


Figure 3.1: Natural Language Processing [73].

This chapter starts with an introduction to Natural Language Processing in section 3.1. Section 3.2 explains the process of representing text in a numerical vector form while describing the concept of word embeddings. Static and contextualized word embedding models are discussed in section 3.3 and 3.4 respectively. An overview on some of the available NLP libraries is presented in section 3.5.

## 3.1 | Natural Language Processing

Natural language processing is a subfield of linguistics, computer science, information engineering and artificial intelligence, which is devoted to the engineering of computational models and processes to give the ability of human language understanding to computers. [74]

Human language is extremely complex and rarely precise, to understand it is to understand not only the words, but the concepts and how they are linked together in order to create meaning. This makes NLP one of the most difficult tasks in computer science.

Figure 3.2 shows the classification of NLP, which consists in two major components, Natural Language Understanding (NLU) and Natural Language Generation (NLG) [74].

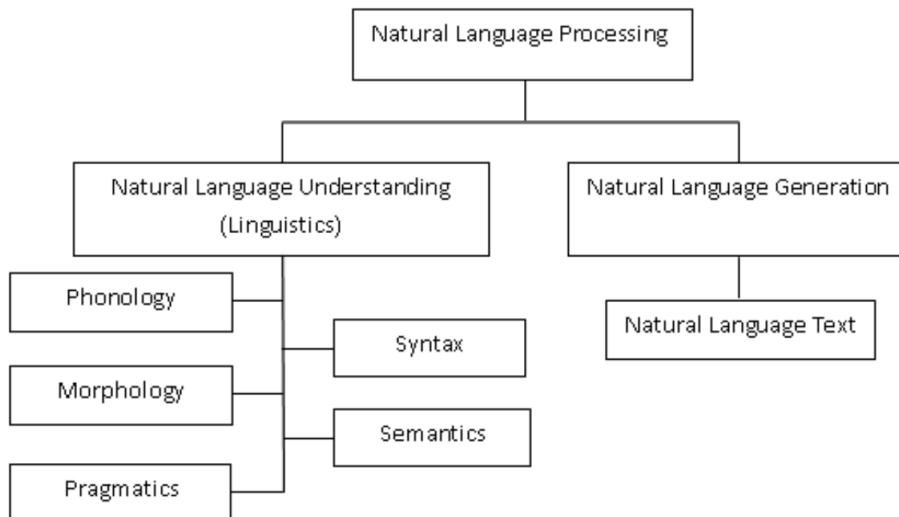


Figure 3.2: Classification of NLP. [74]

Natural Language Understanding is the process of understanding text. It is related to the science of Linguistic that studies the meaning of languages, context and various forms of language.

Natural Language Generation is the process of generating text, sentences and paragraphs that are meaningful from an internal representation [74].

### 3.1.1 | Important NLP Terminologies

**Phonology:** The part of Linguistics which refers to the systematic arrangement of sound.

**Morphology:** In linguistics, morphology is the study of words, how they are formed, and their relationship to other words in the same language. The different parts of the word represent the smallest units of meaning known as Morphemes.

**Lexical:** In Lexical the focus is the interpretation of the meaning of individual words.

**Syntax:** Syntax refers to the study of the grammatical structure of the sentence.

**Semantic:** Semantic processing determines the possible meanings of a sentence by pivoting on the interactions among word-level meanings in the sentence.

**Discourse:** Discourse focuses on the properties of the text as a whole that convey meaning by making connections between component sentences.

**Pragmatic:** Pragmatic is a subfield of linguistics that studies the ways in which the context of a sentence contributes to the meaning [74].

### 3.1.2 | Core Areas

The field of NLP can be divided in two broad sub-areas: core areas and application areas. The core areas address fundamental problems such as language modeling, morphological processing, parsing and semantic processing. Language modeling underscores quantifying associations among naturally occurring words. Morphological processing deals with the segmentation of meaningful components of words and the identification of the true parts of speech of words used. Parsing consists in the building of sentence diagrams as possible precursors to semantic processing. Semantic processing attempts to distill meaning of words, phrases, and higher level components in text [75].

### 3.1.3 | Application Areas

The application areas address topics such as extraction of useful information from text (e.g named entities and relations), translation of text, summarization of written documents, automatic answering of questions, chat bots, email spam detection and many others [75].

## 3.2 | Numerical Representation of Text

Machine learning algorithms and most of all deep learning architectures are incapable of processing strings of text, this is because they require numbers as an input. [76] A human can easily tell that the word "dog" and the word "cat" are identical, since they both represent an animal, however a computer would assume that they are completely different things since all the letters in those words are different.

### 3.2.1 | Word Embeddings

The dominant approach to solve this problem is the usage of word embeddings, which is a type of word representation that allows words with similar meaning to have a similar representation by mapping a set of words, or phrases in a vocabulary, to vectors of numerical values. For example, the word "happy" can be represented as a vector of 4 dimensions [0.24, 0.45, 0.11, 0.49] and "sad" has a vector of [0.88, 0.78, 0.45, 0.91]. The reason for this vectors to exist is so that a machine learning algorithm can perform linear algebra operations on numbers (vectors) instead of words [77].

Word embedding methods learn a real-valued vector representation for a predefined fixed size vocabulary from a corpus of text [78].

A vector representation of a word may be a one-hot encoded vector where 1 stands for the position where the word exists and 0 everywhere else.

As an example, the sentence "Word Embeddings are Word converted into numbers" can be converted to the following dictionary using the one-hot encoded vector representation : ['Word', 'Embeddings', 'are', 'Converted', "Word", 'into', 'numbers'] .

Using this representation the word "numbers" in the one-hot encoded vector is [0,0,0,0,0,1] and for the word "converted" is [0,0,0,1,0,0]. This is considered to be the most simple method to represent words in vector forms [76].

The following image showcases the given example.

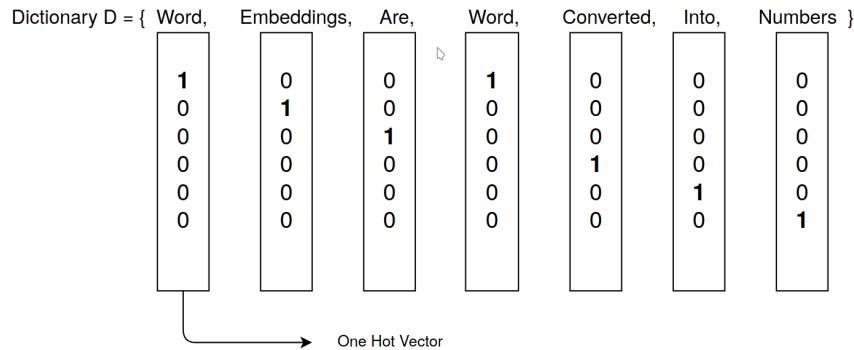


Figure 3.3: Example of text representation by one-hot vector.

### 3.3 | Static Word Embedding Models

This section introduces some common static word embedding models to learn word embeddings from text.

Static word embedding have the fundamental problem which is they generate the same embedding, in different contexts, for the same word, failing to capture the polysemy of the word. This is due to the fact that each word has a single vector, regardless of context. [79]

As an example, having these two phrases:

- "The Apple Company is the one who produces iPhones."
- "This apple is delicious."

