

# The Elements of End-to-end Deep Face Recognition: A Survey of Recent Advances

HANG DU\*, Shanghai University, China  
HAILIN SHI\*, JD AI Research, China  
DAN ZENG†, Shanghai University, China  
XIAO-PING ZHANG, Ryerson University, Canada  
TAO MEI, JD AI Research, China

Face recognition is one of the most popular and long-standing topics in computer vision. With the recent development of deep learning techniques and large-scale datasets, deep face recognition has made remarkable progress and been widely used in many real-world applications. Given a natural image or video frame as input, an end-to-end deep face recognition system outputs the face feature for recognition. To achieve this, a typical end-to-end system is built with three key elements: face detection, face alignment, and face representation. The face detection locates faces in the image or frame. Then, the face alignment is proceeded to calibrate the faces to the canonical view and crop them with a normalized pixel size. Finally, in the stage of face representation, the discriminative features are extracted from the aligned face for recognition. Nowadays, all of the three elements are fulfilled by the technique of deep convolutional neural network. In this survey article, we present a comprehensive review about the recent advance of each element of the end-to-end deep face recognition, since the thriving deep learning techniques have greatly improved the capability of them. To start with, we present an overview of the end-to-end deep face recognition. Then, we review the advance of each element, respectively, covering many aspects such as the to-date algorithm designs, evaluation metrics, datasets, performance comparison, existing challenges, and promising directions for future research. Also, we provide a detailed discussion about the effect of each element on its subsequent elements and the holistic system. Through this survey, we wish to bring contributions in two aspects: first, readers can conveniently identify the methods which are quite strong-baseline style in the subcategory for further exploration; second, one can also employ suitable methods for establishing a state-of-the-art end-to-end face recognition system from scratch.

Additional Key Words and Phrases: Deep learning, convolutional neural network, face recognition, face detection, face alignment, face representation.

## 1 INTRODUCTION

Face recognition (FR) is an extensively studied topic in computer vision. Among the existing technologies of human biometrics, face recognition is the most widely used one in real-world applications. With the great advance of deep convolutional neural networks (DCNNs), the deep learning based methods have achieved significant improvements on various computer vision tasks, including face recognition. In this survey, we focus on 2D image based end-to-end deep face recognition which takes the general images or video frames as input, and extracts the deep feature of each face as output. We provide a comprehensive review of the recent advances of the elements of end-to-end deep face recognition. Specifically, an end-to-end deep face recognition system is composed of three key elements: face detection, face alignment, and face representation. In the following, we give a brief introduction of each element.

---

\*Equal contribution. This work was performed at JD AI Research.

†Corresponding author.

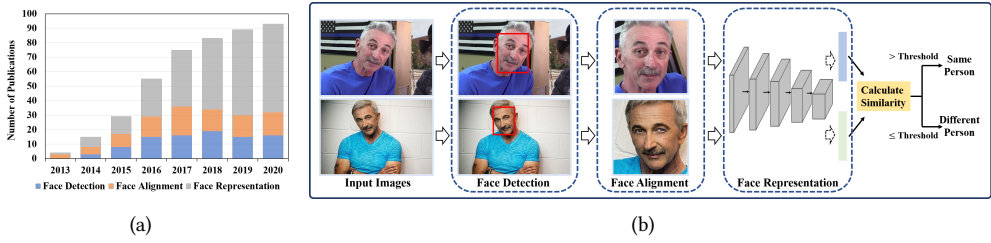


Fig. 1. (a) The publication trend of the elements of end-to-end deep face recognition from 2013 to 2020. (b) The standard pipeline of end-to-end deep face recognition system. First, the face detection stage aims to localize the face region on the input image. Then, the face alignment is proceeded to normalize the detected face to the canonical view. Finally, the face representation devotes to extracting features for recognition.

Face detection is the first step of end-to-end face recognition. It aims to locate the face regions in the still images or video frames. Before the deep learning era, one of the pioneering works for face detection is Viola-Jones [230] face detector, which utilizes AdaBoost classifiers with Haar features to build a cascaded structure. Later on, the subsequent approaches explore the effective hand-craft features [8, 162, 169] and various classifiers [17, 123, 150] to improve the detection performance. One can refer to [286, 299] for a thorough survey of traditional face detection methods.

Next, face alignment refers to calibrate the detected face to the canonical view and crop it to a normalized pixel size, in order to facilitate the subsequent task of face representation computing. It is an essential intermediate procedure for face recognition system. Generally, the facial landmark localization is necessary for face alignment, while some approaches can directly generate aligned face from the input one. Most traditional works of facial landmark localization focus on either generative methods [36, 37] or discriminative methods [153, 345], and there are several exhaustive surveys about them [99, 247, 358].

In the face representation stage, the discriminative features are extracted from the aligned face images for recognition. This is the final and core step of face recognition. In early studies, many approaches calculate the face representation by projecting face images into low-dimensional subspace, such as Eigenfaces [227] and Fisherfaces [13]. Later on, handcrafted local descriptors based methods [3, 131] prevail in this area. For a detailed review of these traditional methods, one can refer to [7, 231, 307]. In the last few years, the face representation benefits from the development of DCNNs and witnesses great improvements for high performance face recognition.

This survey focuses on reviewing and analyzing the recent advances in each element. An important fact is that, the performance of face recognition depends on the contribution of all the elements (*i.e.*, face detection, alignment and representation). In other words, inferiority in any one of the elements will become the bottleneck and harm the final performance. In order to establish high-performance end-to-end face recognition system, it is necessary to understand every element of the holistic framework and their intrinsic connection. A number of face recognition surveys have been published in the past twenty years. The main differences between our survey and the existing ones are summarized as follows.

- **The relationship between the elements and whole.** We provide the thorough discussion about the effect of each element on its subsequent one and the holistic system, which are overlooked in the existing surveys. From the existing experiments and detailed analysis, we can conclude the performance of the holistic system depends on the three elements. Therefore,

it is necessary to review them together for helping the readers who aim to establish state-of-the-art face recognition system from scratch.

- **More recently published works.** The publications in the last three years (2018-2020) are much more than all those published before 2018 (as illustrated in Fig. 1(a)). In view of the rapid development of face detection, face alignment and face representation in the past few years, this survey covers the recently published articles. By doing so, we provide the up-to-date review of the elements, and large number of newly presented methods.
- **New analysis for future work.** Based on the up-to-date review, we conclude the promising trends from the newest frontier, and several insightful thoughts of each element as well as the holistic system, to enlighten the future research.

Specifically, there are certain surveys [7, 231, 307] about face recognition who, however, do not cover deep learning based methods since they are published early before the deep learning era; besides, another set of surveys focus on 3D face recognition [16, 201] and specific tasks [49, 356]. Instead, we focus on the 2D face recognition which is the most needed in practical applications. For deep learning based 2D face recognition, there are a small number of articles that fulfil relevant survey, which differ from this paper in many ways. Among them, Ranjan *et al.* [177] do not include the recent techniques that rapidly evolved in the past few years. In fact, the number of published works has been increasing dramatically during these years (as shown in Fig. 1(a)). Wang and Deng [244] present a systematic review about deep face representation rather than the end-to-end face recognition. More recently, Insaf *et al.* [2] provide a review of 2D and 3D face recognition from the traditional to deep-learning era, while the scope is still limited in the face representation. In summary, the end-to-end face recognition, covering all the elements of the pipeline, needs to be systematically reviewed, while seldom of the existing survey articles attach importance to this task.

Therefore, we systematically review the deep learning based approaches of each element in the end-to-end face recognition, respectively. The review of each element covers many aspects: algorithm designs, evaluation metrics, datasets, performance comparisons, remaining challenges, and promising directions for future research. We hope this survey could bring helpful thoughts for better understanding of the big picture of end-to-end face recognition and deeper exploration in a systematic way. Specifically, the main contributions can be summarized as follows:

- We provide a comprehensive survey of the elements of end-to-end deep face recognition. We review the recent advances of each element, respectively, and present elaborated categorizations of them to make the readers understand them in a systematic way.
- We review the three elements from many aspects: algorithm designs, evaluation metrics, datasets, and performance comparison. Moreover, we point out the effect of each element on its subsequent elements and the holistic system.
- We collect the existing challenges and promising directions for each element and its subcategories to facilitate the future research, and further discuss the major challenges and future trends from the view of the holistic framework.

## 2 OVERVIEW

A typical end-to-end deep face recognition system includes three basic elements: face detection, face alignment, and face representation, as shown in Fig. 1(b). First, face detection localizes the face region on the input image. Then, face alignment is proceeded to normalize the detected face into the canonical layout. Finally, face representation devotes to extracting discriminative features from the aligned face. The features are used to calculate the similarity between them, in order to make the decision that whether the faces belong to the same identity.

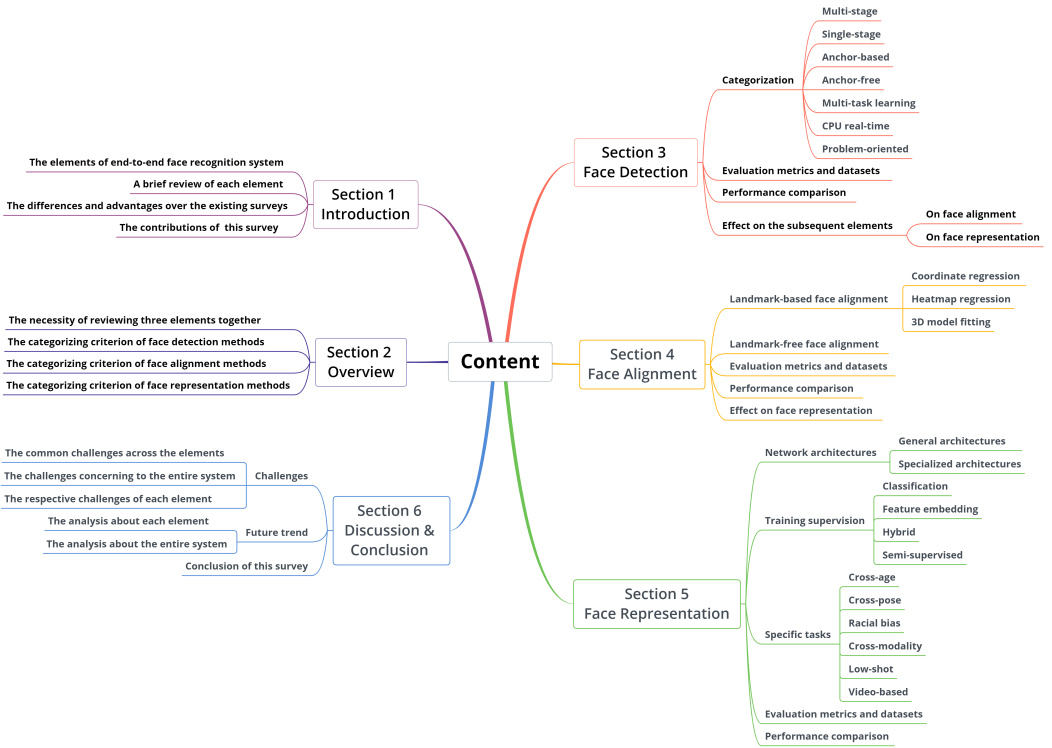


Fig. 2. The structure of this survey. The left parts (Section 1, 2, 6) refer to the functional contents that provide overall introduction and discussion. The right parts (Section 3, 4, 5) refer to the technical contents that provide the detailed reviewing of three elements.

