

0.1 Face Recognition

Face Recognition (FR) is a thoroughly debated and extensively researched task in the Computer Vision community for more than two decades [34], popularized in the early 1990s with the introduction of the Eigenfaces [41] or Fisherfaces [33] 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 [9, 2] that, even though, presented considerable results, weren't distinctive or compact. Beginning in 2010, methods based on learnable filters arose [52, 25], 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 [34] 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) [43].

CNNs are an Artificial Neural Network (ANN) that exhibit a better performance on image or video-based tasks compared to other methods [24]. They were greatly hailed in 2012, after the AlexNet [21] victory, by a great margin, in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC). Just two years later, DeepFace [38] 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 [43].

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 [22].

0.2 A Face Recognition System

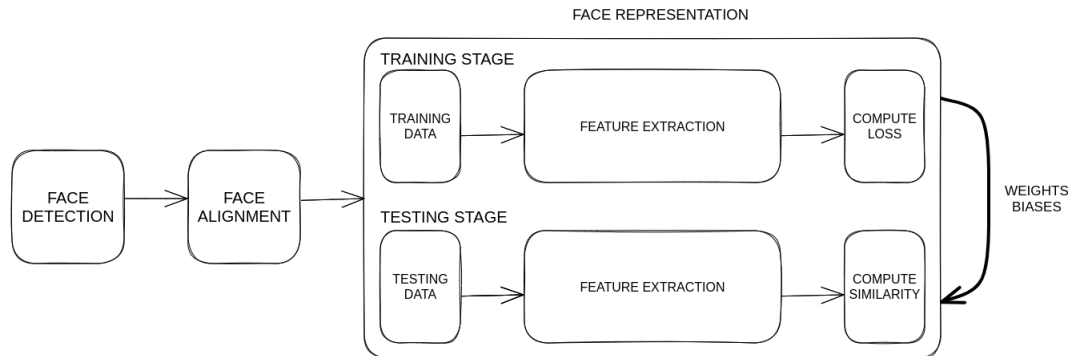


Figure 1: A typical face recognition pipeline, guided by the approach in [43].

According to Ranjan *et al.* [34], 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 [13], such as **military, public security and daily life**.

By and large, all end-to-end automatic face recognition systems follow a sequential and modular¹ pipeline (Figure 1) composed of three pillar stages [43]: 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 spatial transformations to the data in order to normalize the faces' pictures (or frames, in the case where a video is used) to a standardized view. 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 implementation². Face detection is achievable through classical approaches [42, 5] or deep methods, among them is [12] and the widely applied [57]. Face alignment, once again, can be accomplished through traditional measures [10, 31] or more modern ones, namely [19] or the aforementioned [57] which concurrently performs detection and alignment. To conclude, the face representation module is no exception,

¹ 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.

² For a deeper and extensive study, please refer to: [53] in the case of classic face detection approaches and [32] for deep learning based methods; [45] addresses traditional face alignment methods and is complemented with [13] for more up-to-date techniques; and [22] tackles classic face representation (add the following if needed) while X supplements the deep learning ones.

and can also be divided in two groups, regarding the methodology used. Some conventional systems were already mentioned, such as [33, 41], and the deep learning ones are the object 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.

0.2.1 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 [13, 34].

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) [13, 34]. 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 [13] distributes the methods, depending upon their architecture or purpose 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.