In this case, the word "Apple" has two different meanings, being one a company and the other a fruit, however for static word embedding models, words only have one single meaning, and therefore the word representation for "Apple" would be the same for both cases. [80]

#### 3.3.1 | Word2Vec

Developed by Tomas Mikolov, et al. at Google in 2013, Word2Vec is a two-layer neural network that processes text by "vectorizing" words with the purpose of grouping vectors of similar words together in vectorspace. The way Word2Vec detects those similarities is by creating vectors that are distributed numerical representations of word features, without human intervention.

In a regular one-hot encoded vector, all words have the same distance between each other, even though their meanings are completely different.

Using Word2Vec, the resulting vector is able to maintain context.

Word2Vec is capable of making accurate guesses, based on past appearances, of a word's meaning.

The output of Word2Vec is a vocabulary in which each item has a vector attached to it, which can be fed into a deep-learning net or simply queried to detect relationships between words.

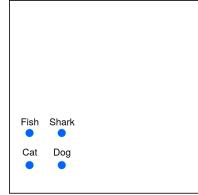


Figure 3.4: One-hot encoding resulting vector. [81]

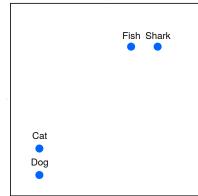


Figure 3.5: Word2Vec encoding resulting vector. [81]

Word2Vec is composed of two different models, CBOW (Continuous Bag of words) which predicts a word given the context and Skip-Gram which predicts context given a word. [79] [82]

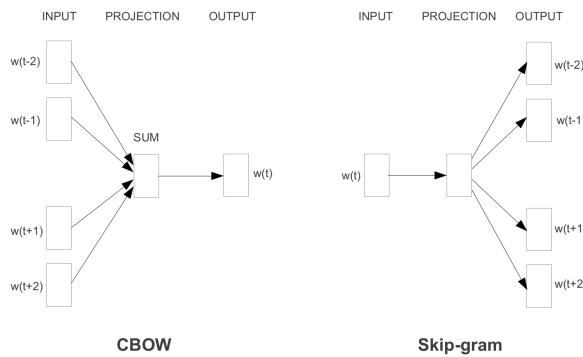


Figure 3.6: CBOW model and Skip-Gram model. [79]

### 3.3.2 | GloVe

GloVe stands for Global Vectors for Word Representation and was a new approach created by Pennington et all. in 2014 [83] to generate word embeddings with unsupervised learning. Glove main goals are to create word vectors that capture meaning in the vector space and to take advantage of global count statistics instead of using only local information.

The problem with Word2Vec is that it only takes local information into account, and does not consider global context. This means that the semantics learnt for a given word are only affected by the surrounding words.

GloVe works by aggregating global word-to-word co-occurrence matrix from a corpus of text. This means that if two words keep appearing together in a corpus of text they either share a linguistic or a semantic similarity. Simply put, similar words will be placed together in

the high-dimensional space. Therefore GloVe can be seen like an extension to the Word2Vec model.

### 3.3.3 | FastText

FastText, created by Facebook's AI Research (FAIR) lab in 2016, is a fast text classifier based on the skipgram model used for efficient learning of word representations and sentence classification. Popular models like word2Vec and GloVe are based on continuous word representations that create vectors directly from words in a sentence while ignoring the morphology of words, this is done by assigning a distinct vector to each word, fastText uses a different approach treating each word as bag of characters n-grams. A vector representation is associated to each character n-gram and words are represented as the sum of these representations. This allows fastText to work with rare words not seen in the training data since the word is broken down into n-grams to get the corresponding embeddings [84].

Using the word "where" as an example and n=3, the representation of this word in a fastText model is <wh, whe, her, er, re> and the special sequence <where>. The angular brackets serve as boundary symbols to distinguish the n-gram of a word from the word itself, this means that if the word "her" was part of the vocabulary it would be represented as <her>, which allows the preservation of the meaning of shorter words and the understanding of suffixes and prefixes.

## 3.4 | Contextualized Word Embedding Models

Contextualized words embeddings aim at capturing word semantics in different contexts to address the issue of polysemous and the context-dependent nature of words [80]. Using the example given in section 3.3, these models would be able to distinguish the different meaning of the word "apple" given the two different sentences.

### 3.4.1 | Context2vec

Context2Vec is an unsupervised model capable of learning efficiently generic context embedding of wide sentential contexts, using a bidirectional LSTM.

A large plain text corpora is utilized in order to learn a neural model capable of embedding entire sentential contexts and target words in the same low-dimensional space, which is optimized to reflect inter-dependencies between targets and their entire sentential context as a whole.

In contrast to word2vec that use context modeling mostly internally and considers the target word embeddings as their main output, the focus of context2vec is the context representation. Context2vec achieves its objective by assigning similar embeddings to sentential contexts and their associated target words [85].

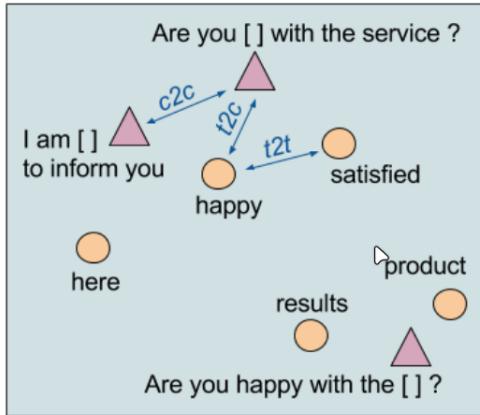


Figure 3.7: A 2D illustration of context2vec’s embedded space and similarity metrics. Triangles and circles denote sentential context embeddings and target word embeddings, respectively [85].

Sentential Context	Closest target words
This [ ] is due	item, fact-sheet, offer, pack, card
This [ ] is due not just to mere luck	offer, suggestion, announcement, item, prize
This [ ] is due not just to mere luck, but to outstanding work and dedication	award, prize, turnabout, offer, gift
[ ] is due not just to mere luck, but to outstanding work and dedication	it, success, this, victory, prize-money

Figure 3.8: Closest target words to various sentential contexts, illustrating context2vec’s sensitivity to long range dependencies, and both sides of the target word [85].

### 3.4.2 | ELMo

ELMo (Embeddings from Language Models) is a NLP model with context-aware representation, it understands different meanings for the same word since it takes into account the surrounding words unlike traditional word embedding models such as Word2Vec and GLoVe. In order to achieve this, ELMo attributes an embedding for each word after looking at the entire context in which it is used, instead of using fixed embeddings for each word. Therefore, the same word might have different word vectors under different contexts.

This NLP models both syntax and semantics of word use and how these uses vary across linguistic context. The word vectors are learned through the usage of internal states of a deep bidirectional LSTM algorithm, trained on a large corpus of text. Bidirectional implies that the algorithm takes into account the words before and the words after it in both directions. LSTM (Long Short-Term Memory) is one type of neural network that is able to retain data in memory for long periods of time, allowing it to learn longer-term dependencies. This language model can predict both the next word and the previous word and it is a character based model allowing the network to use morphological clues to form robust representations for out-of-vocabulary tokens not presented during training. [86]

Below an image showcasing an example of the differences between GLOVe that is a non-context aware model and ELMo biLM (bidirectional Language Model) that is context aware.

Source		Nearest Neighbors
GloVe	play	playing, game, games, played, players, plays, player, Play, football, multiplayer
biLM	Chico Ruiz made a spectacular <u>play</u> on Alusik 's grounder {...}	Kieffer , the only junior in the group , was commended for his ability to hit in the clutch , as well as his all-round excellent play .
	Olivia De Havilland signed to do a Broadway <u>play</u> for Garson {...}	{...} they were actors who had been handed fat roles in a successful <u>play</u> , and had talent enough to fill the roles competently , with nice understatement .

