

1 2 9 0



UNIVERSIDADE D  
COIMBRA

David Alexandre Mendes Carreira

IMAGE RECOGNITION FOR ONLINE  
STUDENT IDENTIFICATION

Thesis submitted to the University of Coimbra in fulfilment of the requirements of the Master's Degree in Engineering Physics under the scientific supervision of PhD David Portugal and MsC José Faria and presented to the Physics Department of the Faculty of Sciences and Technology of the University of Coimbra.

June 2023



# Contents

Acknowledgments . . . . .	i
Abstract . . . . .	ii
Resumo . . . . .	iii
<b>1 Introduction</b>	<b>1</b>
1.1 Context . . . . .	1
1.2 Dissertation structure . . . . .	2
<b>2 State of The Art</b>	<b>4</b>
2.1 History of AI . . . . .	4
2.2 Face Recognition . . . . .	9
2.3 A Face Recognition System . . . . .	10
2.3.1 Face Detection . . . . .	11
2.3.2 Face Alignment . . . . .	15
2.3.3 Face Representation . . . . .	15
2.4 Face Representation Pipeline . . . . .	15
2.4.1 Convolutional Neural Networks . . . . .	15
<b>3 Methodology</b>	<b>18</b>
3.1 Joint Face Detection and Alignment . . . . .	19
3.2 Face Representation . . . . .	19
<b>4 Results</b>	<b>20</b>



## Acknowledgments

# Abstract

## Resumo

# Chapter 1

## Introduction

### 1.1 Context

The outbreak of the COVID-19 pandemic tested the entire world on several levels and changed the concept of what is "normal" thereafter. The devastating health, economic and social consequences that COVID caused, spanned a need to develop novel solutions, for almost every aspect of our lives, that facilitate the adaptation to the new world we're living in.

Educational systems were no exception. In the midst of the pandemic, governments around the world forced institutions to shut down and stop the customary in-person regimen of teaching. By April 2020, most universities transitioned to an adapted remote learning [70] that lacked proper support due to the unanticipated nature of the events, leading to new challenges, in particular, the legitimacy of moments of evaluation performed remotely. To counter this problem, different approaches can be taken, namely, changing the method of evaluation, suppressing it altogether [4] or, when possible, implement a continuous monitoring solution such as TrustID [ref?](#). However, there are still unresolved issues that must be addressed in order to implement an end-to-end solution capable of assuring the success of such systems.



One core aspect of them is the face verification task, therefore, the data obtained directly influences the performance. Due to the purpose of the application and expected devices to be used, what is obtained can be classified as from an unconstrained nature. Even though the capture of image is consensual, there is no way of controlling the conditions of capturing the visual data and consequent results. This can be attributed to the fact that it is anticipated that the system will be executed in a laptop or a smartphone, thus the capture device might not be ideal. The more probable input method will be a webcam or the smartphone's front facing camera, so a high variation in pose, resolution, illumination, etc. is not unforeseeable.

Another detail that must be regarded, is the processing power available to execute the system. It is common for the equipments used to have a deficiency of it<sup>1</sup>, which is not suitable for high-demanding applications, as its improved accuracy comes at the cost of increased computational overhead, which can make real-time continuous monitoring unfeasible.

In conclusion, the method of choice must take the aforesaid into consideration and be a trade-off between accuracy and computational strain, while also being invariant, to a certain degree, to the posed challenges of capturing the required data.

## 1.2 Dissertation structure

This dissertation will be divided into different chapters that partitions themselves into sections and subsections. Chapter one relates to the introduction of the dissertation, it will present the context and motivation behind the problem and structure of the document. The document continues to the second chapter, it starts with an overview of the History of AI, carries out a survey about the topic's

---

<sup>1</sup> According to the February 2023 Steam hardware survey, roughly 5% of its users do not have a dedicated GPU.

State-Of-The-Art, presented and summarized through the step-by-step analysis of the pipeline of a Face Recognition system, and ends with a comparison table of the discussed methods. In chapter number three, the implemented methods and experiments are described. The forth chapter will present and discuss the results. Finally, chapter five, will draw conclusions of the work achieved in the past several months and prospects for the future.

# Chapter 2

## State of The Art

### 2.1 History of AI

The following sections present a broad overview of the history of Artificial Intelligence (AI) by presenting important articles in order for the reader to be able to have a notion of the progress that has been made over the past decades, the hardships encountered and how important AI is in our lives.

#### **Philosophy**

On October 1950, in his article *Computing Machinery and Intelligence*, Alan Turing questioned: "Can machines think?" [62]. At the time, the question was too meaningless to answer since not only the theory but also the technology available weren't developed enough. Nonetheless, Turing still predicted that in the future there would be computers that could, effectively, display human-like intelligence and discernment under the conditions proposed on the aforementioned article.

#### **Relevant events to the birth of AI**

The breakthroughs of AI are predominant, and its importance in our everyday life is undeniable, but the theory behind it has several early roots. The interest in

the area grew immensely with, for example, all the Turing's theoretical research, the proposal of the first mathematical Artificial Neuron model in 1943 by Warren McCulloch and Walter Pitts (based of binary inputs and output) [43] and in 1949 Donald Hebb revolutionized the way the artificial neurons were treated by proposing what is known as the Hebb's rule<sup>1</sup>. Taking into consideration the latter two, but specially Hebb's proposals, Belmont Farley and Westley Clark implemented in 1954 one of the first successful Artificial Neural Networks (ANN), also called Perceptron, composed of two layers of 128 artificial neurons with weighted inputs [18]. Over the span of approximately ten years, multiple researches were performed attempting to computerize the human brain. However, only in 1956, during the *Dartmouth Summer Research Project on Artificial Intelligence* [42], was the term "Artificial Intelligence" firstly proposed by John McCarthy *et al.*, beginning what is now considered to be the birth of AI [78].

### **The fading of general interest**

The succeeding two decades following the Dartmouth conference were filled with important developments, with special emphasis in the works published in 1958 by Frank Rosenblatt (generalized the Farley and Clark training to multi-layer networks rather than only two) [53], the 1959 General Problem Solver implemented by Allen Newel *et al.* (a program intended to work as a universal problem solver that was capable of solving exercises such as the Towers of Hanoi<sup>2</sup>) [46] and the ELIZA a natural language processing tool program developed by Joseph Weizenbaum between 1964 and 1966 [68]. Unfortunately, part of the interest and development around AI met an unforeseen fade after criticisms about the exagger-

---

<sup>1</sup> "When an axon of cell A is near enough to excite a cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells such that A's efficiency, as one of the cells firing B, is increased." [22], meaning that when two neurons fire together their relation is strengthened.