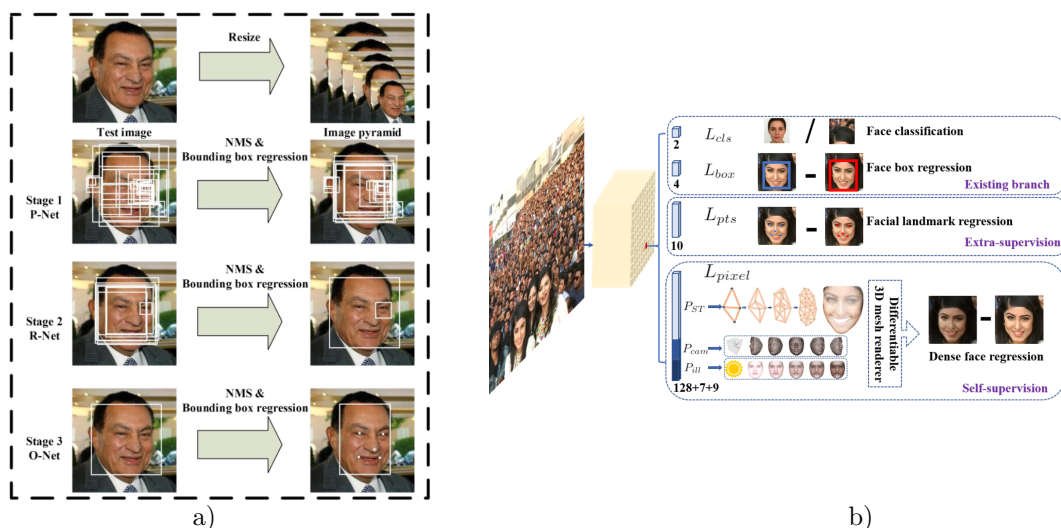


Figure 2: Comparison between **a)** MTCNN: multi-stage, CPU real-time and multi-task learning, and **b)** RetinaFace: single-stage, anchor-based, CPU real-time and multi-task learning. MTCNN [57] proposes a series of bounding boxes then, through a series of refinement stages, the best solution and landmarks are found. RetinaFace [12] 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 [12] 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 [35]. However, these methods can be slower and have a more complex way of training [49].

Single-stage approaches [12] 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 [29]. Finally, the methods included in this class are more efficient, but can incur in compromised accuracy, when compared to multi-stage.

Anchor-based techniques [30, 12, 55] 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.* [30], 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 [18] which attempts to predict faces by processing each pixel as a bounding box, or CenterFace [49] 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 [13].

Multi-task learning are all the methodologies that conjointly performs other tasks, namely facial landmark³ localization, during face classification and bounding box regression [13]. CenterFace [49] is one example, and so it is the widely implemented MTCNN [57], which correlated bounding boxes and face landmarks. RetinaFace [12] 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) [13]. MTCNN [57], Faceboxes [58], CenterFace [49] or RetinaFace [12] 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 [13]. PyramidBox [39] is an example that solves the partial occluded and blurry faces, and HR [17] 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 [34] that groups the methods in region⁴ or sliding-window⁵ based.

0.2.2 Face Alignment

Face Alignment, or facial landmark detection [7], 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 [13].

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 [14].

³ A facial landmark is a key-point in a face that contributes with important geometric information, namely the eyes, nose, mouth, etc. [14]

⁴ 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.

⁵ Sliding-window approaches centers on using a DCNN to compute a face detection score and bounding box at every location in a feature map.

Furthermore, face alignment, can be accomplished through two categories of methods: landmark-based and landmark-free.

Landmark-based alignment is a category of methods that exploits the facial landmarks with the aim of, through spatial transformations, calibrating the face to an established layout [13]. This can be accomplished through: coordinate regression, heatmap regression or 3D Model Fitting. **Coordinate regression-based** methodologies [14, 28, 57] consider the landmark localization as a numerical objective, i.e. a regression, thus an image is fed to a DCNN and it will output a vector of landmark coordinates. **Heatmap Regression** [11, 46, 8] is different from coordinate regression because, although it is a numerical objective task, the output is not a coordinate vector, but a map of likelihood of landmarks' locations. Finally, **3D Model Fitting** [4, 7, 48] is the category that integrates methods that consider the relation between 2D facial landmarks and the 3D shape of a generic face. The particularity of them is the reconstruction of the 3D face from a 2D face image that is then projected over a plane in order to obtain the landmarks.