Figure 3.9: Nearest neighbors to "play" using GLoVe and context embeddings from a biLM [86].

GLoVe only uses the word "play" as source, therefore the obtained neighbors for that word are spread across several parts of speech however they all focus on the sports-related sense of the word "play". ELMo biLM uses the entire sentence as source, this means that it is able to understand the context of the word, therefore in both cases, the biLM is able to disambiguate both the part of speech and word sense in the source sentence [86].

## 3.5 | Available NLP libraries

### 3.5.1 | SpaCy

SpaCy is a free, open-source library for advanced natural language processing written in Python and Cython published by Explosion AI. It was designed specifically for production use and to help in the building of applications that process and "understand" large volumes of text data. Some use cases for this specific library are to build information extraction or natural language understanding systems, or to pre-process text for deep learning. [87]

This NLP library was chosen for the development of the text processing phase of the practical work, not only because it provides a well written documentation and being simple to use but also because it offers many useful features such as:

- **Tokenization** : The segmentation of text into words, punctuation, etc
- **Part-of-Speech Tagging** : The assignment of word types to tokens, like verb, noun, etc
- **Similarity** : The comparison between different words, phrases or text documents and how similar they are.
- **Lemmatization** : The assignment of base forms of words.

### 3.5.2 | Natural Language ToolKit

Developed by Steven Bird, Edward Loper and Ewan Klein in the Department of Computer and Information Science at the University of Pennsylvania, NLTK (Natural Language ToolKit) is a suite of open source program modules, tutorials, problem sets and a leading platform for building Python programs to work with human language data. NLTK covers symbolic and

statistical natural language processing, and is interfaced to annotated corpora. This library provides easy-to-use interfaces such as WordNet, along with a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning [88].

### 3.5.3 | Stanford Core NLP

Developed at Stanford University, Core NLP is library written in Java, however with wrappers for different languages, including Python. This library is fast and some of its components can be integrated to NLTK which boosts efficiency. CoreNLP enables users to derive linguistic annotations for text, including token and sentence boundaries, parts of speech, named entities, numeric and time values, dependency and constituency parses, coreference, sentiment, quote attributions, and relations [89].

### 3.5.4 | Gensim

Gensim ("Generate Similar") is a Natural Language Processing open-source library for unsupervised topic modeling (a technique to extract the underlying topics from large volumes of text) and for natural language processing. This python-cython library specializes in finding the semantic similarity between two documents through vector space modeling and topic modeling toolkit. It is capable of building document or word vectors, corpora, performing topic identification, performing document comparison (retrieving semantically similar documents) and analysing plain-text documents for semantic structure. In terms of producing word embedding, gensim allows for the usage of Word2Vec and fastText [90].

### 3.5.5 | Others

Other NLP libraries not so common are Flair [91], Polyglot [92], CogCompNLP [93], TextBlob [94].



# CHAPTER 4

---

## Initial Work

---

In this chapter a fundamental tool used was imageAI, which is a Computer Vision Python library that allows developers the ability to easily use state-of-the-art AI features . It supports algorithms for image prediction, custom image prediction, object detection, video detection, video object tracking and image predictions trainings. ImageAI also supports object detection, video detection and object tracking using RetinaNet, YOLOv3 and TinyYOLOv3 trained on COCO dataset. In terms of Machine Learning algorithms imageAI supports 4 trained on the imageNet-1000 dataset. [95]

With the goal of finding the best performing neural network or algorithm for image recognition and for object detection a few test runs were made. These test runs consist on feeding each of the neural networks and algorithms available for image recognition and object detection with one picture manually chosen beforehand. Each neural network and algorithm makes five guesses on what the image represents with a prediction probability that ranges in an interval between [0,100]. This prediction probability represents the certainty of the neural network and algorithm in its guess.

In sections 4.1 and 4.2 an overview of three example test runs are presented and analysed.

### 4.1 | Image Recognition test runs

The imageAI library allows the usage of 4 deep learning neural networks for image recognition which are DenseNet, inceptionV3, ResNet50, and SqueezeNet. The explanation on how these neural networks function is presented in chapter 2. The analysis of the results is in section 4.1.4

#### 4.1.1 | Test Run Number 1

For the first run an image of a dog (breed saluki) was analysed by all 4 neural networks.



Figure 4.1: First picture to be analysed.

The obtained results can be seen in the figure below.

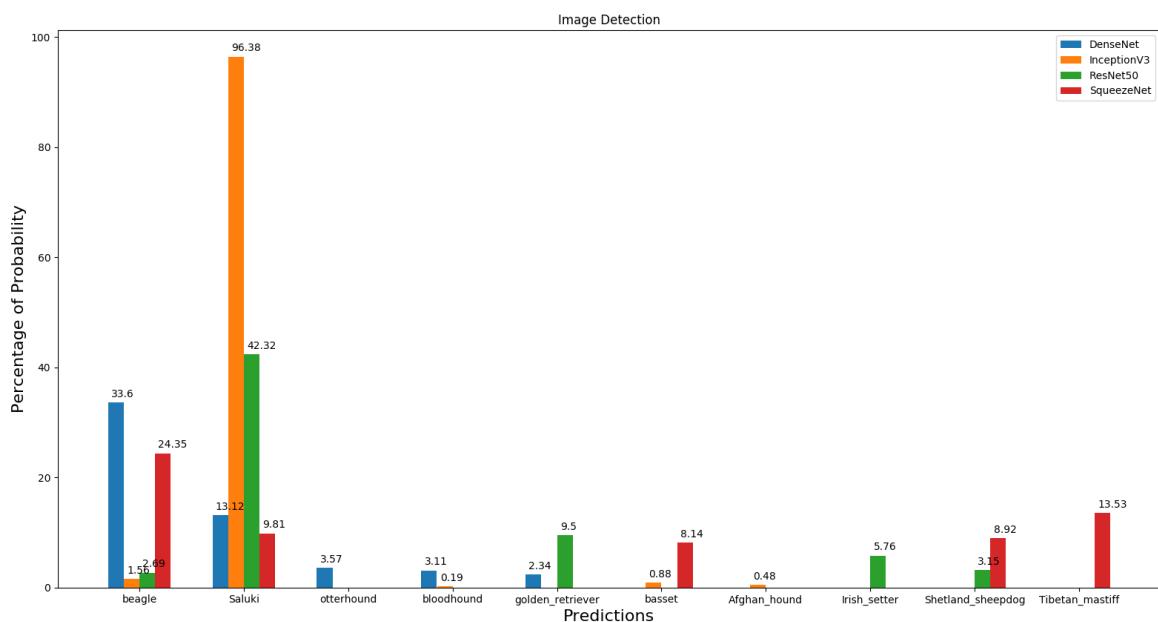


Figure 4.2: Results obtained.

#### 4.1.2 | Test Run Number 2

For the second run an image of a car was analysed by all 4 neural networks.



Figure 4.3: Second picture to be analysed.

The obtained results can be seen in the image below.