<sup>2</sup> *The Towers of Hanoi* is a game with 3 stacks of increasingly smaller disks. The goal is to stack them one at a time, so that they are arranged in a decreasing radius manner.

ated public funding [21] and the Marvin Minsky and Seymour Papert 1969 book *The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain* [45] that reported on the problems of the Perceptron network. The overall sentiments regarding this topic of research was of doubt and fear of no progress, mainly due to the spending and two issues raised by Minsky and Papert: the ANN couldn't solve linear inseparable problems<sup>3</sup> and there were limitations due to a lack of sufficient computing power to handle the processing of multi-layer large networks.

## A better approach

Minsky and Papert raised important questions, but it shouldn't have discouraged other researchers from further trying, since they failed to acknowledge alternative approaches that had already solved those exact problems. As previously stated, the model proposed by McCulloch and Pitts, later improved by the Farley-Clark implementation and, finally, Rosenblatt, couldn't handle linearly inseparable classes. A possible solution for cases like this started being studied in the 1960s [29, 54] and, although it didn't produce relevant results, in 1965 Alexey Ivakhnenko and Valentin Lapa [27] were, indeed, successful in implementing what is nowadays considered to be the first deep learning network of its kind [56]. In 1971 Ivakhnenko also published an article describing a deep learning network with 8 layers that was already able to create hierarchical internal representations [28].

The years progressed, in 1979 Kuniyuki Fukushima introduced the first Convolutional Neural Network (CNN) in a structural sense, due to its similarity to the architecture of modern ones of this category. Ten years later, Yann LeCun *et al.* applied for the first time a revolutionizing training algorithm called Backpropagation to a CNN [33], creating what is now a pillar for most of the modern competition winning networks in computer vision [56] and employing the term "convolution"

---

<sup>3</sup> That is, if two sets  $X$  and  $Y$  in  $\mathbb{R}^d$  can't be divided by a hyperplane such that the elements of  $X$  and  $Y$  stay on opposing sides, then we're dealing with linear inseparable classes [17]

for the first known time [37]. He also introduced the MNIST (**M**odified **N**ational **I**nstitute of **S**tandards and **T**echnology) dataset, a collection of handwritten digits [35], that to this day is still one of the most famous benchmarks in Machine Learning. Backpropagation can be traced back many decades, but the modern version was first described by Seppo Linnainmaa (1970) [38], implemented for the first time by Stuart Dreyfus (1973) [15] and, finally in 1986, David Rumelhart *et al.* popularized it in the Neural Network's (NN) domain by demonstrating the growing usefulness of internal representations [55].

## The importance of Convolutional Neural Networks

The study on Neural Networks continued and there were improvements on all types of architectures [23, 69] with special highlight to pioneering Neural Networks processed by GPUs<sup>4</sup> (standard NN in 2004 by [47] and CNN in 2006 by [7]). But there's a well deserved particular attention related to the developments of CNNs due to their great performance in image related tasks when compared to others networks, as proven by LeCun in his 1998 paper [35]. Some relevant examples: in 2003 the MNIST record was broken by Patrice Simard *et al.* [57], achieving an error rate of 0.4% (whereas a non-convolutional neural network by the same authors took the second place with 0.7%); three years later, the same benchmark had a new set low of 0.39% by Marc'Aurelio Ranzato *et al.* [51]; in 2009 a CNN by Yang *et al.* was the first network of this type to win an official international competition (TRECVID) [74]; a GPU implementation of a CNN [10] achieved superhuman vision performance in a competition (IJCNN 2011) in a *German Traffic Sign Recognition Benchmark* with a 0.56% error rate (0.78% for the best human performance, 1.69% for the second-best neural network contestant and 3.86% for the best non-neural method [58]). This last example conjoined with non-convolutional methods [49, 12] and the previously cited [7, 47], reinforces how fundamental GPUs were to further develop neural networks. To supplement even

---

<sup>4</sup> Graphics Processing Unit

more the importance of CNNs and GPUs, only a year later, Alex Krizhevsky *et al.* proposed a Deep CNN trained by GPUs that was the first one to win the ImageNet Large Scale Visual Recognition Challenge (ILSVRC), achieving an error rate of 15.3% while the second place obtained 26.2% [31].

The year of 2012 was very important for Deep Learning, CNNs and Computer Vision, due to all the attention brought to many researches on this topic after several systems of this kind won image analysis competitions ([9, 11] and the very important previously mentioned [31]), beginning what's considered to be the start of the new wave, we're currently in, of interest in Artificial Intelligence, specially in the aforesaid topics [37].

## 2.2 Face Recognition

Face Recognition (FR) is a thoroughly debated and extensively researched task in the Computer Vision community for more than two decades [50], popularized in the early 1990s with the introduction of the Eigenfaces [63] or Fisherfaces [48] approaches. These methods projected faces in a low-dimensional subspace assuming certain distributions, but lacked the ability to handle uncontrolled facial changes that broke said assumptions, henceforth, bringing about face recognition approaches through local-features [8, 2] that, even though, presented considerable results, weren't distinctive or compact. Beginning in 2010, methods based on learnable filters arose [75, 36], but unfortunately revealed limitations when nonlinear variations were at stake.

Earlier methods for FR worked appropriately when the data was handpicked or generated on a constrained environment, however, they didn't scale adequately in the real world where there are large fluctuations in, particularly, pose, age, illumination, background scenario, the presence of facial occlusion [50] and many unimaginable more. These shortcomings can be dealt with by using Deep Learning, a framework of techniques that solves the nonlinear inseparable classes problem [ref.](#), more specifically a structure called Convolutional Neural Network (CNN) [65].

CNNs are an Artificial Neural Network (ANN) that exhibit a better performance on image or video-based tasks compared to other methods [35]. They were greatly hailed in 2012, after the AlexNet [31] victory, by a great margin, in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC). Just two years later, DeepFace [59] revolutionized the benchmarks scores by achieving state-of-the-art results that approached human performance, reinforcing even further the importance of Deep Learning and shifting the research path to be taken [65].

Given what has been stated so far and the proven robustness, performance, and overall results in computer vision [ref. won competitions](#), the methods discussed in this dissertation will therefore deal exclusively with Deep Learning approaches.



For more information on other methods, please refer to [32].

## 2.3 A Face Recognition System

According to Ranjan *et al.* [50], the goal of a FR system is to find, process and learn from a face, gathering as much information as possible, and as a result, it is one of the most widely implemented biometric system solutions, in light of its versatility when facing real world application [16], such as **military, public security and daily life**.