The structure of this survey is illustrated in Fig. 2. We structure the body sections (Section 3, 4, 5) with respect to the three elements, each of which is a research topic that covers abundant literatures in computer vision. We give an overview of the three elements briefly in this section, and dive into each of them in the following body sections.

## 2.1 Face Detection

Face detection is the first procedure of the face recognition system. Given an input image, the face detection aims to find all the faces in the image and give the coordinates of bounding box with a confidence score. The major challenges of face detection contain varying resolution, scale, pose, illumination, occlusion, *etc.* In Section 3, we provide a categorization of the deep learning based face detection methods from multiple dimensions, which includes multi-stage, single-stage, anchor-based, anchor-free, multi-task learning, CPU real-time and problem-oriented methods. It is worth noting that there exist overlapping techniques between the categories, because, the categorization is built up from multiple perspectives.

**Differences to the existing survey of face detection.** Minaee *et al.* [160] review face detection methods from the beginning of deep learning era, and categorize them by design of network architecture. Compared with them, our categorizing criterion covers poly-aspects. Specifically, we provide a multiple-dimension categorization, to discuss the face detection methods from many different perspectives, which will help us to better understand the developing line and conclude





Fig. 3. Visualization of facial landmarks of different versions. The 4-point and 5-point landmarks are often used for face alignment.

the future trend. Since face detection state of the art is relatively advanced, such comprehensive categorization is necessary for readers.

## 2.2 Face Alignment

In the second stage, face alignment aims to calibrate the detected face to the canonical view. Since human face appears with the regular structure, in which the facial parts (eyes, nose, mouth, *etc*) have constant arrangement, the alignment of face is of great benefit to the subsequent feature computation for face recognition. For most existing methods of face alignment, the facial landmarks, or so-called facial keypoints (as shown in Fig. 3), are indispensable, because they are involved as the reference for similarity transformation or affine transformation. So, the facial landmark localization is a prerequisite for face alignment. The DCNNs based facial landmark localization methods can be divided into three subcategories: coordinate regression, heatmap regression and 3D model fitting based approaches. Without relying on the facial landmarks, several approaches can directly output aligned face from the input by learning the transformation parameters. We will review these methods in Section 4.

**Differences to the existing survey of face alignment.** Previous surveys of face alignment [99, 247, 358] only focus on reviewing the facial landmark localization methods. Since the landmark-free face alignment is also a kind of methods to generate aligned images for face recognition, we further collect them in this survey.

## 2.3 Face Representation

As the key step of face recognition system, face representation devotes to learning deep face model and using it to extract features from aligned faces for recognition. The features are used to calculate the similarity of the matched faces. In Section 5, we provide a review of deep learning based methods for discriminative face features, and retrospect these methods with respect to the network architecture and the training supervision. For network architecture, we introduce the general architectures which are designed for a wide range of computer vision tasks, and the special architectures which are specialized for face representation. As for training supervision, we mainly introduce four schemes, including the classification, feature embedding, hybrid and semi-supervised schemes. Additionally, we present several specific face recognition scenes, including cross domain, low-shot learning and video based scenarios.

**Differences to the existing survey of face representation.** This survey aims to provide the readers with a better understanding of the end-to-end face recognition. Recently, Wang and Deng [244] present a systematic review about deep face recognition, in which they mainly focus on deep face representation, and the categorization of training loss is sub-optimal. For instance, they sort the supervised learning of deep face representation by Euclidean-distance based loss, angular/cosine-margin-based loss, softmax loss and its variations; while, in fact, almost all the angular/cosine-margin-based losses are implemented as the variation of softmax loss rather than an individual set. In contrast, we suggest a more reasonable categorization with three subcategories, *i.e.*, classification, feature embedding, and hybrid methods (in Section 5.2).

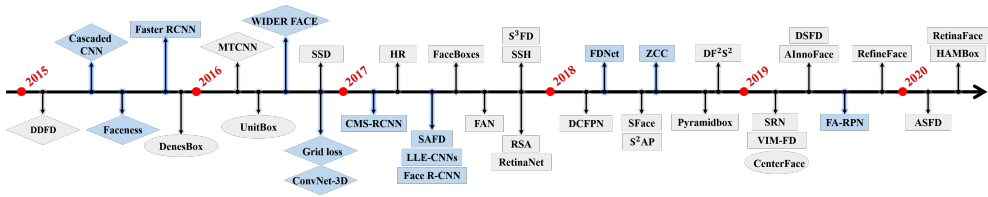


Fig. 4. The development of representative face detection methods. The blue and gray represent multi-stage and single-stage methods; according to the anchor usage, the rectangle, oval, and diamond denote anchor-based, anchor-free and other methods. One can refer to Table 1 for the references of these methods.

### 3 FACE DETECTION

Face detection is the first step of end-to-end face recognition system, which aims to locate the face regions from the input images. In this section, first, we categorize and make comparison of the existing deep learning methods for face detection. Next, we introduce several popular datasets of face detection and the common metrics for evaluation. Finally, we provide a performance comparison of state-of-the-art face detection methods and detailed discussion about the effect of face detection on its subsequent elements.

#### 3.1 Categorization of Face Detection

In order to present the deep face detection methods with a clear categorization, we group them with seven sets, *i.e.*, multi-stage, single-stage, anchor-based, anchor-free, multi-task learning, CPU real-time, and problem-oriented methods (in Table 1). These sets are not necessarily exclusive, because we establish the categorization from multiple perspective. Fig. 4 is the development of representative methods for face detection.

**3.1.1 Multi-stage methods.** Following the coarse-to-fine manner or the proposal-to-refine strategy, multi-stage based detectors first generate a number of candidate boxes, and then refine the candidates by one or more additional stages. The first stage employs sliding window to propose the candidate bounding boxes at a given scale, and the latter stages reject the false positives and refine the remaining boxes. In such regime, the cascaded architecture [119, 193, 301, 314] is naturally an effective solution for the coarse-to-fine face detection.

Face detection can be considered as a specific objective of general object detection. Thus, many works [27, 63, 93, 97, 124, 166, 176, 208, 304, 348] inherit the remarkable achievements from the general object detectors. For instance, Faster R-CNN [181] is a classic and effective detection framework which employs a region proposal network to generate region proposals with a set of dense anchor boxes in the first stage, and then refines the proposals in the second stage. Based on the proposal-to-refine scheme, many works have dedicated to improve the modeling of the refinement stage [93, 97, 208, 304, 348] and the proposal stage [27, 75, 124, 166, 202], and achieved great progress for accurate face detection. Apart from the modeling, how to train the multi-stage detector is another interesting topic. To tackle the issue of inferior optimization for multi-stage detectors, a joint training strategy [174] is designed for both Cascaded CNN [119] and Faster R-CNN to achieve end-to-end optimization and better performance.

**3.1.2 Single-stage methods.** The single-stage methods accomplish the candidate classification and bounding box regression from the feature maps directly, without the dependence on proposal stage.

A classic structure of single stage comes from a general object detector named Single Shot multibox Detector (SSD) [136]. It runs much faster than the multi-stage ones while maintaining

Table 1. The categorization of deep face detection methods.

Category	Description	Method
Multi-stage	Detectors first generate candidate boxes, then the following one or more stages refine the candidates.	Faceness [287], HyperFace [176], STN [27], ConvNet-3D [124], SAFD [75], CMS-RCNN [348], Wan <i>et al.</i> [232], Jiang <i>et al.</i> [97], DeepIR [208], Grid loss [170], Face R-CNN [93], Face R-FCN [253], ZCC [347], FDNet [304], FA-RPN [166], Cascaded CNN [119], MTCNN [314], Qin <i>et al.</i> [174], LLE-CNNs [63], PCN [193], PPN [301]
Single-stage	Detectors accomplish face classification and bounding box regression from feature maps at once.	DDFD [59], DenseBox [89], UnitBox [294], HR [85], Faceboxes [321], SSH [165], S <sup>3</sup> FD [322], DCFPN [323], FAN [242], FANet [313], RSA [143], S <sup>2</sup> AP [202], PyramidBox [221], DF <sup>2</sup> S <sup>2</sup> [223], SFace [241], DSFD [120], RefineFace [318], SRN [32], PyramidBox++ [125], CenterFace [279], VIM-FD [330], ISRN [320], AlInnoFace [308], ASFD [303], RetinaFace [41], HAMBox [145]
Anchor-based	Detectors deploy a number of dense anchors on the feature maps, and then proceed the classification and regression on these anchors.	Wan <i>et al.</i> [232], Face Faster RCNN [97], RSA [143], Face R-CNN [93], FDNet [304], DeepIR [208], SAFD [75], SSH [165], S <sup>3</sup> FD [322], DCFPN [323], Faceboxes [321], FAN [242], FANet [313], PyramidBox [221], ZCC [347], S <sup>2</sup> AP [202], DF <sup>2</sup> S <sup>2</sup> [223], SFace [241], RetinaFace [41], DSFD [120], RefineFace [318], SRN [32], VIM-FD [330], PyramidBox++ [125], FA-RPN [166], ISRN [320], AlInnoFace [308], Group Sampling [161], HAMBox [145], ASFD [303].
Anchor-free	Detectors directly find faces without preset anchors.	DenseBox [89], UnitBox [294], CenterFace [279]
Multi-task learning	Detectors jointly learn the classification and bounding box regression with additional tasks (e.g., landmark localization) in one framework.	STN [27], ConvNet-3D [124], HyperFace [176], MTCNN [314], Face R-CNN [93], RetinaFace [41], DF <sup>2</sup> S <sup>2</sup> [223], FLDet [355], PyramidBox++ [125], CenterFace [279]
CPU real-time	Detectors can run on a single CPU core in real-time for VGA-resolution images.	Cascade CNN [119], STN [27], MTCNN [314], DCFPN [323], Faceboxes [321], PCN [193], RetinaFace [41], FLDet [355], FBI [98], PPN [301], CenterFace [279]
Problem-oriented	Detectors aim to solve specific challenges in face detection, such as tiny faces, occluded faces, rotated and blurry faces.	HR [85], SSH [165], S <sup>3</sup> FD [322], Bai <i>et al.</i> [9], PyramidBox [221], Grid loss [170], FAN [242], LLE-CNNs [63], PCN [193], Group Sampling [161]

comparable accuracy. Based on SSD, many studies [98, 221, 321–323] develop deep face detectors those are robust to different scales of face. As for the backbone architecture, many face detectors resort to the feature pyramid network (FPN) [127] which consists of a top-down architecture with skip connections and merges the high-level and low-level features for detection. The high-level feature maps provide rich semantic information, while the low-level layers supplement more local information. The feature fusion preserves the advantages from both sides, and brings great progress in detecting objects with a wide range of scales. Therefore, many single-stage face detectors [32, 41, 120, 125, 165, 221, 223, 242, 313, 320] are developed with the advantage of FPN. Not only handling the scale issue in face detection via FPN, but also these methods attempt to solve the inherent shortcomings of original FPN such like the conflict of receptive field.

Although the single-stage methods have the advantage of high efficiency, their detection accuracy is below that of the two-stage methods. It is partially because the imbalance problem of positives and negatives brought by the dense anchors, whereas the proposal-to-refine scheme is able to alleviate this issue. Accordingly, RefineDet [319] sets up an anchor refinement module in its network to remove large number of negatives. Inspired by RefineDet, SRN [32] presents a selective two-step classification and regression method; the two-step classification is performed at the low-level layers to reduce the search space of classifier, and the two-step regression is performed at high-level layers to obtain accurate location. Later on, many works [308, 318, 320, 330] improve SRN with several effective techniques, such as training data augmentation, improved feature extractor and training supervision, anchor assignment and matching strategy, multi-scale test strategy, *etc.*

Most aforementioned methods need to preset anchors for face detection, while some representative detectors of single-stage, such as DenseBox [89], UnitBox [294] and CenterFace [279], fulfil the detection without preset anchors. We will present them as anchor-free type in the next subsection.