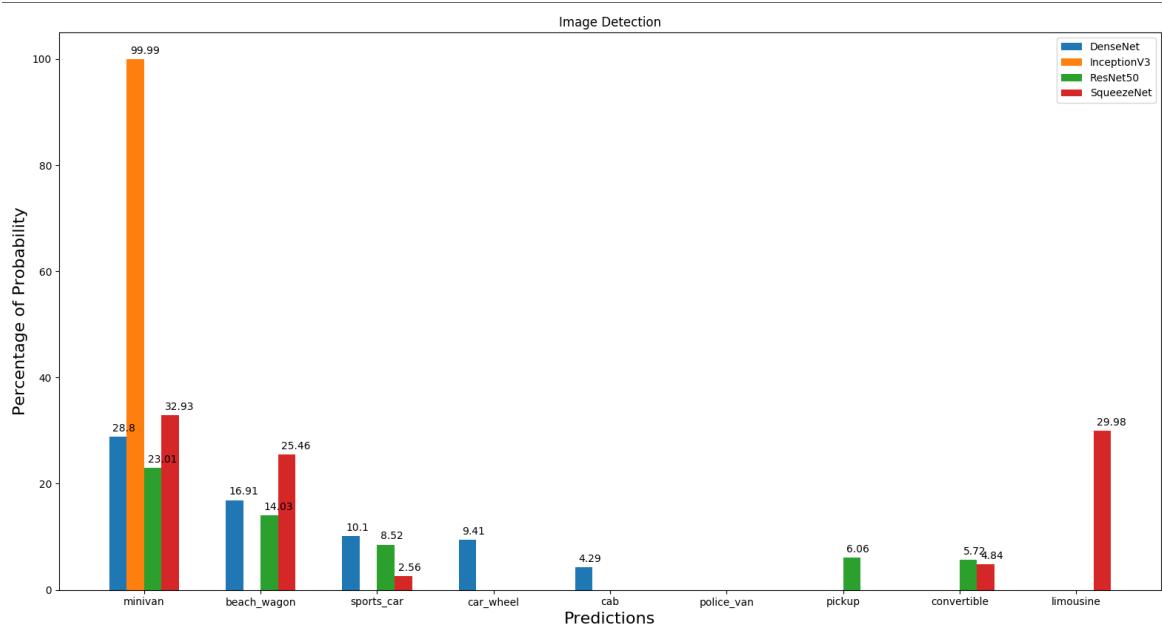


Figure 4.4: Results obtained.

### 4.1.3 | Test Run Number 3

For the last run an image of an espresso was analysed by all 4 neural networks.



Figure 4.5: Third picture to be analysed.

The obtained results can be seen in the figure below

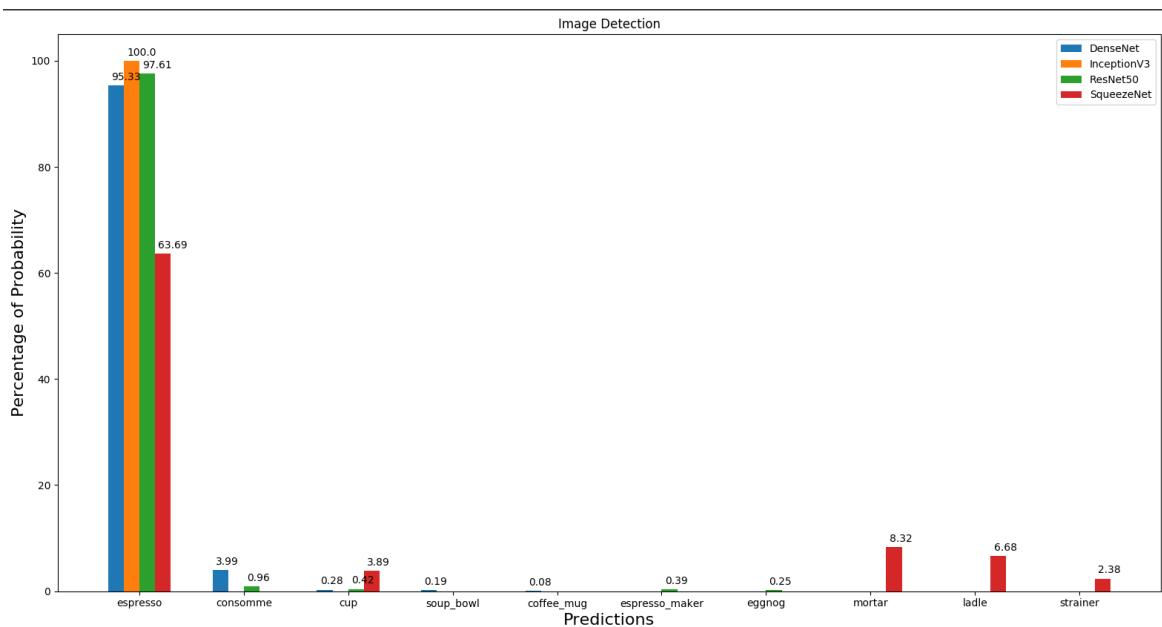


Figure 4.6: Results obtained.

#### **4.1.4 | Results analysis**

## 4.2 | Object Recognition test runs

ImageAI provides 3 different models trained on the COCO dataset for object recognition that are able to identify up to 80 of the most common objects in everyday life. The models provide include RetinaNet, YOLOv3 and TinyYOLOv3. [95]

The explanation on how these algorithms work is presented in chapter 2. The analysis of the results is in section 4.2.4

### 4.2.1 | Test Run Number 1



Figure 4.7: First picture to be analysed.

#### 4.2.1.1 | RetinaNet Results



Figure 4.8: RetinaNet Detections.

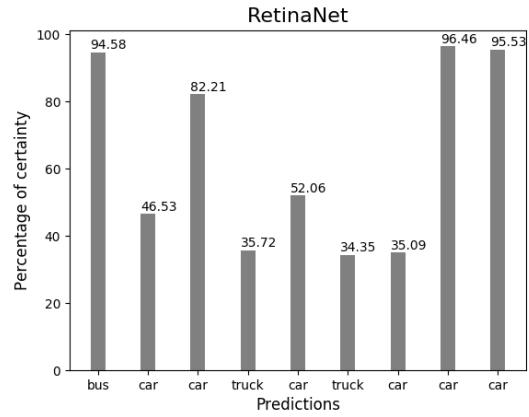


Figure 4.9: RetinaNet Detections.

#### 4.2.1.2 | Yolo Results



Figure 4.10: YOLO Detections.

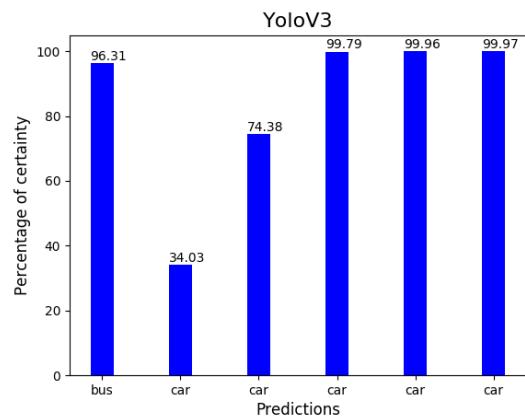


Figure 4.11: YOLO Detections.

#### 4.2.1.3 | TinyYolo Results



Figure 4.12: TinyYolo Detections.

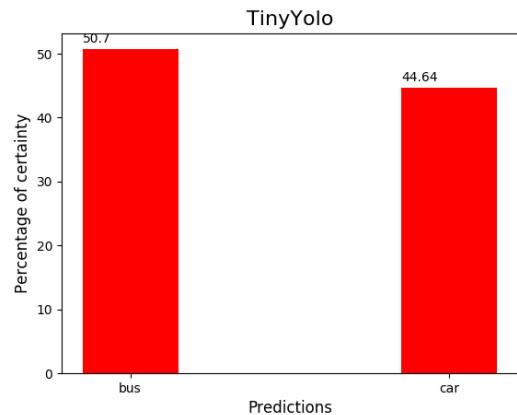


Figure 4.13: TinyYolo Detections.