By and large, all end-to-end automatic face recognition systems follow a sequential and modular<sup>5</sup> pipeline (Figure 1) composed of three pillar stages [65]: face detection, face alignment and face representation. First an image or video feed is used as an input then, as the name suggests, the **face detection** module is responsible for finding a face. Next, the **face alignment** phase applies transformations to the data, such as crop and/or rotation, in order to normalize the faces' pictures (or frames, in the case where a video is used) to standardized coordinates. Finally, the **face representation** stage, makes use of deep learning techniques to learn discriminative features that will allow the recognition.

All three stages have their individual importance and methods of implementation<sup>6</sup>. Face detection is achievable through classical approaches [64, 5] or deep methods, among them is [14] and the widely applied [81]. Face alignment, once again, can be accomplished through traditional measures [13, 41] or more modern ones, namely [26] or the aforementioned [81] which concurrently performs detec-

---

<sup>5</sup> Sequential because each stage relies on the output from the previous ones, and modular in the sense that each stage employs its own method and it can be modified to better adapt to specific tasks.

<sup>6</sup> For a deeper and extensive study, please refer to: [76] in the case of classic face detection approaches and [44] for deep learning based methods; [67] addresses traditional face alignment methods and is complemented with [16] for more up-to-date techniques; and [32] tackles classic face representation **(add the following if needed) while X supplements the deep learning ones**

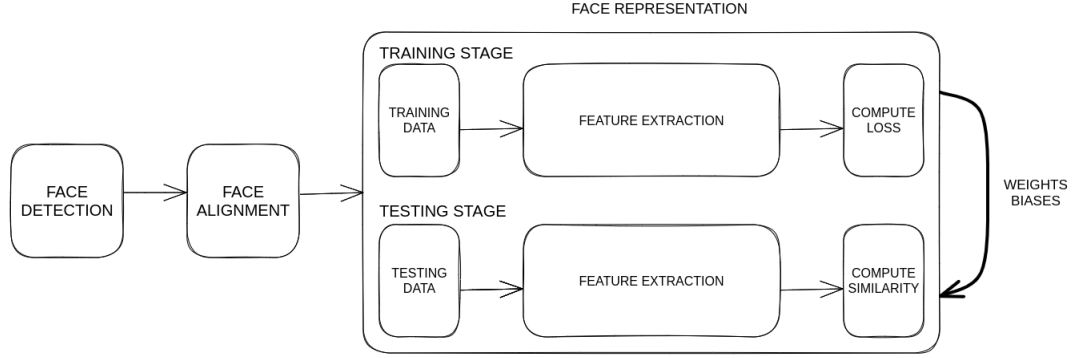


Figure 1: A typical face recognition pipeline, guided by the approach in [65].

tion and alignment. To conclude, the face representation module is no exception, and can also be divided in two groups, regarding the methodology used. Some conventional systems were already mentioned, such as [48, 63], and the deep learning ones are the objective of discussion of this dissertation and will be reviewed along the following sections, therefore, the focus will be on describing, with particular interest, the face representation stage.

### 2.3.1 Face Detection

Face detection is the first step in any automatic facial recognition system. Given an input image to a face detector module, it is in charge of detecting every face in the picture and returning bounding-boxes coordinates, for each one, with a certain confidence score [16, 50].

Previously employed traditional face detectors [cite here](#) are incapable of detecting facial information when faced with challenges such as variants in image resolution, age, pose, illumination, race, occlusions or accessories (masks, glasses, makeup) [16, 50]. The progress in deep learning and increasing GPU power led DCNNs to become a viable and reliable option that solves said problems in face detection.

These techniques can be included in different categories. A more analytical perspective [16] distributes the methods, depending upon their architecture or pur-

pose of application, over seven categories: multi-stage, single-stage, anchor-based, anchor-free, multi-task learning, CPU real-time and, finally, problem-oriented. Additionally, being as the face detection problem can be seen as a specific task in a general object detection situation, it is no surprise that several works inherit from them and, therefore, some bases are referenced throughout the next list.

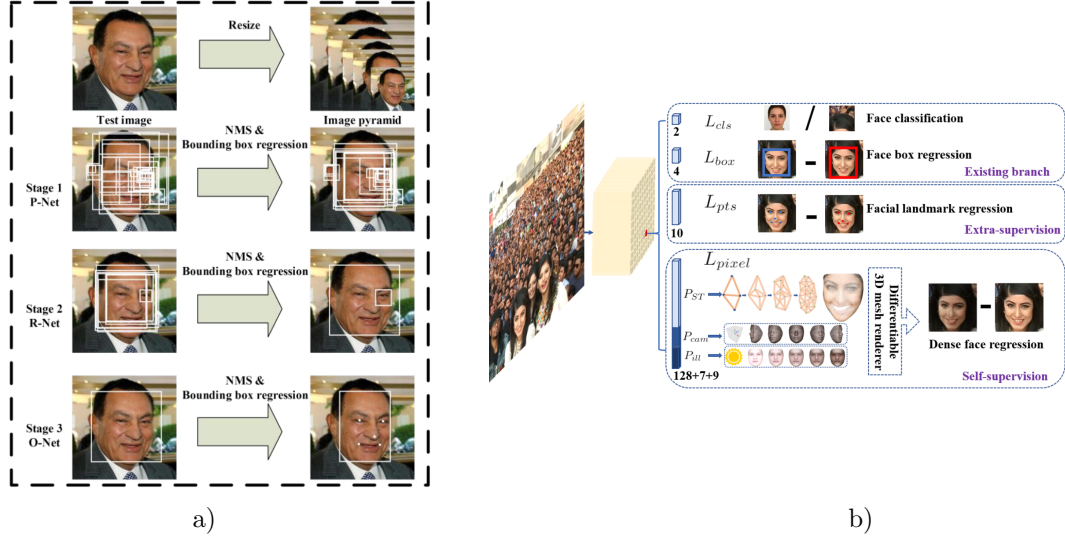


Figure 2: Comparison between a) MTCNN: multi-stage, CPU real-time and multi-task learning, and b) RetinaFace: single-stage, anchor-based, CPU real-time and multi-task learning. MTCNN [81] proposes a series of bounding boxes then, through a series of refinement stages, the best solution and landmarks are found. RetinaFace [14] accomplishes, in a single-stage, face classification and bounding box regression by evaluating anchors, landmark localization and dense 3D projection for facial correspondence.