**3.1.3 Anchor-based and anchor-free methods.** As shown in Table 1, most current face detectors are anchor-based due to the long-time development and superior performance. Generally, we preset the anchors on the feature maps, then fulfil the classification and bounding box regression on these anchors one or more times, and finally output the accepted ones as the detection results. Therefore, the anchor allocation and matching strategy are crucial to the detection accuracy. Most anchor-based

methods focus on the algorithms along this direction, such as scale compensation [145, 322], max-out background label [322], expected max overlapping score [347], group sampling by scale [161], *etc.* However, the settings (*e.g.*, scale, stride, ratio, number) of anchors need to be carefully tuned for each particular dataset, limiting their generalization ability. Besides, the dense anchors increase the computational cost and bring the imbalance problem of positive and negative anchors.

Anchor-free methods [116, 224, 346] attract growing attention in general object detection. As for face detection, certain pioneering works have emerged in recent years. DenseBox [89] and UnitBox [294] attempt to predict the pixel-wise bounding box on face. CenterFace [279] regards face detection as a generalized task of keypoint estimation, which predicts the facial center point and the size of bounding box in feature map. In brief, the anchor-free detectors get rid of the preset anchors and achieve better generalization capacity. Regarding to the detection accuracy, it needs further exploration for better robustness to false positives and stability in training process.

**3.1.4 Multi-task learning methods.** Generally, the multi-task learning methods are designed for solving a problem together with other related tasks by sharing the visual representation. Here, we introduce the multi-task learning methods that train the face detector with the associated facial tasks or auxiliary supervision branches to enrich the feature representation and detection robustness. Many approaches [27, 89, 124, 279, 305, 314, 355] have explored the joint learning of face detection and facial landmark localization. Among them, MTCNN [314] is the most representative one, which exploits the inherent correlation between facial bounding boxes and landmarks. Subsequently, HyperFace [176] fuses the low-level features and high-level features to simultaneously conduct four tasks, including face detection, facial landmark localization, gender classification and pose estimation. Based on RetinaNet [128], RetinaFace [41] integrates face detection, facial landmark localization and dense 3D face regression in one framework. From the multi-task routine, we can see that the face detectors can benefit from the associated facial tasks. Moreover, certain methods [93, 125, 223, 241] exploit auxiliary supervision branches, such as segmentation branch, anchor-free branch, *etc.* These branches are used to boost the training of face detection.

**3.1.5 CPU real-time methods.** Although state-of-the-art face detectors have achieved great success in accuracy, their efficiency is not enough for real-world applications, especially on non-GPU devices. According to the demand of inference speed on CPU, we collect the CPU real-time face detectors [27, 41, 193, 279, 301, 321, 323, 355] here for convenient retrieval. These detectors are able to run at least 20 frames per second (FPS) on a single CPU with VGA-resolution input images. We provide a table in the supplemental material which shows the running efficiency of them, among which the lightweight backbone [41, 279], rapidly digested convolutional layer [321, 323], knowledge distillation [98] and region-of-interest (RoI) convolution [27] are the common practices.

**3.1.6 Problem-oriented methods.** We highlight some problem-oriented methods which are designed against a variety of specific challenges in face detection. Detecting faces with a wide range of scale is a long-existing challenge in face detection. A group of methods [85, 161, 165, 221, 322] are designed for scale-invariant face detection, including scale selection, multi-scale detection, dense anchor setting, scale balancing strategy, *etc.* The partially visible faces (*i.e.*, with occlusion) is another issue that harms the detection recall. The existing solutions [63, 170, 242, 287] resort to the facial part arrangement, anchor-level attention and data augmentation by generation, *etc.* Likewise, the in-plane rotation is an existing factor that impedes face detection. To tackle this problem, PCN [193] calibrates the candidates against the rotation progressively.

Table 2. Statistics of popular datasets for face detection.

Datasets	Year	#Image	#Face	# of faces per image	Description
Training					
ALFW [111]	2011	21,997	25,993	1.18	Training source for face detection.
WIDER FACE [288]	2016	16K	199K	12.43	The largest face detection dataset.
Test					
FDDB [229]	2010	2,845	5,171	1.82	A classic face detection benchmark.
AFW [353]	2012	205	473	2.31	Multiple facial annotations.
PASCAL faces [280]	2014	851	1,335	1.57	Large facial variations.
MALF [281]	2015	5,250	11,931	2.27	Fine-grained evaluation.
WIDER FACE [288]	2016	16K	194K	12.12	The largest face detection dataset.
MAFA [63]	2017	30,811	35,806	1.16	Masked face detection.

### 3.2 Evaluation Metrics and Datasets

**3.2.1 Metrics.** Like the general object detection, average precision (AP) is a widely used metric for evaluating the face detection methods. AP is derived from the precision-recall curve. To obtain precision and recall, Intersection over Union (IoU) is used to measure the overlap of the predicted bounding box ( $Box_p$ ) and the ground-truth ( $Box_{gt}$ ), which can be formulated as

$$IoU = \frac{area(Box_p \cap Box_{gt})}{area(Box_p \cup Box_{gt})}. \quad (1)$$

The prediction of face detector includes a predicted bounding box and its confidence score. The confidence score is used to determine whether to accept this according to the confidence threshold. Then, an accepted prediction can be regarded as true positive (TP) if the IoU is larger than a preset IoU threshold (usually 0.5 for face detection). Otherwise, it will be regarded as a false positive (FP). After determining the TP and FP, the precision-recall curve can be drawn by varying the confidence threshold. AP is computed as the mean precision at a series of uniformly-spaced discrete recall levels [57]. Apart from AP, the receiver operating characteristic (ROC) curve is also adopted as the metric, such as the evaluation in FDDB [229]; frames per second (FPS) is used to measure the runtime efficiency of detectors.

**3.2.2 Datasets.** We introduce several widely used datasets for face detection. The statistics of them are given in Table 2. Among them, FDDB [229] is a classic dataset of unconstrained face detection which includes low resolution, occlusion and difficult pose variations. It is noteworthy that FDDB uses ellipse as ground-truth instead of rectangular box. The images in PASCAL faces dataset [280] are taken from the Pascal person layout dataset [58]. WIDER FACE [288] provides a large number of training data and a challenging test benchmark with large data variations.

### 3.3 Performance Comparison

Table 3 shows the performance of the existing face detectors on WIDER FACE validation and test subsets. From the viewpoint of subcategory, we can observe that the single-stage methods with anchor-based mechanism (e.g., RefineFace [318], HAMBox [145]) dominate the state-of-the-art performance. For many real-world applications, MTCNN [314], Faceboxes [321], and RetinaFace [41] are the widely used face detectors for building a face recognition system, since they can achieve good balance between the detection accuracy and efficiency.

### 3.4 Effect on the Subsequent Elements

Face detection is the very first procedure in the end-to-end face recognition system, and thereby plays the role of *input* towards face alignment and face representation. The quality of detection bounding box directly influences on the performance of the subsequent alignment. There are two possible cases, *i.e.*, the loss of facial region and the excessively residual context region in the

Table 3. The performance of state-of-the-art methods on the WIDER FACE [288] validation and test subsets. The evaluation metric is AP.

Method	Publication	Subcategory	WIDER FACE Val.			WIDER FACE Test		
			Easy	Medium	Hard	Easy	Medium	Hard
Faceness-WIDER [288]	CVPR'16	Multi-stage	0.713	0.634	0.345	0.716	0.604	0.315
MSC-CNN [288]	CVPR'16	Multi-stage	0.691	0.664	0.424	0.711	0.636	0.400
CMS-RCNN [348]	DLB'17	Multi-stage	<b>0.899</b>	<b>0.874</b>	<b>0.624</b>	<b>0.902</b>	<b>0.874</b>	<b>0.643</b>
Face R-CNN [93]	arXiv'17	Multi-stage, Anchor-based	0.937	0.921	0.831	0.932	0.916	0.827
Face R-FCN [253]	arXiv'17	Multi-stage, Anchor-based	0.947	0.935	0.874	0.943	0.931	0.876
ZCC [347]	CVPR'18	Multi-stage, Anchor-based	0.949	0.933	0.861	0.949	0.935	0.865
FDNet [304]	arXiv'18	Multi-stage, Anchor-based	<b>0.959</b>	<b>0.945</b>	<b>0.879</b>	<b>0.950</b>	<b>0.939</b>	0.878
FA-RPN [166]	CVPR'19	Multi-stage, Anchor-based	0.949	0.941	0.894	0.945	0.937	<b>0.891</b>
MTCNN [314]	SPL'16	Multi-stage, CPU real-time, Multi-task learning	0.848	0.825	0.598	0.851	0.820	0.607
HR [85]	CVPR'17	Single-stage	0.925	0.910	0.806	0.923	0.910	0.819
SSH [165]	ICCV'17	Single-stage, Anchor-based	0.931	0.921	0.845	0.927	0.915	0.844
SFD [322]	ICCV'17	Single-stage, Anchor-based	0.937	0.925	0.859	0.935	0.921	0.858
FAN [242]	arXiv'17	Single-stage, Anchor-based	0.952	0.940	0.900	0.946	0.936	0.885
PyramidBox [221]	ECCV'18	Single-stage, Anchor-based	0.961	0.950	0.889	0.956	0.946	0.887
SRN [32]	AAAI'19	Single-stage, Anchor-based	0.964	0.952	0.901	0.959	0.948	0.896
VIM-FD [330]	arXiv'19	Single-stage, Anchor-based	0.967	0.957	0.907	0.962	0.953	0.902
DSFD [120]	CVPR'19	Single-stage, Anchor-based	0.964	0.957	0.904	0.960	0.953	0.900
ISRN [320]	arXiv'19	Single-stage, Anchor-based	0.967	0.958	0.909	0.963	0.954	0.903
AlnoFace [308]	arXiv'19	Single-stage, Anchor-based	0.971	0.961	0.918	0.965	0.957	0.912
RefineFace [318]	TPAMI'20	Single-stage, Anchor-based	0.971	0.962	0.920	0.965	0.958	0.914
HAMBox [145]	CVPR'20	Single-stage, Anchor-based	0.970	0.964	<b>0.933</b>	0.960	0.955	<b>0.923</b>
ASFD [303]	arXiv'20	Single-stage, Anchor-based	<b>0.972</b>	<b>0.965</b>	0.925	<b>0.967</b>	<b>0.962</b>	0.921
Faceboxes [321]	IJCB'17	Single-stage, Anchor-based, CPU real-time	0.840	0.766	0.395	0.839	0.763	0.396
DF <sup>2</sup> S <sup>2</sup> [223]	arXiv'18	Single-stage, Anchor-based, Multi-task learning	0.969	0.959	0.912	<b>0.963</b>	0.954	0.907
PyramidBox++ [125]	arXiv'19	Single-stage, Anchor-based, Multi-task learning	0.965	0.959	0.912	0.956	0.952	0.909
CenterFace [279]	arXiv'19	Single-stage, Anchor-free, CPU real-time, Multi-task learning	0.935	0.924	0.875	0.932	0.921	0.873
RetinaFace [41]	CVPR'20	Single-stage, Anchor-based, CPU real-time, Multi-task learning	<b>0.971</b>	<b>0.961</b>	<b>0.918</b>	<b>0.963</b>	<b>0.958</b>	<b>0.914</b>