## 4.2.2 | Test Run Number 2



Figure 4.14: Second picture to be analysed.

#### 4.2.2.1 | RetinaNet run 2 results

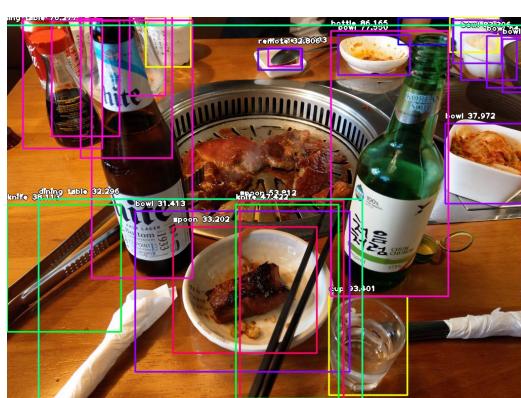


Figure 4.15: RetinaNet Detections.

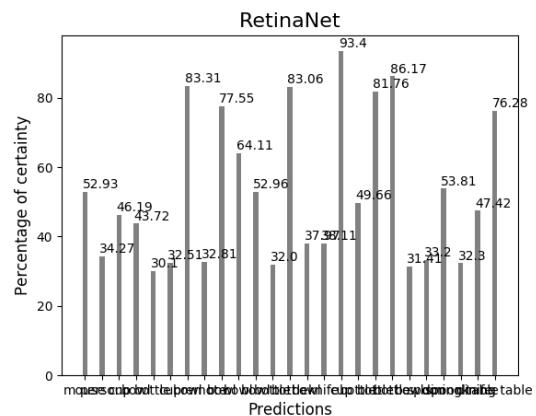


Figure 4.16: RetinaNet Detections.

#### 4.2.2.2 | Yolo run 2 results

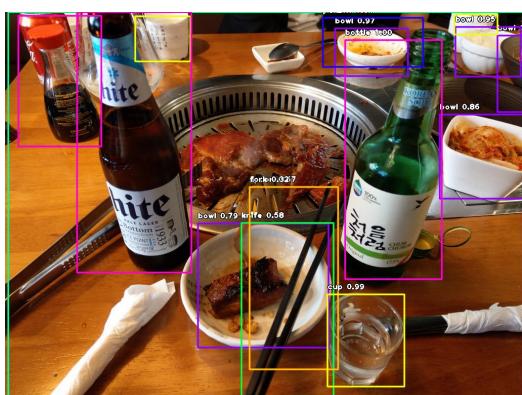


Figure 4.17: YOLO Detections.

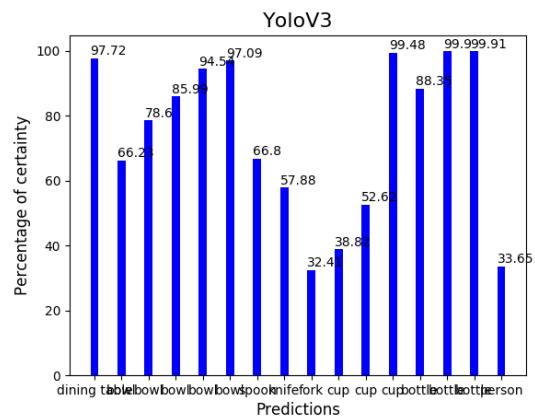


Figure 4.18: YOLO Detections.

#### 4.2.2.3 | TinyYolo run 2 results

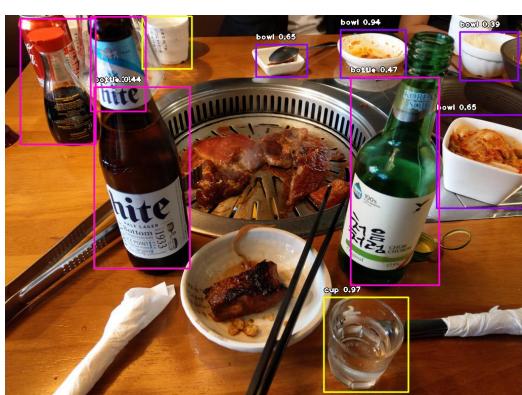


Figure 4.19: TinyYolo Detections.

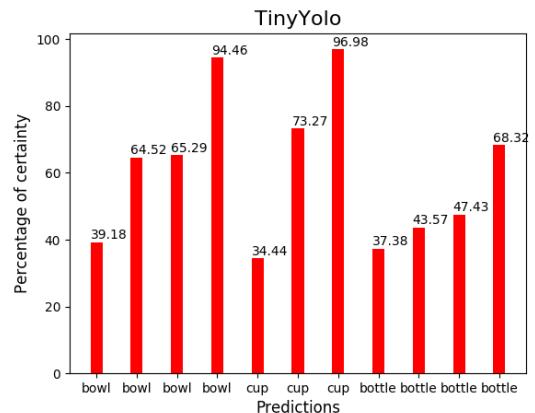


Figure 4.20: TinyYolo Detections.

#### 4.2.3 | Test Run Number 3



Figure 4.21: Third picture to be analysed.

##### 4.2.3.1 | RetinaNet Results

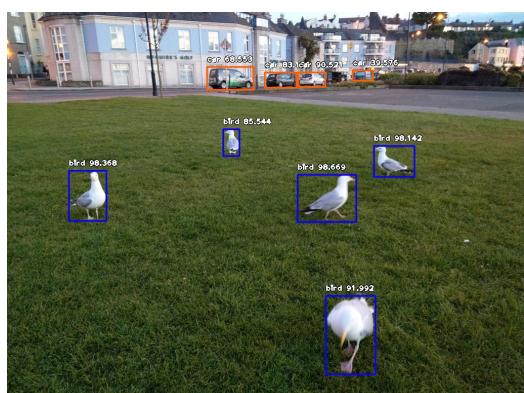


Figure 4.22: RetinaNet Detections.

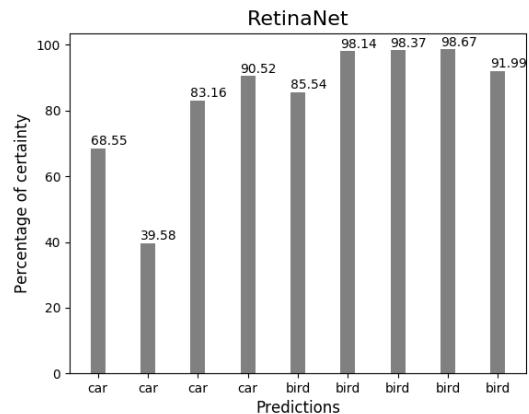


Figure 4.23: RetinaNet Detections.

#### 4.2.3.2 | YOLO Results

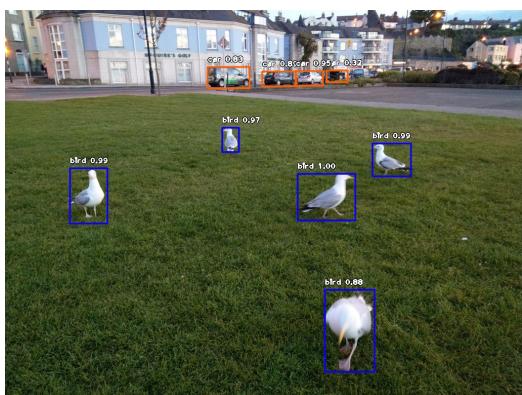


Figure 4.24: YOLO Detections.