**Multi-stage** methods [14] include all the coarse-to-fine facial detectors that work in similar manner to the following two phases. First, bounding box proposals are generated by sliding a window through the input. Then, over one or several subsequent stages, false positives are rejected and the approved bounding boxes are refined. To complement, one widely applied object detection protocol that inspired face detection methods and perfectly describes the steps mentioned above is Faster R-CNN [52]. However, these methods can be slower and have a more complex way of training [72].

**Single-stage** approaches [14] are the ones that perform classification and bounding box regression without the necessity of a proposal stage, producing highly dense face locations and scales. This structure takes inspiration, once again, from general object detectors, for example, the Single Shot MultiBox detector, commonly referred to as SSD [39]. Finally, the methods included in this class are more efficient, but can incur in compromised accuracy, when compared to multi-stage.

**Anchor-based** techniques [40, 14, 79] detect faces by predefining anchors with different settings (scales, strides, number, etc.) on the feature maps, then performing classification and bounding box regression on them until an acceptable output is found. As proven by Liu and Tang *et al.* [40], the choice of anchors highly influences the results of prediction. Hence, it is necessary to fine-tune them on a situation-by-situation basis, otherwise, there is a limitation in generalization. Furthermore, higher densities of anchors directly generate an increase in computational overhead.

**Anchor-free** procedures, obviously, do not need predefined anchors in order to find faces. Alternatively, these methods address the face detection by using different techniques. For example, DenseBox [25] which attempts to predict faces by processing each pixel as a bounding box, or CenterFace [72] that treats face detection as a key-point estimation problem by predicting the center of the face and bounding boxes. Even so, relating to the accuracy of anchor-free approaches, there's still room for improvement for false positives and stability in the training stage [16].

**Multi-task learning** are all the methodologies that conjointly performs other tasks, namely facial landmark<sup>7</sup> localization, during face classification and bound-

---

<sup>7</sup> A facial landmark is a key-point in a face that contributes with important geometric information, namely the eyes, nose, mouth, etc. [19]

ing box regression [16]. CenterFace [72] is one example, and so it is the widely implemented MTCNN [81], which correlated bounding boxes and face landmarks. RetinaFace [14] is another state-of-the-art approach, it mutually detects faces, respective landmarks and performs dense 3D face regression.

**CPU real-time** methods, as the name suggests, include the detectors that can run on a single CPU core, in real-time, for VGA-resolution input images. A face detector can achieve great results in terms of accuracy, but for real world applications, its use can be too computational heavy, therefore, can't be deployed in real time (specially in devices that do not have a GPU) [16]. MTCNN [81], Faceboxes [82], CenterFace [72] or RetinaFace [14] are examples of this category.

**Problem-oriented** is a category that includes the detectors that are projected to resolve a wide range of specific problems, for example, faces that are tiny, partially occluded, blurred or scale-invariant face detection [16]. PyramidBox [60] is an example that solves the partial occluded and blurry faces, and HR [24] tackles the tiny faces challenge.

Although this distribution can create some overlap among the categories, it is superior due to the simplicity of inferring what defines each category and being a more fine-grained way of classifying techniques when compared to others, namely the dual categorical division by [50] that groups the methods in region<sup>8</sup> or sliding-window<sup>9</sup> based.

---

<sup>8</sup> Region-based approaches creates thousands of generic object-proposals for every image, and subsequently, a DCNN classifies if a face is present in any of them.

<sup>9</sup> Sliding-window approaches centers on using a DCNN to compute a face detection score and bounding box at every location in a feature map.

### **2.3.2 Face Alignment**

Face Alignment is the second stage of the face recognition pipeline, and has the objective of calibrating the detected face to a canonical layout, through landmark-based or landmark-free approaches, in order to leverage the core final stage of face representation [16].

Despite the fact that traditional face alignment methods are very accurate, that only occurs in constrained circumstances. Therefore, once again, to address that issue, deep learning-based methods are the solution to perform an accurate facial landmark localization that realistically scales to real world scenarios [19].

This step in the face recognition process can be accomplished through standalone methods that process the detected face from the previous stage, or through joint detection and alignment methods, as previously mentioned in the multi-task learning definition.

**Landmark-based alignment**

**Landmark-free alignment**

**Joint Face Detection and Alignment**

### **2.3.3 Face Representation**

## **2.4 Face Representation Pipeline**

### **2.4.1 Convolutional Neural Networks**

There are several types of Neural Networks architectures, but Convolutional Neural Networks (CNNs or Convnets) are probably the most widely implemented model overall [73, 37] with successful applications in the domains of Computer Vision [31, 59, 61, 80] or Natural Language Processing[1, 66, 71]. In the CNN

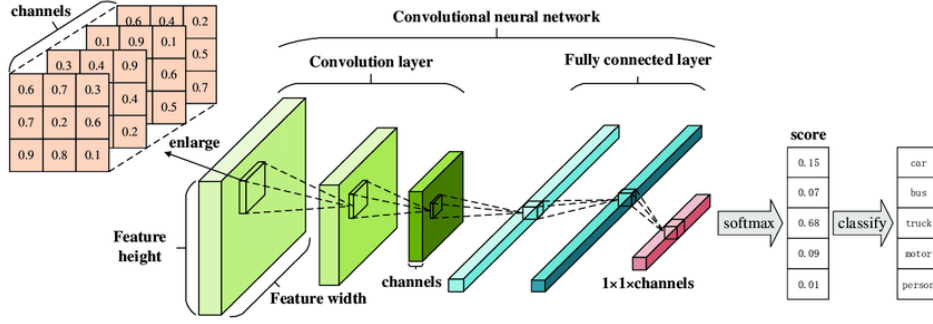


Figure 3: Architecture of a Convolutional Neural Network [30].

category itself there are different variants, but they all abide the fundamental structure of a feedforward hierarchical multi-layer network (Figure 3). Feedforward because the information only flows in a singular direction without cycling [77], hierarchical because the higher complexity internal representations are learned from lower ones [34, 83] and multi-layer because it is composed of a series of stages, blocks or layers: the raw data is fed to an input layer, forwarded to a sequence of intercalating convolutional and pooling layers, transmitted to a stage of one or more fully-connected layers [34, 73, 20, 3]. The convolutional layer is designed to extract feature representations by being composed of kernels (or filter banks [34]) that compute feature maps through element-wise product, to which is applied a nonlinear activation function [20, 73]. Next is the pooling layer, that's responsible for reducing the spatial size of the input data [20] and joining identical features [34]. Finally, the fully connected layers, and their core function is to perform high logic and generate semantic information [20].