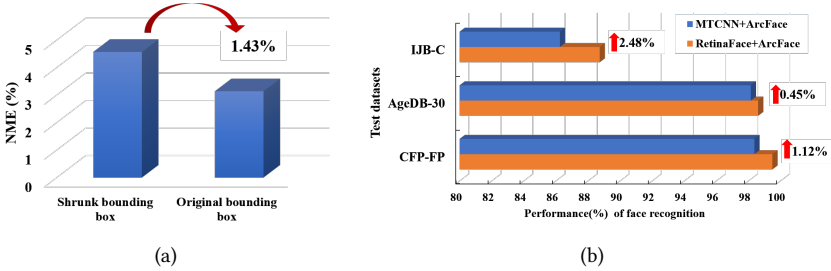


Fig. 5. The accuracy of face detection can influence on the subsequent elements, *i.e.*, face alignment and face representation. (a) Inaccurately detected bounding boxes will bring the performance degradation to facial landmark localization [275]. (b) A more robust face detector can further improve the recognition accuracy [41].

bounding box, both of which are the adverse factors to the subsequent process. Certain relevant literature shows solid evidence of that face detection influence on the face alignment and face recognition. Firstly, the quality of detected bounding boxes has a significant impact to facial landmark localization. For example, Xiong *et al.* [275] compare the performance (Fig. 5(a)) of facial landmark localization on the correct face bounding boxes and shrunk face bounding boxes, indicating that the inaccurately detected bounding boxes will bring the performance degradation to the landmark localization. Moreover, as shown in Fig. 5(b), RetinaFace [41] compares the recognition accuracy after using different face detection methods, which proves that a robust face detector can further improve the face recognition accuracy. In summary, face detection has significant impact to both face alignment and face representation. It is indispensable to consider the effect of face detection when establishing high-performance face recognition system.

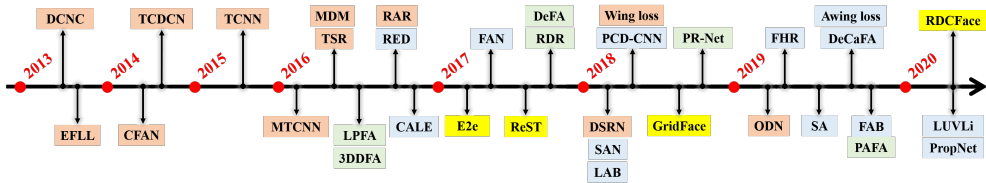


Fig. 6. The development of representative methods for face alignment. The orange, blue, green, and yellow represent coordinate regression, heatmap regression, 3D model fitting, and landmark-free face alignment methods, respectively. One can refer to Table 4 for the references of these methods.

Table 4. The categorization of face alignment methods.

Category		Description	Method
Landmark-based Face Alignment	Coordinate regression	Take the landmark coordinates as the target of regression, and learn the nonlinear mapping from the input face image to the landmark coordinates.	DCNC [212], EFL [344], CFAN [312], TCDCN [331], RAR [273], MDM [226], TSR [149], JEA [277], SIR [271], TCNN [268], DSRN [159], SBR [53], Wing loss [61], AAN [296], ODN [350], HyperFace [176], MTCNN [314], RetinaFace [41], FLDet [355], CenterFace [279], RDN [133]
	Heatmap regression	Output the likelihood response maps of each landmark.	CALE [19], RED [173], Yang <i>et al.</i> [100], JMFA [44], FAN [20], LAB [264], SAN [51], FALGCN [158], DU-Net [222], Guo <i>et al.</i> [70], PCD-CNN [113], RCN(L+ELT) [81], HR-Net [240], Zhang <i>et al.</i> [311], SA [147], FHR [219], Awing loss [249], DeCaFA [38], HSLE [357], FAB [207], KDN [29], Dong <i>et al.</i> [52], LaplaceKL [182], LUVLi [114], PropagationNet [91]
	3D model fitting	Infer the 3D face shape from 2D image, and then project it to the image plane to obtain 2D landmarks.	LPFA [101], 3DDEFA [351], FacePoseNet [25], PIFASCNN [102], DeFA [141], RDR [272], Bhagavatula <i>et al.</i> [14], Zhang <i>et al.</i> [309], PR-Net [60], PAPA [121]
Landmark-free Face Alignment		Directly output aligned faces without the explicit use of landmark.	Hayat <i>et al.</i> [76], E2e [340], ReST [263], GridFace [343], Wei <i>et al.</i> [256], RDC-Face [332]

## 4 FACE ALIGNMENT

Given the detected face, face alignment aims to calibrate unconstrained faces to the canonical layout for facilitating the downstream tasks of recognition and analysis. In this section, we review the mainstream routines for face alignment, including landmark-based face alignment, and landmark-free face alignment. Fig. 6 shows the development of representative methods for face alignment.

### 4.1 Landmark-based Face Alignment

Landmark-based face alignment utilizes the spatial transformation to calibrate faces to the pre-defined canonical layout by involving the facial landmarks as the reference. Therefore, the facial landmark localization is the core task of landmark-based alignment. We sort the existing landmark-based alignment methods into three subcategories, *i.e.*, coordinate regression based methods, heatmap regression based methods and 3D model fitting based methods.

**4.1.1 Coordinate regression.** The coordinate regression based methods regard the landmark coordinates as the numerical objective of the regression via neural networks. In other words, they focus on learning the nonlinear mapping from the face image to the landmark coordinate vectors.

Following the coarse-to-fine manner, most methods employ cascaded regression [149, 212, 312, 344] or recurrent neural network (RNN) [226, 273] to progressively refine the prediction of landmark coordinate. Besides, the multi-task learning is also a common routine to facilitate landmark localization with the related facial tasks, such as face detection [41, 176, 279, 314, 355] and facial attribute recognition [277, 331]. Moreover, many regression methods employ the L1, L2, or smoothed L1 loss functions, which are effective but, nonetheless, sensitive to outliers. To handle this problem, Wing loss [61] amplifies the impact of the samples with small or medium range errors. The above methods study the facial landmark localization on still images. For video face landmark localization,

how to leverage the temporal information across frames becomes necessary. TSTN [132] develops a two-stream architecture, which locates the landmark from a single frame and captures the temporal consistency for refinement. Besides, SBR [53] proposes to exploit the optical flow coherency of detected landmarks when training with video data.

**4.1.2 Heatmap regression.** In contrast to the coordinate regression, the heatmap regression based methods output likelihood response maps of each landmark. The early exploration [19] studies how to aggregate the score maps and refine the prediction with DCNNs. Later on, Newell *et al.* [168] design stacked hourglass (HG) network to generate heatmap for human pose estimation, which has achieved great success. As the facial landmark localization task is similar to the human pose estimation, many works [20, 44, 91, 100, 249, 311] adopt the stacked HG network for facial landmark localization and greatly improve the state-of-the-art performance.

The dense pixel-wise classification by the fully convolutional network (FCN) is an effective way for the heatmap regression task. The HG structure can be regarded as an instance of the fully convolutional network. Beyond the HG structure, a number of effective network architectures [38, 51, 113, 158, 240] are newly designed for heatmap regression. Among them, DeCaFA [38] utilizes stacked U-nets to preserve the spatial resolution, and landmark-wise attention maps to extract local information around the current estimation. High-resolution network (HR-Net) [240] is designed to maintain the high-resolution representation and shows its advantage for landmark-kind tasks.

The above-mentioned wing loss, which is designed for the coordinate regression, however, does not guarantee the convergence for the heatmap regression, due to the imbalance pixel number of foreground and background. To address this issue, Wang *et al.* [249] propose adaptive wing loss to penalize more on foreground pixels than on background pixels; similarly, PropNet [91] presents a focal wing loss which adjusts the loss weight of samples in each mini-batch.

Some facial landmarks have ambiguous definition, such as those on cheek, leading to inconsistent annotations by different annotators. Besides, the landmarks in occluded facial regions also cause imprecise annotations. Many methods [29, 114, 147, 147, 264, 357] devote to these two issues. Facial boundary heatmap [264] is a good choice to provide the facial geometric structure for reducing the semantic ambiguities. Regarding the semantic ambiguities as noisy annotation, Liu *et al.* [147] provide another path to estimate the real landmark location with a probabilistic model. More recently, KDN [29] and LUVLi [114] propose to estimate the uncertainty of predictions. The uncertainty can be used to identify the images in which the face alignment fails.

**4.1.3 3D model fitting.** Considering the explicit relationship between 2D facial landmarks and 3D face shape, the 3D model fitting based methods reconstruct the 3D face shape from 2D image, and then project it onto the image plane to obtain the 2D landmarks. Compared with the regular 2D methods which estimate a set of landmarks, 3D model fitting based methods are able to fit faces with 3D model of thousands of vertexes and align them with large poses.

Since the cascaded regression is an effective manner to estimate model parameters, some methods [101, 141, 351] combine the cascaded CNN regressor with a dense 3D Morphable Model (3DMM) [15] to estimate the 3D face shape. Despite many advantages, the cascaded CNNs often suffer from the lack of end-to-end training. As a roundabout, Jourabloo *et al.* [102] attempt to fit a 3D face model through a single CNN, which consists of several blocks to adjust the 3D shape and projection matrix according to the features and predictions from the previous blocks.

Although the above methods take great advantages from 3DMM, the diverse facial shape would lead to inaccurate 2D landmark location, especially when the 3D shape coefficients are sparse. To tackle this problem, RDR [272] proposes to fit 3D faces by a dynamic expression model and use a recurrent 3D-2D dual learning model to alternatively refine 3D face model and 2D landmarks. Beyond regressing the parameters of 3D face shape, Faster-TRFA [14] and FacePoseNet [25] estimate



Table 5. Statistics of popular facial landmark datasets. “-” refers to none official protocol for splitting the training and test set.

Datasets	Year	# Total	# Training	# Test	# Point	Description
Multi-PIE [67]	2008	755,370	-	-	68	The largest facial dataset in controlled condition.
LFPW [12]	2010	2,845	-	-	35	Images taken from uncontrolled setting.
ALFW [111]	2011	24,386	20,000	4,386	21	A large-scale facial landmark dataset.
AFW [353]	2012	473	-	-	6	Multiple facial annotations.
HELEN [117]	2012	2,330	2,000	330	194	Providing dense landmark annotations.
COFW [21]	2013	1,852	1,345	507	29	Containing occluded faces.
300-W [185]	2013	3,837	3,148	689	68	The most frequently used dataset of facial landmark.
300-VW [192]	2015	114	50	64	68	A video facial landmark dataset.
Menpo [298]	2017	28,273	12,014	16,259	68	Containing both semi-frontal and profile faces.
WFLW [264]	2018	10,000	7,500	2,500	98	Multiple annotations and large variations.
JD-landmark [144]	2019	15,393	13,393	2,000	106	Covering large facial variations.

the warping parameters of rendering a different view of a general 3D face model. Besides, some methods [60, 309] aim to directly regress the landmarks from the 3D coordinates of face shape.