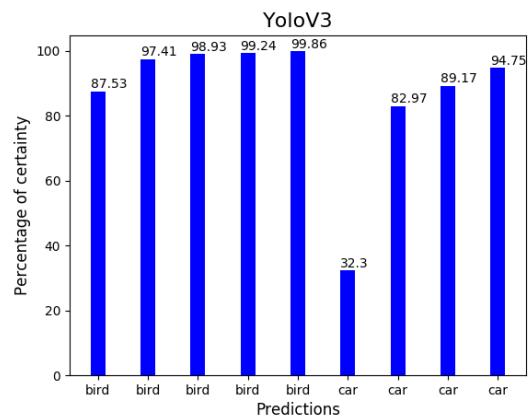


Figure 4.25: YOLO Detections.

#### 4.2.3.3 | TinyYOLO Results

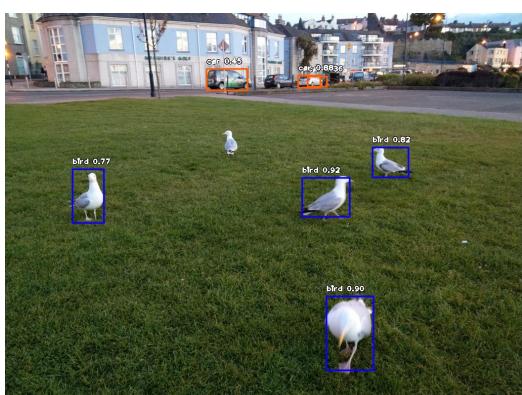


Figure 4.26: TinyYolo Detections.

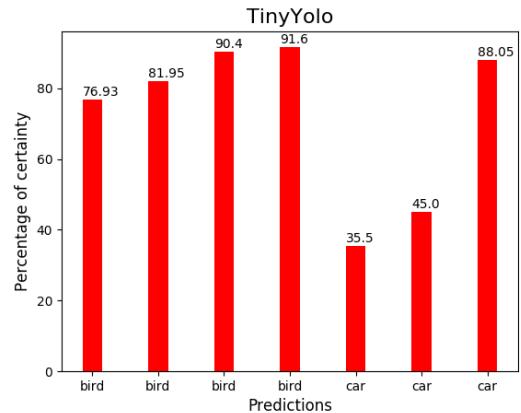


Figure 4.27: TinyYolo Detections.

#### **4.2.4 | Results analysis**

# CHAPTER 5

---

## ImageCLEF

---

### 5.1 | imageCLEF challenge

ImageCLEF launched in 2013 and was proposed by Mark Sanderson and Paul Clough from the Department of Information Studies from the University of Sheffield, with the goal of providing support for the evaluation of 1) language-independent methods for the automatic annotation of images with concepts, 2) multimodal information retrieval methods based on the combination of visual and textual features, and 3) multilingual image retrieval methods, so as to compare the effect of retrieval of image annotations and query formulations in several languages.

Currently, ImageCLEF main goal is to support the advancement of the field of visual media analysis, indexing, classification, and retrieval, by developing the necessary infrastructure for the evaluation of visual information retrieval systems operating in both monolingual, cross-language and language-independent contexts. ImageCLEF aims at providing reusable resources for such benchmarking purposes.

### 5.2 | Tasks

#### 5.2.1 | LMRT sub-task



---

## Bibliography

---

- [1] Jin Zhang. *The Information Retrieval Series*. 2008, p. 300. ISBN: 9783540751472.
- [2] Ricardo Ferreira Ribeiro and Ieeta Deti. “Object Recognition with Convolutional Neural Networks”. In: ().
- [3] Shivang Agarwal and Jean Ogier. “Recent Advances in Object Detection in the Age of Deep Convolutional Neural Networks”. In: (2019). arXiv: [arXiv:1809.03193v2](https://arxiv.org/abs/1809.03193v2).
- [4] Aastha Tiwari, Anil Kumar Goswami, and Mansi Saraswat. “Feature Extraction for Object Recognition and Image Classification”. In: *International Journal of Engineering Research & Technology (IJERT)* 2.10 (2013), pp. 1238–1246.
- [5] *The Computer Vision Pipeline, Part 4: feature extraction / Manning*. URL: <https://freecodecamp.manning.com/the-computer-vision-pipeline-part-4-feature-extraction/> (visited on 02/18/2020).
- [6] MathWorks. *What Is Artificial Intelligence? / KurzweilAI*. URL: <https://www.mathworks.com/discovery/artificial-intelligence.html%20http://www.kurzweilai.net/what-is-artificial-intelligence> (visited on 03/06/2020).
- [7] MathWorks. *What Is a Neural Network? - MATLAB & Simulink*. URL: <https://www.mathworks.com/discovery/neural-network.html> (visited on 03/06/2020).
- [8] MathWorks. *What Is Deep Learning? / How It Works, Techniques & Applications - MATLAB & Simulink*. 2019. URL: <https://www.mathworks.com/discovery/deep-learning.html> (visited on 03/05/2020).
- [9] Adrien Kaiser. *What is Computer Vision? / Hayo*. URL: <https://hayo.io/computer-vision/> (visited on 01/22/2020).
- [10] Jason Brownlee. *A Gentle Introduction to Computer Vision*. URL: <https://machinelearningmastery.com/what-is-computer-vision/> (visited on 01/22/2020).
- [11] Ilija Mihajlovic. *Everything You Ever Wanted To Know About Computer Vision*. URL: <https://towardsdatascience.com/everything-you-ever-wanted-to-know-about-computer-vision-heres-a-look-why-it-s-so-awesome-e8a58dfb641e> (visited on 01/22/2020).
- [12] Xin Feng et al. “Computer vision algorithms and hardware implementations: A survey”. In: *Integration* 69.August (2019), pp. 309–320. ISSN: 01679260.

- [13] Limarc Ambalina. *What is Image Annotation? – An Intro to 5 Image Annotation Services - By Limarc Ambalina*. URL: <https://hackernoon.com/what-is-image-annotation-an-intro-to-5-image-annotation-services-yt6n3xfj> (visited on 01/22/2020).
- [14] Jason Brownlee. *A Gentle Introduction to Object Recognition With Deep Learning*. URL: <https://machinelearningmastery.com/object-recognition-with-deep-learning/> (visited on 01/22/2020).
- [15] Thomas S. Huang. “Can the world-wide web bridge the semantic gap?” In: *Image and Vision Computing* 30.8 (2012), pp. 463–464. ISSN: 02628856.
- [16] Tsung Yi Lin et al. “Microsoft COCO: Common objects in context”. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 8693 LNCS.PART 5 (2014), pp. 740–755. ISSN: 16113349. arXiv: 1405.0312.
- [17] Y. Takamitsu and Y. Orita. “Effect of glomerular change on the electrolyte reabsorption of the renal tubule in glomerulonephritis (author’s transl)”. In: *Japanese Journal of Nephrology* 20.11 (1978), pp. 1221–1227. ISSN: 03852385.
- [18] Connecting Language et al. “Visual Genome - connected language and Vision using crowdsourced dense image annotations”. In: (2015).
- [19] Alina Kuznetsova et al. “The Open Images Dataset V4: Unified image classification, object detection, and visual relationship detection at scale”. In: (2018), pp. 1–20. arXiv: 1811.00982.
- [20] Mark Everingham et al. “The pascal visual object classes (VOC) challenge”. In: *International Journal of Computer Vision* 88.2 (2010), pp. 303–338. ISSN: 09205691.
- [21] Ivan Culjak et al. “A brief introduction to OpenCV”. In: *MIPRO 2012 - 35th International Convention on Information and Communication Technology, Electronics and Microelectronics - Proceedings* (2012), pp. 1725–1730.
- [22] OpenCV Team. *About*. URL: <https://opencv.org/about/> (visited on 01/23/2020).
- [23] Martín Abadi et al. *TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems*. Tech. rep.
- [24] John D. Dignam et al. “Eukaryotic gene transcription with purified components”. In: *Methods in Enzymology* 101.C (1983), pp. 582–598. ISSN: 15577988.
- [25] A. Vedaldi and B. Fulkerson. *VFeat: An Open and Portable Library of Computer Vision Algorithms*. 2008. (Visited on 01/23/2020).
- [26] BoofCV Team. *BoofCV*. URL: <https://boofcv.org/index.php?title=Main%7B%5C-%7DPage> (visited on 01/23/2020).