Landmark-free alignment, on the other hand, integrates the approaches that do not rely on landmarks as a reference to align the face, in contrast, these type of methods incorporate the alignment into a DCNN that gives, as a result, an aligned face [13]. An example of an end-to-end method that does not depend on facial landmarks is RDCFace [59], and it rectifies distortions, applies alignment transformations and executes face representation. Hayat et al. [16] proposes a method that deals with extreme head poses. The process to register faces in an image with high pose variance can be quite challenging and often demands complex pre-processing, namely landmark localization, therefore, to address that, a DCNN is employed that does not rely on landmark localization and concomitantly register and represent faces.

As can be seen from the previous section, this step in the face recognition process can be accomplished, very sporadically, through standalone methods that process the detected face from the previous stage, but generally joint detection and alignment methods (and sometimes even face representation), previously referenced in the multi-task learning definition, are the optimal choice [7].

0.2.3 Face Representation

Finally, Face Representation is the last stage of the Face Recognition process. It is responsible for processing the aligned face from the previous stage and mapping the produced feature representation to a feature space, in which features from the same person are closer together and those that are different stand further apart from each other [13].

According to the literature [13, 26, 34, 36, 43], there's a consensus about how Face Recognition can be performed in two settings of operation: face verification and face identification. This distinction is only made possible due to the approaches available in the Face Representation stage that can leverage one, the other or both.

Face verification, also referred to as **face authentication**, is a one-to-one match, and it's the action of verifying if the query face matches the identity that's being claimed. **These principles are used in biometric systems such as self-service immigration clearance using E-passport.** [26]

Face identification, also called **face recognition**, is a one-to-many correlation process that compares a query face to a database of faces and associates it to the corresponding match (or matches). **A typical use case is to identify someone in a watchlist or surveillance videos.** [26]

The overall pipeline comes to a conclusion in this module, however, in reality, it goes further than that. As can be seen in (Figure 1), due to its importance for the face recognition problem, it's highlighted the inherent pipeline of the Face Representation stage, henceforth, it shall be discussed in depth in the next section.

0.3 Face Representation Pipeline

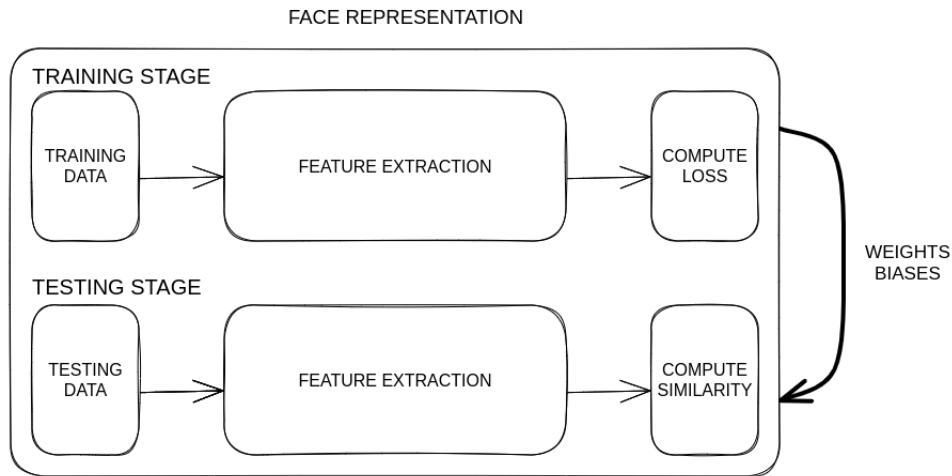


Figure 3: Face Representation pipeline, guided by the approach in [43].

As shown in (Figure 3), Face Representation is a two-step module composed of a training and testing stage. So as to be capable of performing face recognition, in either a verification or identification manner, a face representation system needs to learn robust, invariant and discriminative features that can distinguish identities [34].