## 4.2 Landmark-free Face Alignment

Landmark-free face alignment methods integrate the alignment transformation processing into DCNNs and output aligned face without relying on facial landmarks. This set of methods generally employ the spatial transformer network (Spatial-TN) [95] for geometric warping, where the transformation parameters are learned via end-to-end training. Based on Spatial-TN, Hayat *et al.* [76] and Zhong *et al.* [340] propose to optimize the face alignment with a subsequent module of face representation jointly. Since the facial variations are quite complex with various factors, some methods [263, 343] are designed to improve the deformation ability of Spatial-TN. Besides, the radial distortion of face images is another common problem, which is brought by the wide-angle cameras. RDCFace [332] presents a cascaded network which learns the rectification against the radial lens distortion, the face alignment transformation, and the face representation in an end-to-end manner.

## 4.3 Evaluation Metrics and Datasets

We introduce the commonly used evaluation metrics and datasets for face alignment. As presented in the following part of this subsection, most landmark-based methods employ the quantitative metrics, such as normalized mean error. Besides, landmark-free methods employ the evaluation oriented to face recognition, and we will describe their metrics in the face representation section.

**4.3.1 Metrics.** The widely used evaluation metric is to measure the point-to-point Euclidean distance by normalized mean error (NME), which can be defined as

$$NME = \frac{1}{M} \sum_{k=1}^M \frac{\|p_k - g_k\|_2}{d}, \quad (2)$$

where  $M$  is the number of landmarks,  $p_k$  and  $g_k$  represent the prediction and ground-truth coordinates of the face landmarks,  $k$  denotes the index of landmarks, and  $d$  refers to the normalized distance which is used to alleviate the abnormal measurement caused by different face scales and large pose. There are four types of normalized distance for computing NME, *i.e.*, the geometric mean of the width and height of the face bounding box, the distance between the outer corners of eyes, the distance between the pupils, and the diagonal of the face bounding box.

The cumulative errors distribution (CED) curve is also used as an evaluation criterion. CED is a distribution function of NME. The vertical axis of CED represents the proportion of test images that have an error value less than or equal to the error value on the horizontal axis. The area under

Table 6. Performance of facial landmark localization methods on the 300W, WFLW-All, ALFW-Full, and COFW datasets. The evaluation metric is NME (%). For 300W test set, two types of NME normalization (*i.e.*, inter-pupil normalization and inter-ocular normalization) are used. For WFLW-All dataset, the inter-ocular normalization is applied. For ALFW-Full dataset, the diagonal of face bounding box is adopted as the normalization factor. For COFW dataset, the inter-pupil normalization is applied. “-” indicates that the authors do not report the performance with the corresponding protocol.

Method	Publication	Subcategory	300-W (inter-pupil normalization)			300-W (inter-ocular normalization)			WFLW	ALFW	COFW
			Com. subset	Chall. subset	Full set	Com. subset	Chall. subset	Full set			
CFAN [312]	ECCV'14	Coordinate regression	5.50	16.78	7.69	-	-	-	-	10.94	8.38
TCDCN [331]	ECCV'14	Coordinate regression	4.80	8.60	5.54	-	-	-	-	7.60	8.05
MDM [226]	CVPR'16	Coordinate regression	4.83	5.88	10.14	-	-	-	-	-	6.26
RAR [273]	ECCV'16	Coordinate regression	4.12	8.35	4.94	-	-	-	-	7.23	6.03
TSR [149]	CVPR'17	Coordinate regression	4.36	7.56	4.99	-	-	-	-	2.17	-
SIR [271]	AAAI'18	Coordinate regression	4.29	8.14	5.04	-	-	-	-	-	-
DSRN [159]	CVPR'18	Coordinate regression	4.12	9.68	5.21	-	-	-	-	1.86	-
Wing loss [61]	CVPR'18	Coordinate regression	<b>3.01</b>	<b>6.01</b>	<b>3.60</b>	-	-	-	<b>5.11</b>	<b>1.47</b>	5.44
CPM + SBR [53]	CVPR'18	Coordinate regression	-	-	-	<b>3.28</b>	7.58	<b>4.10</b>	-	2.14	-
ODN [350]	CVPR'19	Coordinate regression	-	-	-	3.56	<b>6.67</b>	4.17	-	1.63	<b>5.30</b>
RDN [133]	TPAMI'20	Coordinate regression	3.31	7.04	4.23	-	-	-	-	2.06	5.82
3DDFA [351]	CVPR'16	3D model fitting	6.15	10.59	7.01	-	-	-	-	5.60	-
PIFASCNN [102]	ICCV'17	3D model fitting	5.43	9.88	6.30	-	-	-	-	4.45	-
DeFA [141]	ICCVW'17	3D model fitting	5.37	9.38	6.10	-	-	-	-	-	-
RDR [272]	ICCV'17	3D model fitting	5.03	8.95	5.80	-	-	-	-	4.41	-
PAFA [121]	BMVC'19	3D model fitting	<b>3.42</b>	<b>5.73</b>	<b>3.87</b>	-	-	-	-	<b>1.51</b>	<b>3.55</b>
RCN+ [81]	CVPR'18	Heatmap regression	4.20	7.78	4.90	-	-	-	-	2.17	-
PCD-CNN [113]	CVPR'18	Heatmap regression	-	-	-	3.67	7.62	4.44	-	2.40	5.77
SAN [51]	CVPR'18	Heatmap regression	-	-	-	3.34	6.60	3.98	-	1.91	-
HR-Net [240]	CVPR'18	Heatmap regression	-	-	-	2.87	5.15	3.32	-	1.57	-
LAB [264]	CVPR'18	Heatmap regression	3.42	6.98	4.12	2.98	5.19	3.49	5.27	<b>1.25</b>	3.92
DU-Net [222]	ECCV'18	Heatmap regression	-	-	-	2.90	5.15	3.35	-	-	-
SA [147]	CVPR'19	Heatmap regression	3.45	6.38	4.02	-	-	-	-	1.60	-
HG-HSLE [357]	ICCV'19	Heatmap regression	3.94	7.24	4.59	2.85	5.03	3.28	-	-	-
Awing loss [249]	ICCV'19	Heatmap regression	3.77	6.52	4.31	2.72	4.52	3.07	<b>4.36</b>	-	4.94
LaplaceKL [182]	ICCV'19	Heatmap regression	<b>3.28</b>	7.01	<b>4.01</b>	-	-	-	-	1.97	-
DeCaFA [38]	ICCV'19	Heatmap regression	-	-	-	2.93	5.26	3.39	4.62	-	-
LUVLi [114]	CVPR'20	Heatmap regression	-	-	-	2.76	5.16	3.23	4.37	1.39	-
PropNet [91]	CVPR'20	Heatmap regression	3.70	<b>5.75</b>	4.10	<b>2.67</b>	<b>3.99</b>	<b>2.93</b>	<b>4.05</b>	-	<b>3.71</b>

the curve (AUC) also provides a reference of how the algorithm performs at a given error:

$$AUC_{\alpha} = \int_0^{\alpha} f(e)de, \quad (3)$$

where  $\alpha$  is the given error corresponding to the upper bound of integration calculation,  $e$  is the progressive normalized errors and  $f(e)$  refers to the CED curve. Larger AUC indicates better performance. Based on CED curve, failure rate can be used to measure the robustness of an algorithm, which denotes the percentage of samples in the test set whose NME is larger than a threshold.

**4.3.2 Datasets.** The facial landmark datasets can be sorted by the constrained condition and in-the-wild condition. The statistics of these datasets are given in Table 5. CMU Multi Pose, Illumination, and Expression (Multi-PIE) [67] is the largest facial dataset in constrained condition, which provides 337 subjects with 15 predefined poses, 19 illumination conditions and 6 facial expressions. The annotated facial landmarks are 68 points for frontal faces and 39 points for profile ones.

In addition, more in-the-wild datasets [12, 21, 111, 117, 144, 185, 192, 264, 298, 353] are proposed for facial landmark localization. Among them, 300-W [185] is the most frequently used dataset, which follows the annotation configuration of Multi-PIE and re-annotates the images in LFPW, AFW, HELEN, and iBug [184]. Besides, Menpo [298] is a large-scale facial landmark dataset with more difficult cases for facial landmark localization. JD-landmark [144] annotates face images with 106 facial landmarks, providing more structural information of facial components. 300-VW [192] provides 50 video clips for training and 64 for test of landmark localization in video.

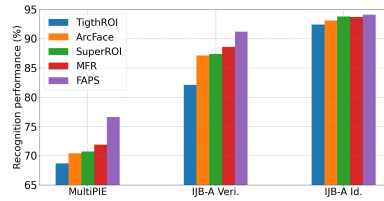


Fig. 7. Appropriate face alignment policy is beneficial to face recognition in many situations [278]. The indicated choices of alignment policy are different in number of used facial landmarks, cropping size of face image, and vertical shift. Among them, ArcFace [41] employs a 5-point alignment template, and MFR [22] utilizes a 25-point one. TigthROI [278] involves few external facial feature (*e.g.*, jaw-line, ears, part of hair), which lacks useful facial features. SuperRol [278] uses large cropping size, which potentially covers irrelevant background. FAPS [278] is designed to search an optimal face alignment template. The latter three policies use 68 landmarks which provide adequate information for computing the affine transformation matrix.

#### 4.4 Performance Comparison

Table 6 shows the comparison of state-of-the-art facial landmark localization methods on various test datasets, including 300-W, WLFW-ALL, ALFW-Full, and COFW. Among coordinate regression based methods, Wing loss [61] is a simple but effective approach which has been widely used. More recently, the heatmap regression based methods attract more attention, since they can obtain the leading performance by maintaining facial structure information throughout the models.

#### 4.5 Effect on the Face Representation

Face alignment is the intermediate procedure. The study of how face alignment influences on face representation is vital for tuning the recognition system to attain its maximum effect. For landmark-based face alignment, a set of inaccurate facial landmarks will harm the alignment and then impede the following feature computation as well. Specifically, human faces appear in the images with similar layout, and such layout can be regarded as a template in spatial coordinates. In fact, the alignment is accomplished mostly by warping the face to the predefined coordinates according to the predicted landmarks. Then, the face representation model learns identity feature from facial images with such layout. Once the predicted landmarks are inaccurate, the facial image will drift away from the predefined coordinates, which is unexpected layout for the face representation model. Guo *et al.* [70] and Deng *et al.* [41] both compare the widely used MTCNN [314] and their methods, and find that poor landmark localization will bring shift variation, while robust face alignment can boost recognition accuracy, especially for the cross-pose face recognition. Besides, as discussed in certain studies [172, 189, 278], the configuration of face alignment process (so-called face alignment policy), including the number of used facial landmarks, the cropping size of face image, and the vertical shift, greatly influences on the performance of face recognition. As shown in Fig. 7, the results from [278] indicate that the proper face alignment policy is beneficial to face recognition in many situations. Moreover, a moderate degree of spatial transformation is required in the alignment processing [256]. Both limited and excessive transformation will bring disturbance.

### 5 FACE REPRESENTATION

Subsequent to face alignment, the face representation stage aims to map the aligned face images to a feature space, where the features of the same identity are close and those of different identities are far apart. In practical applications, there are two major tasks of face recognition, *i.e.*, face verification and face identification. The face verification refers to predict whether a pair of face

Table 7. The categorization of face representation learning.