- [27] Jian Guo et al. “GluonCV and GluonNLP: Deep Learning in Computer Vision and Natural Language Processing”. In: (2019), pp. 1–6. arXiv: 1907.04433.
- [28] Armaan Merchant. *Neural Networks Explained – Data Driven Investor – Medium*. 2018. URL: <https://medium.com/datadriveninvestor/neural-networks-explained-6e21c70d7818> (visited on 03/04/2020).
- [29] Kjell Magne Fauske. *What Is Deep Learning? / How It Works, Techniques & Applications - MATLAB & Simulink*. 2019. URL: <https://www.mathworks.com/discovery/deep-learning.html> (visited on 03/05/2020).
- [30] Keiron O’Shea and Ryan Nash. “An Introduction to Convolutional Neural Networks”. In: November (2015). arXiv: 1511.08458.
- [31] Forrest N. Iandola et al. “SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5MB model size”. In: (2016), pp. 1–13. arXiv: 1602.07360.
- [32] Sik-Ho Tsang. *Review: SqueezeNet (Image Classification) - Towards Data Science*. URL: <https://towardsdatascience.com/review-squeezezenet-image-classification-e7414825581a> (visited on 01/23/2020).
- [33] Kaiming He et al. “Deep residual learning for image recognition”. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition* 2016-December (2016), pp. 770–778. ISSN: 10636919. arXiv: 1512.03385.
- [34] Raimi Karim. *Illustrated: 10 CNN Architectures*. URL: <https://towardsdatascience.com/illustrated-10-cnn-architectures-95d78ace614d> (visited on 03/12/2020).
- [35] Bharath Raj. *A Simple Guide to the Versions of the Inception Network*. URL: <https://towardsdatascience.com/a-simple-guide-to-the-versions-of-the-inception-network-7fc52b863202> (visited on 01/23/2020).
- [36] Christian Szegedy et al. “Rethinking the Inception Architecture for Computer Vision”. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition* 2016-December (2016), pp. 2818–2826. ISSN: 10636919. arXiv: 1512.00567.
- [37] Sik-Ho Tsang. *Review: Inception-v3 — 1st Runner Up (Image Classification) in ILSVRC 2015*. URL: <https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c> (visited on 01/23/2020).
- [38] Christian Szegedy et al. *Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning*. 2016. arXiv: 1602.07261 [cs.CV].
- [39] Tsung Yi Lin et al. “Focal Loss for Dense Object Detection”. In: *Proceedings of the IEEE International Conference on Computer Vision* 2017-October (2017), pp. 2999–3007. ISSN: 15505499. arXiv: 1708.02002.

- [40] Tsung-Yi Lin et al. *Feature Pyramid Networks for Object Detection*. 2016. arXiv: 1612.03144 [cs.CV].
- [41] Joseph Redmon and Ali Farhadi. “YOLOv3: An Incremental Improvement”. In: (2018). arXiv: 1804.02767.
- [42] Lilian Weng. “Object Detection Part 4: Fast Detection Models”. In: *lilianweng.github.io/lil-log* (2018).
- [43] Zhang Yi, Shen Yongliang, and Zhang Jun. “An improved tiny-yolov3 pedestrian detection algorithm”. In: *Optik* 183.January (2019), pp. 17–23. ISSN: 00304026.
- [44] Karen Simonyan and Andrew Zisserman. *Very Deep Convolutional Networks for Large-Scale Image Recognition*. 2014. arXiv: 1409.1556 [cs.CV].
- [45] Wei Liu et al. “SSD Net”. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 9905 LNCS (2016), pp. 21–37. ISSN: 16113349. arXiv: [arXiv:1512.02325v5](https://arxiv.org/abs/1512.02325v5).
- [46] Lilian Weng. “Object Detection for Dummies Part 3: R-CNN Family”. In: *lilianweng.github.io/lil-log* (2017).
- [47] Lilian Weng. “Object Detection for Dummies Part 1: Gradient Vector, HOG, and SS”. In: *lilianweng.github.io/lil-log* (2017).
- [48] Ross Girshick. *Fast R-CNN*. 2015. arXiv: 1504.08083 [cs.CV].
- [49] Shaoqing Ren et al. *Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks*. 2015. arXiv: 1506.01497 [cs.CV].
- [50] Kaiming He et al. *Mask R-CNN*. 2017. arXiv: 1703.06870 [cs.CV].
- [51] Paperswithcode. *COCO test-dev Leaderboard / Papers with Code*. URL: <https://paperswithcode.com/sota/object-detection-on-coco> (visited on 03/06/2020).
- [52] Paperswithcode. *ImageNet Leaderboard / Papers with Code*. URL: <https://paperswithcode.com/sota/image-classification-on-imagenet> (visited on 03/06/2020).
- [53] Yudong Liu et al. “CBNet: A Novel Composite Backbone Network Architecture for Object Detection”. In: (2019). arXiv: 1909.03625.
- [54] Mingxing Tan, Ruoming Pang, and Quoc V. Le. “EfficientDet: Scalable and Efficient Object Detection”. In: (2019). arXiv: 1911.09070.
- [55] Shifeng Zhang et al. “Bridging the Gap Between Anchor-based and Anchor-free Detection via Adaptive Training Sample Selection”. In: 2 (2019). arXiv: 1912.02424.

- [56] Ross Girshick et al. *Detectron*. 2018. URL: <https://github.com/facebookresearch/detectron> (visited on 03/09/2020).
- [57] Yanghao Li et al. “Scale-Aware Trident Networks for Object Detection”. In: (2019). arXiv: 1901.01892.
- [58] Dhruv Mahajan et al. “Exploring the Limits of Weakly Supervised Pretraining”. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 11206 LNCS (2018), pp. 185–201. ISSN: 16113349. arXiv: 1805.00932.
- [59] Qijie Zhao et al. “M2Det: A Single-Shot Object Detector Based on Multi-Level Feature Pyramid Network”. In: *Proceedings of the AAAI Conference on Artificial Intelligence* 33 (2019), pp. 9259–9266. ISSN: 2159-5399. arXiv: 1811.04533.
- [60] Zhaowei Cai and Nuno Vasconcelos. “Cascade R-CNN: Delving into High Quality Object Detection”. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition* (2018), pp. 6154–6162. ISSN: 10636919. arXiv: 1712.00726.
- [61] Jiaqi Wang et al. “Region proposal by guided anchoring”. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition* 2019-June (2019), pp. 2960–2969. ISSN: 10636919. arXiv: 1901.03278.
- [62] Abhinav Shrivastava et al. *Beyond Skip Connections: Top-Down Modulation for Object Detection*. 2016. arXiv: 1612.06851 [cs.CV].
- [63] Seung Wook Kim et al. “Parallel Feature Pyramid Network for Object Detection”. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 11209 LNCS (2018), pp. 239–256. ISSN: 16113349.
- [64] Saining Xie et al. “Aggregated residual transformations for deep neural networks”. In: *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017* 2017-January (2017), pp. 5987–5995. arXiv: 1611.05431.
- [65] Olga Russakovsky et al. “ImageNet Large Scale Visual Recognition Challenge”. In: *International Journal of Computer Vision* 115.3 (2015), pp. 211–252. ISSN: 15731405. arXiv: 1409.0575.
- [66] Qizhe Xie et al. “Self-training with Noisy Student improves ImageNet classification”. In: (2019). arXiv: 1911.04252.
- [67] Alexander Kolesnikov et al. *Large Scale Learning of General Visual Representations for Transfer*. 2019. arXiv: 1912.11370 [cs.CV].
- [68] Hugo Touvron et al. *Fixing the train-test resolution discrepancy*. 2019. arXiv: 1906.06423 [cs.CV].