Using CNNs for Computer Vision tasks is not an arbitrary choice, but due to the fact that the network design can benefit from the intrinsic characteristics of the input data, consequently performing really well in image related applications [34, 6]. In the first place, images have an array-like structure with numerous elements, namely, each pixel has an assigned value organized in a grid-like manner [73]. In the second place, there's an inherent correlation between local groups of values,

which creates distinguishable motifs [34]. Finally, the local values of images are invariant to location, that is, a certain composition should have the same value independently of the spatial location in the picture [34]. Therefore, the following key, unique features potentiate the previously stated efficient performance [6]:

1. Designed to process multidimensional arrays [34];
2. Shared weights between the same features in different locations;
3. Automatically identifies the relevant features without any human supervision, hence, small amounts of preprocessing [3, 37];
4. Local connections (or receptive fields/sparse connectivity) [3];
5. Pooling layers that reduces the spatial size of the input data.

The ensemble of features 2, 4 and 5 enable an invariance of the network to small shifts, distortions and rotations [20, 34], while 2, 3, 4 and 5 helps to reduce the complexity of the model, and as a result training it is easier [20, 37].



## Chapter 3

### Methodology

### **3.1 Joint Face Detection and Alignment**

### **3.2 Face Representation**

# Chapter 4

## Results

## Chapter 5

## Conclusion

# Bibliography

- [1] Ossama Abdel-Hamid et al. “Convolutional Neural Networks for Speech Recognition”. In: *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 22.10 (Oct. 2014), pp. 1533–1545. ISSN: 2329-9304. DOI: 10.1109/TASLP.2014.2339736 (cit. on p. 15).
- [2] T. Ahonen, A. Hadid, and M. Pietikäinen. “Face Description with Local Binary Patterns: Application to Face Recognition”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 28.12 (2006). Cited By :4611, pp. 2037–2041. ISSN: 0162-8828. DOI: 10.1109/TPAMI.2006.244 (cit. on p. 9).
- [3] Laith Alzubaidi et al. “Review of Deep Learning: Concepts, CNN Architectures, Challenges, Applications, Future Directions”. In: *Journal of Big Data* 8.1 (Mar. 2021), p. 53. ISSN: 2196-1115. DOI: 10.1186/s40537-021-00444-8. (Visited on 02/09/2023) (cit. on pp. 16, 17).
- [4] Maria Barron Rodriguez et al. *Remote Learning During the Global School Lockdown: Multi-Country Lessons*. World Bank, Aug. 2021. DOI: 10.1596/36141. (Visited on 03/13/2023) (cit. on p. 1).
- [5] S. Charles Brubaker et al. “On the Design of Cascades of Boosted Ensembles for Face Detection”. In: *International Journal of Computer Vision* 77.1 (May 2008), pp. 65–86. ISSN: 1573-1405. DOI: 10.1007/s11263-007-0060-1 (cit. on p. 10).

- [6] Weipeng Cao et al. “A Review on Neural Networks with Random Weights”. In: *Neurocomputing* 275 (Jan. 2018), pp. 278–287. ISSN: 0925-2312. DOI: 10.1016/j.neucom.2017.08.040. (Visited on 02/09/2023) (cit. on pp. 16, 17).
- [7] Kumar Chellapilla, Sidd Puri, and Patrice Simard. “High Performance Convolutional Neural Networks for Document Processing”. In: () (cit. on p. 7).
- [8] Chengjun Liu and H. Wechsler. “Gabor Feature Based Classification Using the Enhanced Fisher Linear Discriminant Model for Face Recognition”. In: *IEEE Transactions on Image Processing* 11.4 (Apr. 2002), pp. 467–476. ISSN: 1941-0042. DOI: 10.1109/TIP.2002.999679 (cit. on p. 9).
- [9] D.C. Cireşan et al. “Deep Neural Networks Segment Neuronal Membranes in Electron Microscopy Images”. In: *NIPS* 25 (2012). Export Date: 26 January 2023; Cited By: 92, pp. 2852–2860 (cit. on p. 8).
- [10] Dan Cireşan et al. “A Committee of Neural Networks for Traffic Sign Classification”. In: *The 2011 International Joint Conference on Neural Networks*. July 2011, pp. 1918–1921. DOI: 10.1109/IJCNN.2011.6033458 (cit. on p. 7).
- [11] Dan C. Cireşan et al. “Mitosis Detection in Breast Cancer Histology Images with Deep Neural Networks”. In: *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2013*. Ed. by Kensaku Mori et al. Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, 2013, pp. 411–418. ISBN: 978-3-642-40763-5. DOI: 10.1007/978-3-642-40763-5\_51 (cit. on p. 8).
- [12] Dan Claudiu Cireşan et al. “Deep, Big, Simple Neural Nets for Handwritten Digit Recognition”. In: *Neural Computation* 22.12 (Dec. 2010), pp. 3207–3220. ISSN: 0899-7667. DOI: 10.1162/NECO\_a\_00052. (Visited on 01/25/2023) (cit. on p. 7).
- [13] T.F Cootes et al. “View-Based Active Appearance Models”. In: *Image and Vision Computing* 20.9 (Aug. 2002), pp. 657–664. ISSN: 0262-8856. DOI: 10.1016/S0262-8856(02)00055-0 (cit. on p. 10).

- [14] Jiankang Deng et al. *RetinaFace: Single-stage Dense Face Localisation in the Wild*. May 2019. arXiv: [arXiv:1905.00641](#). (Visited on 04/13/2023) (cit. on pp. 10, 12–14).
- [15] Stuart E. Dreyfus. “The Computational Solution of Optimal Control Problems with Time Lag”. In: *IEEE Transactions on Automatic Control* 18.4 (1973). Cited by: 32, pp. 383–385. DOI: [10.1109/TAC.1973.1100330](#) (cit. on p. 7).
- [16] Hang Du et al. “The Elements of End-to-end Deep Face Recognition: A Survey of Recent Advances”. In: *ACM Computing Surveys* 54.10s (Jan. 2022), pp. 1–42. ISSN: 0360-0300, 1557-7341. DOI: [10.1145/3507902](#). (Visited on 03/07/2023) (cit. on pp. 10, 11, 13–15).
- [17] D. Elizondo. “The Linear Separability Problem: Some Testing Methods”. In: *IEEE Transactions on Neural Networks* 17.2 (Mar. 2006), pp. 330–344. ISSN: 1045-9227. DOI: [10.1109/TNN.2005.860871](#). (Visited on 01/24/2023) (cit. on p. 6).
- [18] B. Farley and W. Clark. “Simulation of Self-Organizing Systems by Digital Computer”. In: *Transactions of the IRE Professional Group on Information Theory* 4.4 (1954), pp. 76–84. DOI: [10.1109/TIT.1954.1057468](#) (cit. on p. 5).
- [19] Zhen-Hua Feng et al. *Wing Loss for Robust Facial Landmark Localisation with Convolutional Neural Networks*. Comment: 11 pages, 6 figures, 6 tables. Oct. 2018. arXiv: [arXiv:1711.06753](#). (Visited on 04/14/2023) (cit. on pp. 13, 15).
- [20] Jiuxiang Gu et al. “Recent Advances in Convolutional Neural Networks”. In: *Pattern Recognition* 77 (May 2018), pp. 354–377. ISSN: 0031-3203. DOI: [10.1016/j.patcog.2017.10.013](#). (Visited on 02/09/2023) (cit. on pp. 16, 17).