Category		Description	Method
Network Architectures	General	The basic and universal designs for common visual recognition tasks.	AlexNet [112], VGGNet [198], GoogleNet [217], ResNet [77], Xception [34], DenseNet [64], AttentionNet [236], SENet [84], SqueezeNet [94], MobileNet [83], ShuffleNet [328], MobileNetV2 [186], Shufflenetv2 [151]
	Specialized	The modified or ensemble designs oriented to face recognition.	HybridDL[213], DeepID series [209, 211, 215], MM-DFR [48], B-CNN [35], ComparatorNet [274], Contrastive CNN [74], PRN [106], AFRN [105], FANFace [282], SparseNet [216], Light-CNN [266], MobileFaceNet [30], Mobiface [56], ShuffleFaceNet [154], Hayat <i>et al.</i> [76], E2e [340], ReST [263], GridFace [343], RDCFace [332], Wei <i>et al.</i> [256], Co-Mining [250], GroupFace [109], MFR [22]
Training Supervision	Classification	Considering the face representation learning as a classification task.	DeepFace [220], DeepID [214], MM-DFR [48], L-softmax [138], NormFace [237], L2-softmax [175], COCO loss [142], SphereFace [137], Ring loss [339], AM-softmax [235], CosFace [239], ArcFace [42], AdaptiveFace [134], Fair loss [130], MV-softmax [251], ArcNeg [295], AdaCos [326], P2SGrad [327], NTP [86], Co-Mining [250], PFE [194], CurricularFace [92], Shi <i>et al.</i> [196], GroupFace [109], MFR [22], RCM loss [270], DUL [26]
	Feature embedding	Optimizing the feature distance according to the label of sample pair.	DeepID2 [209], FaceNet [189], VGG Face [172], Lifted structured [203], N-pair loss [200], Multibatch [218], TPE [187], Smart mining [152], Contrastive CNN [74]
	Hybrid	Applying classification and feature embedding together as the supervisory signals.	DeepID2 [209], DeepID2+ [215], DeepID3 [211], TUA [135], Center loss [259], Marginal loss [45], Range loss [325], DM [199], PRN [106], UniformFace [55], RegularFace [336], UT [342], LBL [352], AFRN [105], Circle loss [210]
	Semi-supervised	Exploiting labeled and unlabeled faces for representation learning.	CDP [302], GCN-DS [285], GCN-VE [284], UIR [293], RoyChowdhury <i>et al.</i> [183]
Specific Tasks	Cross-age	Identifying faces across a wide range of ages.	LF-CNNs [258], CAN [276], AFRN [54], DAL [238], AE-CNN [228], OE-CNN [252], IPC-GANs [254], LMA [5], Dual cGANs [204], AIM [333]
	Cross-pose	Identifying faces across a wide range of poses.	TP-GAN [90], PIM [334], DREAM [23], DA-GAN [335], DR-GAN [225], UV-GAN [43], CAPG-GAN [87], PAMs [155], MPRs [1], MvDN [104]
	Racial bias	Addressing the imbalance race distribution of training datasets.	IMAN [246], RL-RBN [245]
	Cross-modality	Performing face recognition on a pair of images captured by different sensing modalities.	Reale <i>et al.</i> [180], HFR-CNNs [188], TRIVET [140], IDR [79], DVR [267], MC-CNN [47], WCNN [80], NAD [118], ADHFR [205], CFC [78], Mittal <i>et al.</i> [163], ForensicFR [62], TDFL [233], E2EPG [315], CASPG [306], DualGAN [290], PS2-MAN [243], DTFS [317], Cascaded-FS [316], PTFs [310]
	Low-shot	Training and test with the data that has a small number of samples per identity.	SSPP-DAN [82], Guo <i>et al.</i> [72], Choe <i>et al.</i> [33], Hybrid Classifiers [269], Cheng <i>et al.</i> [31], DM [199], Yin <i>et al.</i> [292].
	Video-based	Performing face recognition with video sequences.	TBE-CNN [50], NAN [283], C-FAN [66], FANVFR [146], MARN [65], Rao <i>et al.</i> [178], CFR-CNN [171], ADRL [179], DAC [139]

images belong to the same identity. The face identification can be regarded as an extension of face verification, which aims to determine the specific identity of a face (*i.e.*, probe) among a set of identities (*i.e.*, gallery); moreover, in the case of open-set face identification, a prior task is needed, whose target is predicting whether the face belongs to one of the gallery identities or not.

For both the face verification and face identification, face representation is used to measure the similarity between face images. Therefore, how to learn discriminative face representation is the core target. With the advanced feature learning ability of DCNNs, face representation has made great progress. In the followings, we provide a systematic review of the learning methods of face representation from two major aspects, *i.e.*, network architecture and training supervision.

## 5.1 Network Architectures

The recent improvement of face representation partly benefits from the advance of deep architecture design. We first review the literature of network architecture for face representation learning. According to the designing purpose, we divide them into general architectures and specialized architectures. The general architectures are the basic and universal designs for common visual recognition tasks in the first place, and applied to face representation learning afterward. The specialized architectures include the modified or ensemble designs oriented to face recognition.

**5.1.1 General architectures.** With the advanced feature learning ability of DCNNs [34, 64, 77, 84, 112, 198, 217, 236], face representation has made great progress. Among them, AlexNet [112] obtains the first place in ImageNet [40] competition (ILSVRC) 2012 and achieves significant improvement compared with the traditional methods. Then, VGGNet [198] presents a more generic network, which replaces the large convolutional kernels by the stacked 3×3 ones, enabling the network to grow in depth. In order to enlarge the network without the extra increase of computational budget, GoogleNet [217] develops an inception architecture to concatenate the feature maps that

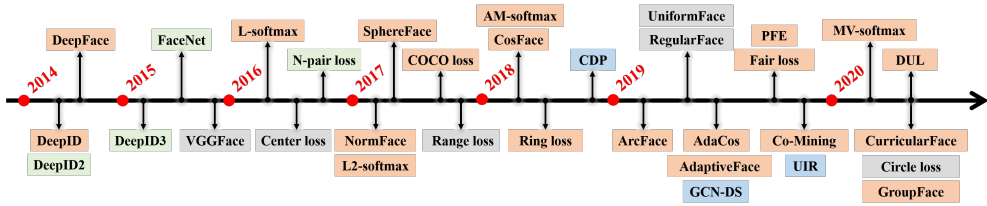


Fig. 8. The development of training supervision for face representation learning. The orange, green, gray and blue represent classification, feature embedding, hybrid, and semi-supervised methods, respectively. One can refer to Table 7 for the detailed references.

are generated by the convolutions of different receptive field. Soon, GoogleNet is applied to face representation learning, namely FaceNet [189]. More recently, ResNet [77] proposes a residual structure to make it possible for training deep networks that have hundreds of layers. ResNet is a modern network that has been widely used on many visual tasks, including face recognition. Additionally, several lightweight neural networks [83, 94, 151, 186, 328] are proposed to achieve the trade-off between speed and accuracy. All of them have been employed as backbone network for representation learning in the face recognition literature after being designed.

**5.1.2 Specialized architectures.** The aforementioned architectures are initially proposed for general visual tasks. Besides, many works develop specialized architectures for face representation learning. At first, many works [48, 209, 213, 214] attempt to assemble multiple convolution networks together for learning multiple local features from a set of facial patches. Given the human face appearing with regular arrangement of facial parts (eyes, nose, mouth, *etc*), such combination of multiple networks with respect to facial part can be more reliable than a single network. Besides, Xie *et al.* [274] design an end-to-end architecture, namely Comparator Network, to measure the similarity of two sets of a variable number of face images. Certain approaches [105, 106] develop feature-pair relational network to capture the relations between a pair of local appearance patches. More recently, FANFace [282] integrates the face representation network and facial landmark localization network, so that the heatmap of landmarks will boost the features for recognition.

In addition, many studies [30, 56, 154, 216, 265, 266] focus on developing the lightweight architecture. To reduce the parameters of deep networks, SparseNet [216] proposes to iteratively learn sparse structures from the previously learned dense models. Light-CNN [266] introduces a max-feature-map (MFM) activation function to gain better generalization ability than ReLU for face recognition; based on MFM, a lightweight architecture is developed that achieves the advantages in terms of runtime efficiency and model size. MobileFaceNet [30] replaces the global average pooling layer in the original MobileNet [186] with a global depth-wise convolution layer so the output feature can be improved by the spatial importance in the last layer.

It is worth noting that, in some landmark-free face alignment methods [76, 256, 263, 332, 340, 343] which have been presented in Section 4.2, the network can be optimized with respect to the objective of face representation learning and face alignment jointly.

## 5.2 Training Supervision

Besides network architectures, the training supervision also plays a key role for learning face representation. The objective of supervision for face representation learning is to encourage the faces of same identity to be close and those of different identities to be far apart in the feature space.

Following the convention of representation learning, we categorize the existing methods of training supervision for face representation into supervised scheme, semi-supervised scheme,

and unsupervised scheme. Although there are certain deep unsupervised learning methods [71, 129, 195, 255] for face clustering, in this review, we focus on the supervised and semi-supervised ones which comprise the major literature of state-of-the-art face recognition. Fig.8 shows the development of training supervision for face representation learning. In the supervised scheme, we can further categorize the existing works into three subsets, *i.e.*, classification, feature embedding and hybrid methods. The classification methods accomplish face representation learning with a  $N$ -way classification objective, regarding each of the  $N$  classes as an identity. The feature embedding methods aim to optimize the feature distance between samples with respect to the identity label, which means maximizing the inter-person distance and minimizing the intra-person distance. Besides, several works employ both classification and feature embedding routine to jointly train the network, namely hybrid methods. As for the semi-supervised scheme, several works exploit the labeled and unlabeled faces for representation learning.

**5.2.1 Classification scheme.** The classification based deep face representation learning is derived from the general object classification task. Each class corresponds to an identity that contains a number of faces of the same person. The softmax loss function is the most widely used supervision for classification task, which consists of a fully-connected (FC) layer, the softmax function and the cross-entropy loss. For face representation learning, DeepFace [220] and DeepID [214] are the pioneers of utilizing softmax to predict the probability over a large number of identities of training data. Their training loss function can be formulated as follows:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{W_{y_i}^T x_i + b_{y_i}}}{\sum_{j=1}^c e^{W_j^T x_i + b_j}}, \quad (4)$$

where  $N$  is the batch size,  $c$  is the number of classes (identities),  $y_i$  is the ground-truth label of sample  $x_i$ ,  $W_{y_i}$  is the ground-truth weight vector of sample  $x_i$  in the FC layer, and  $b_j$  is the bias term. The term inside the logarithm is the predicted probability on the ground-truth class. The training objective is to maximize this probability. Based on the softmax loss function, NormFace [237] and COCO loss [142] study the necessity of the normalization operation and apply  $L_2$  normalization constraint on both features and weights with omitting the bias term  $b_j$ . To effectively train with the normalized features, a scale factor is adopted to re-scale the cosine similarity between the features and the weights. Specifically, the normalized softmax loss function can be reformulated as

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{s \cos(\theta_{y_i})}}{e^{s \cos(\theta_{y_i})} + \sum_{j=1, j \neq y_i}^c e^{s \cos \theta_j}}, \quad (5)$$

where  $\cos(\theta_j)$  derives from the inner product  $W_j^T x_i$  with the  $L_2$  normalization on weights  $W_j = \frac{W_j}{\|W_j\|_2}$  and features  $x_i = \frac{x_i}{\|x_i\|_2}$ , and  $s$  is the scale parameter.