To meet these requirements, the feature extractor must first be trained properly by taking data from previous stages and producing an output that's compared to the desired value using a loss function⁶ [23, 43]. After that, everything is ready for the testing stage, where the face recognition *per se* occurs by calculating a similarity score for the feature representation produced by the trained feature extractor, and dictating if the identity belongs to the same person (face verification) or if it matches any identity (face identification) [34].

0.3.1 Datasets - Training and Testing Data

As been discussed throughout this dissertation, Deep Learning techniques can solve the problem of handling unconstrained scenarios, where there are variations in pose, illumination, occlusion, and so forth. To support this, in the past few years, datasets have been developed with this in mind so to be able to provide a large and diverse set of both training data, allowing for adequate regularization to unseen circumstances, and testing data that benchmarks the face recognition system in, as similar as possible, unconstrained real world scenarios [13].

⁶ Also referred to as a cost or objective function. [23]

Large training datasets are essential for face recognition, therefore, there is an abundance of options available. Earlier works [36, 37, 38] resourced private face datasets, but as time passed, more public datasets were released.

CASIA-WebFace [51], composed of 494,414 face images and 10,575 identities, was proposed as a novel dataset to overcome the problem of data dependence in face recognition and improve comparability across different methods. By training on the same dataset, methods can be better evaluated and compared.

VGGFace

CelebA

UMDFaces

MS-Celeb-1M

MegaFace

VGGFace2

UMDFaces-Videos

MS-Celeb-1M Low-shot

0.3.2 Feature Extraction - Deep 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 [50, 27] with successful applications in the domains of Computer Vision [21, 38, 40, 56] or Natural Language Processing[1, 44, 47]. In the CNN category itself there are different variants, but they all abide the fundamental structure of a feedforward hierarchical multi-layer network (Figure 4). Feedforward because the information only flows in a singular direction without cycling [54], hierarchical because the higher complexity internal representations are learned from lower ones [23, 60] 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 [23, 50, 15, 3]. The convolutional layer is designed to extract feature representations by being composed of kernels (or filter banks [23]) that compute feature maps through element-wise product, to which is applied a nonlinear activation function [15, 50]. Next is the pooling layer, that's responsible

for reducing the spatial size of the input data [15] and joining identical features [23]. Finally, the fully connected layers, and their core function is to perform high logic and generate semantic information [15].

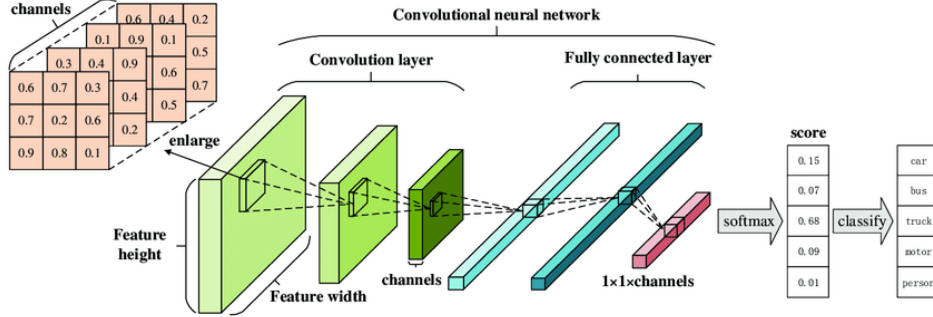


Figure 4: Architecture of a Convolutional Neural Network [20].

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 [23, 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 [50]. In the second place, there's an inherent correlation between local groups of values, which creates distinguishable motifs [23]. 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 [23]. Therefore, the following key, unique features potentiate the previously stated efficient performance [6]:

1. Designed to process multidimensional arrays [23];
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, 27];
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 [15, 23], while 2, 3, 4 and 5 helps to reduce the complexity of the model, and as a result training it is easier [15, 27].

Transfer Learning

0.3.3 Loss