- [21] Michael Haenlein and Andreas Kaplan. “A Brief History of Artificial Intelligence: On the Past, Present, and Future of Artificial Intelligence”. In: *California Management Review* 61 (July 2019), p. 000812561986492. DOI: 10.1177/0008125619864925 (cit. on p. 6).
- [22] Donald Olding Hebb. *The Organization of Behavior: A Neuropsychological Theory*. Wiley, 1949. ISBN: 978-0-471-36727-7 (cit. on p. 5).
- [23] Sepp Hochreiter and Jürgen Schmidhuber. “Long Short-Term Memory”. In: *Neural Computation* 9.8 (Nov. 1997), pp. 1735–1780. ISSN: 0899-7667. DOI: 10.1162/neco.1997.9.8.1735. (Visited on 01/25/2023) (cit. on p. 7).
- [24] Peiyun Hu and Deva Ramanan. “Finding Tiny Faces”. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Honolulu, HI: IEEE, July 2017, pp. 1522–1530. ISBN: 978-1-5386-0457-1. DOI: 10.1109/CVPR.2017.166. (Visited on 04/14/2023) (cit. on p. 14).
- [25] Lichao Huang et al. *DenseBox: Unifying Landmark Localization with End to End Object Detection*. Sept. 2015. arXiv: arXiv:1509.04874. (Visited on 04/13/2023) (cit. on p. 13).
- [26] Xiehe Huang et al. *PropagationNet: Propagate Points to Curve to Learn Structure Information*. Comment: 10 pages, 8 figures, 8 tables, CVPR2020. June 2020. arXiv: arXiv:2006.14308. (Visited on 04/08/2023) (cit. on p. 10).
- [27] A G Ivakhnenko and V G Lapa. “Cybernetic Predicting Devices”. In: () (cit. on p. 6).
- [28] A. G. Ivakhnenko. “Polynomial Theory of Complex Systems”. In: *IEEE Transactions on Systems, Man, and Cybernetics* SMC-1.4 (1971), pp. 364–378. DOI: 10.1109/TSMC.1971.4308320 (cit. on p. 6).
- [29] Roger David Joseph. *Contributions to Perceptron Theory*. Cornell Aeronautical Laboratory, 1960 (cit. on p. 6).



- [30] Xu Kang, Bin Song, and Fengyao Sun. “A Deep Similarity Metric Method Based on Incomplete Data for Traffic Anomaly Detection in IoT”. In: *Applied Sciences* 9 (Jan. 2019), p. 135. DOI: 10.3390/app9010135 (cit. on p. 16).
- [31] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. “ImageNet Classification with Deep Convolutional Neural Networks”. In: *Advances in Neural Information Processing Systems*. Vol. 25. Curran Associates, Inc., 2012. (Visited on 01/26/2023) (cit. on pp. 8, 9, 15).
- [32] Erik Learned-Miller et al. “Labeled Faces in the Wild: A Survey”. In: *Advances in Face Detection and Facial Image Analysis*. Ed. by Michal Kawulok, M. Emre Celebi, and Bogdan Smolka. Cham: Springer International Publishing, 2016, pp. 189–248. ISBN: 978-3-319-25956-7 978-3-319-25958-1. DOI: 10.1007/978-3-319-25958-1\_8. (Visited on 03/09/2023) (cit. on p. 10).
- [33] Y. LeCun et al. “Backpropagation Applied to Handwritten Zip Code Recognition”. In: *Neural Computation* 1.4 (1989), pp. 541–551. DOI: 10.1162/neco.1989.1.4.541 (cit. on p. 6).
- [34] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. “Deep Learning”. In: *Nature* 521.7553 (May 2015), pp. 436–444. ISSN: 0028-0836, 1476-4687. DOI: 10.1038/nature14539 (cit. on pp. 16, 17).
- [35] Yann LeCun et al. “Gradient-Based Learning Applied to Document Recognition”. In: (1998) (cit. on pp. 7, 9).
- [36] Z. Lei, M. Pietikainen, and S.Z. Li. “Learning Discriminant Face Descriptor”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36.2 (2014). Cited By :287, pp. 289–302. ISSN: 0162-8828. DOI: 10.1109/TPAMI.2013.112 (cit. on p. 9).
- [37] Zewen Li et al. “A Survey of Convolutional Neural Networks: Analysis, Applications, and Prospects”. In: *IEEE Transactions on Neural Networks and Learning Systems* 33.12 (Dec. 2022), pp. 6999–7019. ISSN: 2162-2388. DOI: 10.1109/TNNLS.2021.3084827 (cit. on pp. 7, 8, 15, 17).