To further improve the intra-class compactness and inter-class separateness, L-softmax [138] replaces the ground-truth logit  $\cos(\theta_{y_i})$  with  $(-1)^k \cos(m\theta_{y_i}) - 2k$ ,  $\theta_{y_i} \in \left[ \frac{k\pi}{m}, \frac{(k+1)\pi}{m} \right]$ , where  $m$  is the angular margin that being a positive integer, and  $k$  is also an integer that  $k \in [0, m-1]$ . Similar to L-softmax, SphereFace [137] applies an angular margin in the ground-truth logit  $\cos(\theta_{y_i})$  to make the learned face representation to be more discriminative on a hypersphere manifold. However, the multiplicative angular margin in  $\cos(m\theta_{y_i})$  leads to potentially unstable convergence during the training. To overcome the problem, AM-softmax [235] and CosFace [239] present an additive margin penalty to the logit,  $\cos(\theta_{y_i}) + m_1$ , which brings more stable convergence. Subsequently, ArcFace [42] introduces an additive angular margin inside the cosine,  $\cos(\theta_{y_i} + m_2)$ , which corresponds to the geodesic distance margin penalty on a hypersphere manifold. The

following is a unified formulation of AM-softmax, CosFace, and ArcFace:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{s(\cos(\theta_{y_i} + m_2) + m_1)}}{e^{s(\cos(\theta_{y_i} + m_2) + m_1)} + \sum_{j=1, j \neq y_i}^c e^{s \cos \theta_j}}, \quad (6)$$

where  $m_1 < 0$  represents the additive cosine margin of AM-softmax and CosFace, and  $m_2 > 0$  denotes to the additive angular margin of ArcFace. They are easy to be implemented and can achieve better performance than the original softmax loss. Going further with the margin based supervision, AdaptiveFace [134] and Fair loss [130] propose the adaptive margin that being class-wise in the training data. The purpose is to address the imbalance distribution problem in the training dataset.

Resorting to the advantage of hard sample mining strategy [128, 197], some approaches [92, 251, 295] reformulate the negative (non-ground-truth) logit in softmax loss function. For example, MV-softmax [251] proposes to re-weight the negative logit to emphasize the supervision on the mis-classified samples, and thus to improve the representation learning from the negative view. In addition, certain studies [326, 327] deeply analyze the formulation of margin-based softmax loss function from the perspective of classification probability, and propose hyperparameter-free approaches for face representation learning.

More recently, many methods [22, 26, 86, 109, 148, 194, 196, 250, 270, 342] go further with the classification supervision for face representation learning. Some of them [86, 250, 342] focus on the noise-robust face representation learning, and some of the others [148, 270] tackle the issue of performance degradation of low-bit quantified model. Wu *et al.* [270] regard the quantization error as the combination of class error and individual error, and propose a rotation-consistent margin loss to reduce the latter error which is more critical. Besides, PFE [194] and DUL [26] propose to take into account the data uncertainty for modeling deep face representation, preventing from the uncertainty issue caused by low quality face images.

**5.2.2 Feature embedding scheme.** Feature embedding scheme aims to optimize the feature distance according to the label of sample pair. If the pair belong to the same identity, *i.e.*, positive pair, the objective is to minimize the distance or to maximize the similarity; otherwise, *i.e.*, negative pair, to maximize the distance or to minimize the similarity. For instance, contrastive loss [209, 211, 215, 289] direct optimizes the pair-wise distance with a margin that to encourage positive pairs to be close together and negative pairs to be far apart. The loss function to be minimized is written as

$$\mathcal{L}_c = \begin{cases} \frac{1}{2} \|f(x_i) - f(x_j)\|_2^2 & \text{if } y_i = y_j, \\ \frac{1}{2} \max(0, m_d - \|f(x_i) - f(x_j)\|_2)^2 & \text{if } y_i \neq y_j, \end{cases} \quad (7)$$

where  $y_i = y_j$  denotes  $x_i$  and  $x_j$  are positive pair,  $y_i \neq y_j$  denotes negative pair,  $f(\cdot)$  is the embedding function, and  $m_d$  is the non-negative distance margin. The contrastive loss drives the supervision on all the positive pairs and those negative pairs whose distance is smaller than the margin.

FaceNet [189] first applies the triplet loss [190, 257] to deep face representation learning. Different from contrastive loss, the triplet loss encourages the positive pairs to have smaller distance than the negative pairs with respect to a margin,

$$\mathcal{L}_t = \sum_i^N \left[ \|f(x_i^a) - f(x_i^p)\|_2^2 - \|f(x_i^a) - f(x_i^n)\|_2^2 + m_d \right]_+, \quad (8)$$

where  $m_d$  is the distance margin,  $x_i^a$  denotes the anchor sample,  $x_i^p$  and  $x_i^n$  refer to the positive sample and negative sample, respectively. The contrastive loss and triplet loss take into account only one negative sample each time, while negative pairs are abundant in training data and deserve

thorough involvement in training supervision. Therefore, N-pair loss [200] generalizes the triplet loss to the form with multiple negative pairs, and gained further improvement on face recognition.

Compared with the supervision of classification, feature embedding can save the parameters of the FC layer in softmax, especially when the training dataset is in large scale. But the batch size of training samples limits the effectiveness of feature embedding. To alleviate this problem, some approaches [152, 199, 203] propose the hard sample mining strategy to exploit the effective information in each batch, which is crucial to promote the performance of feature embedding.

*5.2.3 Hybrid methods.* The hybrid methods refer to those which apply classification and feature embedding together as the supervisory signals. DeepID series [209, 211, 211] utilize softmax loss and contrastive loss jointly for learning face representation. Later, several methods [45, 55, 259, 336] improve the feature embedding portion within the hybrid scheme, by utilizing either the intra-class or the inter-class constraints. Some methods [325, 342, 352] show the advantage for handling the long-tail distributed data which is a widely-existing problem in FR. Generally, the classification scheme works well on the head data but poorly on the tail data. Compared with classification scheme, the feature embedding scheme is able to provide the complementary supervision on the tail data. Thus, the combination of classification and feature embedding can improve the training on long-tail distributed data. More recently, Sun *et al.* [210] propose a circle loss from a unified perspective of the classification and embedding learning, which integrates the triplet loss with the cross-entropy loss to simultaneously learn deep features with pair-wise labels and class-wise labels.

*5.2.4 Semi-supervised scheme.* The aforementioned methods focus on supervised learning. Constructing labeled dataset requires much of annotation effort, while large amount of unlabeled data is easily available. Therefore, it is an attractive direction that to exploit the labeled and unlabeled data together for training deep models. For semi-supervised face representation learning, assuming the identities of unlabeled data being disjoint with the labeled data, several existing works [284, 285, 302] focus on generating the pseudo labels for unlabeled data. However, these methods assume non-overlapping identities between unlabeled and labeled data, which is generally impractical in real-world scenarios. Consequently, the unlabeled samples of overlapping identity will be incorrectly clustered as a new class by the pseudo-labeling methods. The intra-class label noise in pseudo-labeled data is another problem. To address these issues, RoyChowdhury *et al.* [183] separates unlabeled data into samples of disjoint and overlapping classes via an out-of-distribution detection algorithm. Besides, they design an improved training loss based on uncertainty to alleviate the label noise of pseudo-labeled data.

### 5.3 Specific Face Recognition Tasks

*5.3.1 Cross-domain face recognition.* Here, the term of cross-domain refers to a generalized definition that includes various factors, such like cross-age and cross-pose FR. As deep learning is a data-driven technique, the deep network usually works well on the training domains but poorly on the unseen ones. In real-world applications of face recognition, it is essential to improve the generalization ability of face representation across various domain factors. In the following, we discuss certain aspects of cross-domain FR; also, the current solutions are presented.

**Cross-age:** As the facial appearance has large intra-class variation along with the growing age, identifying faces across wide range of age is a challenging task. For such cross-age FR, there are two directions. The first direction [54, 228, 238, 252, 258, 276] aims to learn age-invariant face representation by decomposing deep face features into age-related and identity-related components. The second direction is based on generative mechanism. In this way, several methods [6, 108, 248] attempt to synthesize faces of target age, but they present imperfect preservation of the original



identities in aged faces. Thus, supplementary methods [5, 204, 254, 333] are designed to improve the identity-preserving ability during the face aging.

**Cross-pose:** In unconstrained conditions, such as surveillance video, the cameras cannot always capture the frontal face image for every appeared subject. Thus, the captured faces have large pose variation from frontal to profile view. However, generating the frontal faces will increase the burden of face recognition system. Cao *et al.* [23] alleviate this issue by transforming the representation of a profile face to the frontal view in the feature space. Another problem is that the number of profile faces are much fewer than frontal ones in the training data. Thus, some generative approaches [43, 87, 225, 335] propose to synthesize identity-preserving faces of arbitrary poses to enrich the training data. Moreover, certain methods [1, 104, 155] develop multiple pose-specific deep models to compute the multi-view face representations.

**Racial bias:** Due to the imbalance distribution of different races in training data, the deep face feature shows favorable recognition performance to the races of large proportion in training data than those of small proportion. Recently, Wang *et al.* [246] construct an in-the-wild face dataset (RFW) with both identity and race annotation, which consists of four racial subsets, *i.e.*, Caucasian, Asian, Indian, and African. Besides, they propose a domain adaptation method to alleviate the racial bias. Later on, RL-RBN [245] sets a fixed margin for the large-proportion races and automatically select an optimal margin for the small-proportion races, in order to achieve balanced performance.

**Cross-modality:** Cross-modality face recognition generally refers to the heterogeneous face recognition, which performs with a pair of input face images captured by different sensing modalities, such as infrared vs. visible, or sketch vs. photo. How to alleviate the domain gaps between different modalities is the major challenge. Besides, the available infrared or sketch images are of very limited number. The existing works mainly handle these two issues. Many methods [62, 140, 163, 180, 188, 233] exploit the transfer learning, *i.e.*, pretraining on the visible-light (VIS) images and finetuning with the infrared or sketch data, to reduce the domain discrepancy. Another set of methods [47, 79, 80, 267] decompose the cross-modality features to the modality-specific and modality-invariant components, and use the latter one for the recognition task. Moreover, recent methods [78, 118, 205, 243, 290, 310, 316, 317, 349] aim to synthesize the common VIS image from infrared or sketch input, and then perform the regular FR in the VIS domain.

**5.3.2 Low-shot face recognition.** Low-shot learning in face recognition focuses on the condition of identification of low-shot face IDs, each of which has a small number of samples. MS-Celeb-1M low-shot learning benchmark [72] is most used, which has about 50 to 100 training samples per ID in the base set and only one training sample per ID in the novel set. The target is to recognize the IDs in both base and novel sets. The key challenge is to correctly recognize the subjects in the novel set which has only one training sample per ID. To tackle this problem, many methods [31, 72, 199, 269, 292] improve the low-shot face recognition with better training supervision or strategy. Besides, face generation [33, 82, 206] is another effective routine for low-shot issue.

**5.3.3 Video face recognition.** The above methods focus on still image-based face recognition. For video face recognition, a common way [28, 50] is to equally consider the importance of each frame and simply average a set of deep features as the template. However, this routine does not consider the different quality of frames and the temporal information across frames. How to obtain an optimal template feature in video is the major challenge of video face recognition. Several methods [65, 66, 146, 283] aggregate the frame-level features with the attention weights or quality scores. Synthesizing representative or high-quality face image from a video sequence is another possibility [171, 178]. Additionally, certain methods [139, 157, 179] model the temporal-spatial information with the attention mechanism and find the focus of video frames.