- [69] Cihang Xie et al. *Adversarial Examples Improve Image Recognition*. 2019. arXiv: 1911.09665 [cs.CV].
- [70] Mingxing Tan and Quoc V. Le. *EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks*. 2019. arXiv: 1905.11946 [cs.LG].
- [71] Mark Sandler et al. *MobileNetV2: Inverted Residuals and Linear Bottlenecks*. 2018. arXiv: 1801.04381 [cs.CV].
- [72] Mingxing Tan et al. *MnasNet: Platform-Aware Neural Architecture Search for Mobile*. 2018. arXiv: 1807.11626 [cs.CV].
- [73] Rite Wiki. *Natural Language Processing - Ryte Wiki - The Digital Marketing Wiki*. URL: [https://en.ryte.com/wiki/Natural\\_Language\\_Processing](https://en.ryte.com/wiki/Natural_Language_Processing) (visited on 08/01/2020).
- [74] Diksha Khurana et al. “Natural Language Processing : State of The Art , Current Trends and Challenges Natural Language Processing : State of The Art , Current Trends and Challenges Department of Computer Science and Engineering Manav Rachna International University , Faridabad-”. In: *arXiv preprint arXiv August 2017* (2018).
- [75] Daniel W. Otter, Julian R. Medina, and Jugal K. Kalita. “A Survey of the Usages of Deep Learning in Natural Language Processing”. In: XX.X (2018), pp. 1–22. arXiv: 1807.10854.
- [76] Analytics Vidhya. *Understanding Word Embeddings: From Word2Vec to Count Vectors*. 2017. URL: <https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/> (visited on 02/06/2020).
- [77] Murat Mustafa. *GloVe / Mustafa Murat ARAT*. URL: <https://mmuratarat.github.io/2020-03-20/glove> (visited on 05/13/2020).
- [78] Jason Brownlee. *What are word embeddings for text?* 2017. URL: <https://machinelearningmastery.com/what-are-word-embeddings/> (visited on 02/06/2020).
- [79] Tomas Mikolov et al. “Efficient estimation of word representations in vector space”. In: *1st International Conference on Learning Representations, ICLR 2013 - Workshop Track Proceedings* (2013), pp. 1–12. arXiv: 1301.3781.
- [80] David Batista. */oorsig/ Language Models and Contextualised Word Embeddings*. 2018. URL: [http://www.davidsbatista.net/blog/2018/12/06/Word%7B%5C\\_%7DEmbeddings/%20http://www.davidsbatista.net//blog/2018/12/06/Word%7B%5C\\_%7DEmbeddings/](http://www.davidsbatista.net/blog/2018/12/06/Word%7B%5C_%7DEmbeddings/%20http://www.davidsbatista.net//blog/2018/12/06/Word%7B%5C_%7DEmbeddings/) (visited on 02/10/2020).
- [81] *Word2Vec Explained Easily - InsightsBot*. URL: <http://www.insightsbot.com/word2vec-explained-easily/> (visited on 02/08/2020).

- [82] A.I. Wiki. *A Beginner’s Guide to Word2Vec and Neural Word Embeddings / Skymind*. URL: <https://pathmind.com/wiki/word2vec> % 20<https://skymind.ai/wiki/word2vec> (visited on 02/08/2020).
- [83] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. “GloVe: Global vectors for word representation”. In: *EMNLP 2014 - 2014 Conference on Empirical Methods in Natural Language Processing, Proceedings of the Conference*. Association for Computational Linguistics (ACL), 2014, pp. 1532–1543. ISBN: 9781937284961.
- [84] Piotr Bojanowski et al. “Enriching Word Vectors with Subword Information”. In: *arXiv preprint arXiv:1607.04606* (2016).
- [85] Oren Melamud, Jacob Goldberger, and Ido Dagan. “context2vec: Learning generic context embedding with bidirectional LSTM”. In: *CoNLL 2016 - 20th SIGNLL Conference on Computational Natural Language Learning, Proceedings* (2 016), pp. 51–61.
- [86] Matthew E. Peters et al. “Deep contextualized word representations”. In: *Proc. of NAACL*. 2018.
- [87] Spacy. *spaCy 101: Everything you need to know · spaCy Usage Documentation*. 2017. URL: <https://spacy.io/usage/spacy-101> (visited on 02/10/2020).
- [88] Edward Loper and Steven Bird. “Nltk”. In: March (2002), pp. 63–70.
- [89] Christopher Manning et al. “The Stanford CoreNLP Natural Language Processing Toolkit”. In: (2015), pp. 55–60.
- [90] Radim Řehůřek and Petr Sojka. “Software Framework for Topic Modelling with Large Corpora”. English. In: *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. <http://is.muni.cz/publication/884893/en>. Valletta, Malta: ELRA, May 2010, pp. 45–50.
- [91] Alan Akbik, Duncan Blythe, and Roland Vollgraf. “Contextual String Embeddings for Sequence Labeling”. In: *COLING 2018, 27th International Conference on Computational Linguistics*. 2018, pp. 1638–1649.
- [92] Rami Al-Rfou, Bryan Perozzi, and Steven Skiena. “Polyglot: Distributed Word Representations for Multilingual NLP”. In: *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*. Sofia, Bulgaria: Association for Computational Linguistics, Aug. 2013, pp. 183–192.
- [93] “CogCompNLP: Your Swiss Army Knife for NLP”. In: *11th Language Resources and Evaluation Conference*. 2018.
- [94] *TextBlob: Simplified Text Processing — TextBlob 0.16.0 documentation*. URL: <https://textblob.readthedocs.io/en/dev/> (visited on 07/28/2020).

- [95] Moses and John Olafenwa. *ImageAI, an open source python library built to empower developers to build applications and systems with self-contained Computer Vision capabilities*. Mar. 2018–. URL: <https://github.com/OlafenwaMoses/ImageAI>.
- [96] MathWorks. *What Is a Machine Learning? - MATLAB & Simulink*. URL: <https://www.mathworks.com/discovery/machine-learning.html> (visited on 03/12/2020).
- [97] Vikas Gupta. *Understanding Feedforward Neural Networks*. 2017. URL: <https://www.learnopencv.com/understanding-feedforward-neural-networks/>.
- [98] Ngoc Giang Nguyen et al. “DNA Sequence Classification by Convolutional Neural Network”. In: *Journal of Biomedical Science and Engineering* 09.05 (2016), pp. 280–286. ISSN: 1937-6871.