- [38] Seppo Linnainmaa. “The Representation of the Cumulative Rounding Error of an Algorithm as a Taylor Expansion of the Local Rounding Errors”. PhD thesis. Master’s Thesis (in Finnish), Univ. Helsinki, 1970 (cit. on p. 7).
- [39] Wei Liu et al. “SSD: Single Shot MultiBox Detector”. In: vol. 9905. Comment: ECCV 2016. 2016, pp. 21–37. DOI: 10.1007/978-3-319-46448-0\_2. arXiv: 1512.02325 [cs]. (Visited on 04/13/2023) (cit. on p. 13).
- [40] Yang Liu et al. *HAMBox: Delving into Online High-quality Anchors Mining for Detecting Outer Faces*. Comment: 9 pages, 6 figures. arXiv admin note: text overlap with 1802.09058 by other authors. Dec. 2019. arXiv: arXiv: 1912.09231. (Visited on 04/13/2023) (cit. on p. 13).
- [41] Brais Martinez et al. “Local Evidence Aggregation for Regression-Based Facial Point Detection”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35.5 (May 2013), pp. 1149–1163. ISSN: 1939-3539. DOI: 10.1109/TPAMI.2012.205 (cit. on p. 10).
- [42] J McCarthy et al. “A PROPOSAL FOR THE DARTMOUTH SUMMER RESEARCH PROJECT ON ARTIFICIAL INTELLIGENCE”. In: () (cit. on p. 5).
- [43] Warren S Mcculloch and Walter Pitts. “A LOGICAL CALCULUS OF THE IDEAS IMMANENT IN NERVOUS ACTIVITY”. In: () (cit. on p. 5).
- [44] Shervin Minaee et al. *Going Deeper Into Face Detection: A Survey*. Mar. 2021 (cit. on p. 10).
- [45] Marvin Minsky and Seymour Papert. *Perceptrons: An Introduction to Computational Geometry*. Cambridge, MA, USA: MIT Press, 1969 (cit. on p. 6).
- [46] Allen Newell, John C Shaw, and Herbert A Simon. “Report on a General Problem Solving Program”. In: *IFIP Congress*. Vol. 256. Pittsburgh, PA. 1959, p. 64 (cit. on p. 5).

- [47] Kyoung-Su Oh and Keechul Jung. “GPU Implementation of Neural Networks”. In: *Pattern Recognition* 37.6 (June 2004), pp. 1311–1314. ISSN: 0031-3203. DOI: 10.1016/j.patcog.2004.01.013 (cit. on p. 7).
- [48] P. N. Belhumeur, J. P. Hespanha, and D. J. Kriegman. “Eigenfaces vs. Fisherfaces: Recognition Using Class Specific Linear Projection”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 19.7 (July 1997), pp. 711–720. ISSN: 1939-3539. DOI: 10.1109/34.598228 (cit. on pp. 9, 11).
- [49] Rajat Raina, Anand Madhavan, and Andrew Y. Ng. “Large-Scale Deep Unsupervised Learning Using Graphics Processors”. In: *Proceedings of the 26th Annual International Conference on Machine Learning*. ICML ’09. New York, NY, USA: Association for Computing Machinery, June 2009, pp. 873–880. ISBN: 978-1-60558-516-1. DOI: 10.1145/1553374.1553486. (Visited on 01/25/2023) (cit. on p. 7).
- [50] Rajeev Ranjan et al. “Deep Learning for Understanding Faces: Machines May Be Just as Good, or Better, than Humans”. In: *IEEE Signal Processing Magazine* 35.1 (Jan. 2018), pp. 66–83. ISSN: 1558-0792. DOI: 10.1109/MSP.2017.2764116 (cit. on pp. 9–11, 14).
- [51] Marc’ aurelio Ranzato et al. “Efficient Learning of Sparse Representations with an Energy-Based Model”. In: *Advances in Neural Information Processing Systems*. Vol. 19. MIT Press, 2006. (Visited on 01/25/2023) (cit. on p. 7).
- [52] Shaoqing Ren et al. *Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks*. Comment: Extended tech report. Jan. 2016. arXiv: arXiv:1506.01497. (Visited on 04/13/2023) (cit. on p. 12).
- [53] F. Rosenblatt. “The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain.” In: *Psychological Review* 65 (1958), pp. 386–408. ISSN: 1939-1471(Electronic),0033-295X(Print). DOI: 10.1037/h0042519 (cit. on p. 5).

- [54] Frank Rosenblatt. *Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms*. Spartan Books, 1962 (cit. on p. 6).
- [55] DE Rumelhart, GE Hinton, and RJ Williams. *Learning Internal Representations by Error Propagation*, in *Parallel Distributed Processing*, DE Rumelhart, JL McClelland Eds. 1986 (cit. on p. 7).
- [56] Jürgen Schmidhuber. “Deep Learning in Neural Networks: An Overview”. In: *Neural Networks* 61 (Jan. 2015), pp. 85–117. ISSN: 08936080. DOI: 10.1016/j.neunet.2014.09.003. (Visited on 01/24/2023) (cit. on p. 6).
- [57] P.Y. Simard, D. Steinkraus, and J.C. Platt. “Best Practices for Convolutional Neural Networks Applied to Visual Document Analysis”. In: *Seventh International Conference on Document Analysis and Recognition, 2003. Proceedings*. Vol. 1. Edinburgh, UK: IEEE Comput. Soc, 2003, pp. 958–963. ISBN: 978-0-7695-1960-9. DOI: 10.1109/ICDAR.2003.1227801. (Visited on 01/25/2023) (cit. on p. 7).
- [58] J. Stallkamp et al. “Man vs. Computer: Benchmarking Machine Learning Algorithms for Traffic Sign Recognition”. In: *Neural Networks*. Selected Papers from IJCNN 2011 32 (Aug. 2012), pp. 323–332. ISSN: 0893-6080. DOI: 10.1016/j.neunet.2012.02.016. (Visited on 01/25/2023) (cit. on p. 7).
- [59] Yaniv Taigman et al. “DeepFace: Closing the Gap to Human-Level Performance in Face Verification”. In: *2014 IEEE Conference on Computer Vision and Pattern Recognition*. Columbus, OH, USA: IEEE, June 2014, pp. 1701–1708. ISBN: 978-1-4799-5118-5. DOI: 10.1109/CVPR.2014.220. (Visited on 02/13/2023) (cit. on pp. 9, 15).
- [60] Xu Tang et al. *PyramidBox: A Context-assisted Single Shot Face Detector*. Comment: 21 pages, 12 figures. Aug. 2018. arXiv: [arXiv:1803.07737](https://arxiv.org/abs/1803.07737). (Visited on 04/14/2023) (cit. on p. 14).

- [61] Jonathan Tompson et al. *Efficient Object Localization Using Convolutional Networks*. Comment: 8 pages with 1 page of citations. June 2015. arXiv: [arXiv:1411.4280](#). (Visited on 02/13/2023) (cit. on p. 15).
- [62] A. M. Turing. “I.—COMPUTING MACHINERY AND INTELLIGENCE”. In: *Mind* LIX.236 (Oct. 1950), pp. 433–460. ISSN: 1460-2113, 0026-4423. DOI: [10.1093/mind/LIX.236.433](#). (Visited on 01/13/2023) (cit. on p. 4).
- [63] Matthew Turk and Alex Pentland. “Eigenfaces for Recognition”. In: *Journal of Cognitive Neuroscience* 3.1 (Jan. 1991), pp. 71–86. ISSN: 0898-929X. DOI: [10.1162/jocn.1991.3.1.71](#). (Visited on 03/07/2023) (cit. on pp. 9, 11).
- [64] P. Viola and M. Jones. “Rapid Object Detection Using a Boosted Cascade of Simple Features”. In: *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*. Vol. 1. Dec. 2001, pp. I–I. DOI: [10.1109/CVPR.2001.990517](#) (cit. on p. 10).
- [65] Mei Wang and Weihong Deng. “Deep Face Recognition: A Survey”. In: *Neurocomputing* 429 (Mar. 2021), pp. 215–244. ISSN: 0925-2312. DOI: [10.1016/j.neucom.2020.10.081](#). (Visited on 02/27/2023) (cit. on pp. 9–11).
- [66] Mingxuan Wang et al. “genCNN: A Convolutional Architecture for Word Sequence Prediction”. In: *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Beijing, China: Association for Computational Linguistics, July 2015, pp. 1567–1576. DOI: [10.3115/v1/P15-1151](#). (Visited on 02/13/2023) (cit. on p. 15).
- [67] Nannan Wang et al. “Facial Feature Point Detection: A Comprehensive Survey”. In: *Neurocomputing* 275 (Jan. 2018), pp. 50–65. ISSN: 0925-2312. DOI: [10.1016/j.neucom.2017.05.013](#) (cit. on p. 10).
- [68] Joseph Weizenbaum. “ELIZA—a Computer Program for the Study of Natural Language Communication between Man and Machine”. In: *Commu-*

- nications of the ACM* 9.1 (Jan. 1966), pp. 36–45. ISSN: 0001-0782. DOI: 10.1145/365153.365168. (Visited on 01/18/2023) (cit. on p. 5).
- [69] J. Weng, N. Ahuja, and T.S. Huang. “Cresceptron: A Self-Organizing Neural Network Which Grows Adaptively”. In: *[Proceedings 1992] IJCNN International Joint Conference on Neural Networks*. Vol. 1. June 1992, 576–581 vol.1. DOI: 10.1109/IJCNN.1992.287150 (cit. on p. 7).
- [70] Chris J Winstead. “Remote Microelectronics Laboratory Education in the COVID-19 Pandemic”. In: *2022 Intermountain Engineering, Technology and Computing (IETC)*. May 2022, pp. 1–6. DOI: 10.1109/IETC54973.2022.9796805 (cit. on p. 1).
- [71] Lingyun Xiang et al. “A Convolutional Neural Network-Based Linguistic Steganalysis for Synonym Substitution Steganography”. In: *Mathematical Biosciences and Engineering* 17.2 (2020), pp. 1041–1058. ISSN: 1551-0018. DOI: 10.3934/mbe.2020055. (Visited on 02/13/2023) (cit. on p. 15).
- [72] Yuanyuan Xu et al. *CenterFace: Joint Face Detection and Alignment Using Face as Point*. Comment: 11 pages, 3 figures. A demo of CenterFace can be available at <https://github.com/Star-Clouds/CenterFace>. Nov. 2019. arXiv: arXiv:1911.03599. (Visited on 04/13/2023) (cit. on pp. 12–14).
- [73] Rikiya Yamashita et al. “Convolutional Neural Networks: An Overview and Application in Radiology”. In: *Insights into Imaging* 9.4 (Aug. 2018), pp. 611–629. ISSN: 1869-4101. DOI: 10.1007/s13244-018-0639-9. (Visited on 02/09/2023) (cit. on pp. 15, 16).
- [74] Ming Yang et al. “Detecting Human Actions in Surveillance Videos”. In: 2009 TREC Video Retrieval Evaluation Notebook Papers. Cited by: 26. 2009 (cit. on p. 7).
- [75] Z. Cao et al. “Face Recognition with Learning-Based Descriptor”. In: *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recog-*

- niton*. June 2010, pp. 2707–2714. ISBN: 1063-6919. DOI: 10.1109/CVPR.2010.5539992 (cit. on p. 9).
- [76] Stefanos Zafeiriou, Cha Zhang, and Zhengyou Zhang. “A Survey on Face Detection in the Wild: Past, Present and Future”. In: *Computer Vision and Image Understanding* 138 (Sept. 2015), pp. 1–24. ISSN: 1077-3142. DOI: 10.1016/j.cviu.2015.03.015 (cit. on p. 10).
  - [77] Andreas Zell. “Simulation Neuronaler Netze”. In: 1994 (cit. on p. 16).
  - [78] Caiming Zhang and Yang Lu. “Study on Artificial Intelligence: The State of the Art and Future Prospects”. In: *Journal of Industrial Information Integration* 23 (Sept. 2021), p. 100224. ISSN: 2452414X. DOI: 10.1016/j.jii.2021.100224. (Visited on 01/11/2023) (cit. on p. 5).
  - [79] Changzheng Zhang, Xiang Xu, and Dandan Tu. *Face Detection Using Improved Faster RCNN*. Feb. 2018. arXiv: arXiv:1802.02142. (Visited on 04/13/2023) (cit. on p. 13).
  - [80] Yu-Dong Zhang et al. “Improved Breast Cancer Classification Through Combining Graph Convolutional Network and Convolutional Neural Network”. In: *Information Processing & Management* 58.2 (Mar. 2021), p. 102439. ISSN: 0306-4573. DOI: 10.1016/j.ipm.2020.102439 (cit. on p. 15).
  - [81] Kaipeng Zhang et al. “Joint Face Detection and Alignment Using Multi-task Cascaded Convolutional Networks”. In: *IEEE Signal Processing Letters* 23.10 (Oct. 2016). Comment: Submitted to IEEE Signal Processing Letters, pp. 1499–1503. ISSN: 1070-9908, 1558-2361. DOI: 10.1109/LSP.2016.2603342. arXiv: 1604.02878 [cs]. (Visited on 04/13/2023) (cit. on pp. 10, 12, 14).
  - [82] Shifeng Zhang et al. *FaceBoxes: A CPU Real-time Face Detector with High Accuracy*. Comment: Accepted by IJCB 2017; Added references; Released codes. Dec. 2018. arXiv: arXiv:1708.05234. (Visited on 04/14/2023) (cit. on p. 14).

- [83] Xinqi Zhu and Michael Bain. *B-CNN: Branch Convolutional Neural Network for Hierarchical Classification*. Comment: 9 pages, 8 figures. Oct. 2017. DOI: 10.48550/arXiv.1709.09890. arXiv: arXiv:1709.09890. (Visited on 02/13/2023) (cit. on p. 16).