## 5.4 Evaluation Metrics and Datasets

**5.4.1 Metrics.** The performance of face recognition is usually evaluated on two tasks: verification and identification, each of which has its corresponding evaluation metrics. Specifically, two sets of samples, *i.e.*, gallery and probe, are required for the evaluation. The gallery refers to a set of faces registered in the face recognition system with known identities, while the probe denotes a set of faces need to be recognized in verification or identification. Before discussing the commonly used evaluation metrics, we first introduce some basic concepts. A face recognition system determines whether to accept the matching of a probe face and a gallery face by comparing their similarity, computed by some measurement between their features, with a given threshold. Specifically, when a probe face and a gallery face are the same identity, a true acceptance (TA) means their similarity is above the threshold, and a false rejection (FR) represents their similarity is below the threshold; if they are different identities, a true rejection (TR) means their similarity is below the threshold, and a false acceptance (FA) means their similarity is above the threshold. These are the basic concepts to build the evaluation metrics in the followings. One can refer to [68, 69] for more details.

**Verification task:** Face verification is often applied in identity authentication system, which measures the similarity of face pairs. One presents his or her face and claims the enrolled identity in the gallery. Then, the system determines whether it accepts the person being the same one of the claimed identity by calculating the similarity between the presented face and the claimed face. Thus, the verification task can be regarded as a one-to-one face matching process. The false accept rate (FAR) and true accept rate (TAR) are used to evaluate the verification performance. FAR is the fraction of impostor pairs with the similarity above the threshold, which can be calculated by  $\frac{FA}{FA+TR}$ ; TAR represents the fraction of genuine pairs with the similarity above the threshold, which can be calculated by  $\frac{TA}{TA+FR}$ . Then, by varying the threshold, the ROC curve can be drawn by many operating points, each of which is determined by a pair of TAR vs. FAR. The ROC curve (with TAR value at selected FAR) and its AUC (*i.e.*, area under curve) are widely used to evaluate the performance for the face verification task.

**Identification task:** Face identification task determines whether a probe face belongs to a enrolled identity in the gallery set. To this end, the probe face needs to be compared with every person in the gallery set. Thus, the identification task can be also referred as one-to- $N$  face matching.

Generally, face identification includes two tasks, *i.e.*, the open-set and closed-set identification. The open-set identification task refers to that the probe face is not necessarily the very identity contained in the gallery set, which is the most general case in practice. The true positive identification rate (TPIR) and false positive identification rate (FPIR) are the most used metrics for the following two situations. The first situation refers to that the probe corresponds to an enrolled identity in the gallery set. This situation is called mate searching, and the probe is called mate probe. The succeeded mate searching represents that the rank of true matching is higher than the target rank, and meanwhile its similarity is above the threshold. In such case, the mate probe is correctly identified as its true identity, and the mate searching is measured by the TPIR which represents the proportion of succeeded trials of mate searching. The second is non-mate searching, in which the probe does not correspond to any enrolled identity (*i.e.*, non-mate probe). The non-mate searching is measured by the FPIR which reports the proportion of non-mate probes wrongly identified as enrolled identity. By fixing the rank and varying the threshold, the ROC curve can be drawn by many operating points, each of which is determined by a pair of TPIR vs. FPIR. The ROC curve (TPIR value at a given FPIR) is used to evaluate performance in the open-set face identification task.

In the closed-set scenario, the identity of each probe face is included in the gallery set. The cumulative match characteristic (CMC) curve is used for evaluating the closed-set face identification. The CMC curve is drawn by the operating points that are determined by a pair of identification

Table 8. The commonly used public datasets for training and testing deep face recognition.

Dataset	Year	# Subject	# Image/Video	# of Img/Vid per Subj	Description
Training					
CASIA-WebFace [289]	2014	10,575	494,414/-	47	The first public large-scale face dataset
VGGFace [172]	2015	2,622	2.6M/-	1,000	Containing large number of images in each subject
CelebA [297]	2015	10,177	202,599/-	20	Rich annotations of attributes and identities
UMDFaces [11]	2015	8,277	367K/-	45	Abundant variation of facial pose
MS-Celeb-1M [73]	2016	100K	10M/-	100	A large-scale public dataset of celebrity faces
MegaFace [107, 167]	2016	672,057	4.7M/-	7	A long-tail dataset of non-celebrity
VGGFace2 [24]	2017	9,131	3.31M/-	363	A high-quality dataset with a wide range of variation
UMDFaces-Videos [10]	2017	3,107	-/22,075	7	A video training dataset collected from YouTube
MS-Celeb-1M Low-shot [72]	2017	20K,1K	1M,1K/-	58,1	Low-shot face recognition
IMDb-Face [234]	2018	57K	1.7M/-	29	A large-scale noise-controlled dataset
QMUL-SurvFace [234]	2018	5,319	220,890/-	41	A low-resolution surveillance dataset
Glint360K [4]	2021	360K	17M/-	47	A large-scale and cleaned dataset
WebFace260M [354]	2021	4M	260M/-	65	The largest public dataset of celebrity faces
Test					
LFW [88]	2007	5,749	13,233/-	2.3	A classic benchmark in unconstrained conditions
YouTube Faces (YTF) [262]	2011	1,595	-/3,425	2.1	Face recognition in unconstrained videos
CUFSS [324]	2011	1,194	2,388/-	2	Photo-sketch face recognition
CASIA NIR-VIS v2.0 [122]	2013	725	17,580/-	24.2	Near-infrared vs. RGB face recognition
IJB-A [110]	2015	500	5,712/2,085	11.4/4.2	Set-based face recognition with large variation
CFP [191]	2016	500	7,000/-	14	Frontal to profile cross-pose face verification
MS-Celeb-1M Low-shot [72]	2016	20K,1K	100K,20K/-	5,20	Low-shot face recognition
MegaFace [107, 167]	2016	690,572	1M/-	1.4	A large-scale benchmark with one million faces
IJB-B [260]	2017	1,845	11,754/7,011	6.37/3.8	Set-based face recognition with full pose variation
CALFW [338]	2017	4,025	12,174/-	3	Cross-age face verification
AgeDB [164]	2017	570	16,516/-	29	Cross-age face verification
SLLFW [46]	2017	5,749	13,233/-	2.3	Improving the difficulty of negative pairs in LFW
CPLFW [337]	2017	3,968	11,652/-	2.9	Cross-pose face verification
Trillion Pairs [39]	2018	1M	1.58M/-	1.6	A large-scale benchmark with massive distractors
IJB-C [156]	2018	3,531	31,334/11,779	6/3	Set-based face recognition with large variation
IJB-S [103]	2018	202	5,656/552	28/12	Real-world surveillance videos
RFW [246]	2018	11,429	40,607/-	3.6	For reducing racial bias in face recognition
DFW [115]	2018	600	7,771/-	13	Disguised face recognition
QMUL-SurvFace [234]	2018	10,254	242,617/-	23.7	Low-resolution surveillance videos

rate vs. rank. The identification rate refers to the fraction of probe faces that are correctly identified as the true identities, thus the CMC curve reports the fraction of the true matching with a given rank, and the identification rate at rank one is the most commonly used indicator of performance. It is noteworthy that the CMC is a special case of the TPIR when we relax the threshold.

**5.4.2 Datasets.** With the development of deep face recognition, another key role to promote face representation learning is the growing datasets for training and test. In the past few years, the face datasets have become large scale and diverse, and the testing scene has been approaching to the real-world unconstrained condition. The statistics of them are presented in Table 8.

**Training data:** Large-scale training datasets are essential for learning deep face representation. The early works often employ the private face datasets, such as Deepface [220], FaceNet [189], DeepID [209]. To make it possible for fair comparison, Yi *et al.* [289] release the CASIA-WebFace dataset, which has been one of the most widely-used training datasets. Afterward, more public training datasets are published to provide abundant face images for training deep face model. Among them, VGGFace [172] and VGGFace2 [24] contain many training samples for each subject. In contrast, MS-Celeb-1M [73], MegaFace [107], IMDb-Face [234] and WebFace260M [354] provide a large number of subjects with relatively less training samples per subject.

**Test data:** As for testing, Labeled Faces in the Wild (LFW) [88] is classic and the most widely used benchmark for face recognition in unconstrained environments. The original protocol of LFW contains 3,000 genuine and 3,000 impostor face pairs, and evaluates the mean accuracy of

Table 9. Performance (%) comparison of face recognition on various test datasets. “Training Data” denotes the number of training face images used by the methods. For the evaluation on MegaFace, “Id.” refers to the rank-1 face identification accuracy with 1M distractors, and “Veri.” refers to the face verification TAR at  $10^{-6}$  FAR. For the evaluation on IJB-B and IJB-C, we report the 1:1 verification TAR ( $@FAR=10^{-4}$ ). The performance with “\*” refers to the evaluation on the refined version of MegaFace [42]. “-” indicates that the authors do not report the performance with the corresponding protocol.

Method	Publication	Subcategory	Training Data	Backbone	LFW	MegaFace		IJB-B	IJB-C	YTF	CALFW	CPLFW	CFP-FP	AgeDB30
						Id.	Veri.							
DeepFace [220]	CVPR'14	Classification	4M	CNN-8	97.35	-	-	-	-	91.4	-	-	-	-
DeepID [220]	CVPR'14	Classification	0.3M	CNN-8	97.45	-	-	-	-	-	-	-	-	-
L-Softmax [138]	ICML'16	Classification	0.5M	VGGNet-18	99.10	67.12	80.42	-	-	-	-	-	-	-
NormFace [237]	ACMMM'17	Classification	0.5M	ResNet-28	99.16	-	-	-	-	-	-	-	-	-
SphereFace [137]	CVPR'17	Classification	0.5M	ResNet-64	99.42	72.72	85.56	-	-	95.0	-	-	-	-
ReST [137]	CVPR'17	Classification	0.5M	CNN-9	99.03	65.16	-	-	-	95.4	-	-	-	-
E2e [340]	SPL'17	Classification	0.7M	ResNet-27	99.33	-	-	-	-	95.0	-	-	-	-
AM-softmax [235]	SPL'18	Classification	0.5M	ResNet-20	98.98	72.47	84.44	-	-	-	-	-	-	-
CosFace [239]	CVPR'18	Classification	5M	ResNet-64	99.73	82.72	96.65	-	-	97.6	-	-	-	-
ComparatorNet [274]	ECCV'18	Classification	3.3M	ResNet-50	99.73	-	-	84.1	88.0	-	-	-	-	-
ArcFace [42]	CVPR'19	Classification	0.5M	ResNet-50	99.53	77.50	92.34	-	-	-	-	-	95.56	95.15
Fair loss [130]	ICCV'19	Classification	0.5M	ResNet-50	99.57	77.45	92.87	-	-	96.2	-	-	-	-
PFE [194]	ICCV'19	Classification	4.4M	ResNet-64	99.82	<b>78.95</b>	92.51	-	93.25	-	-	-	93.34	-
FANFace [282]	AAAI'20	Classification	0.5M	ResNet-50	99.56	78.32	92.83	-	-	96.72	-	-	-	-
TURL [196]	CVPR'20	Classification	4.8M	ResNet-100	99.78	78.60	<b>95.04</b>	-	96.6	-	-	-	98.64	-
RDCFace [332]	CVPR'20	Classification	1.7M	ResNet-50	99.80	-	-	-	-	97.10	-	-	96.62	-
AdaCos [326]	CVPR'19	Classification	2.35M	ResNet-50	99.73	97.41*	-	-	92.4	-	-	-	-	-
P2SGrad [327]	CVPR'19	Classification	2.35M	ResNet-50	99.82	97.25*	-	-	92.3	-	-	-	-	-
AdaptiveFace [134]	CVPR'19	Classification	5M	ResNet-50	99.62	95.02*	95.61*	-	-	-	-	-	-	-
ArcFace [42]	CVPR'19	Classification	5.8M	ResNet-100	99.82	98.35*	98.48*	94.2	95.6	97.7	95.45	92.08	98.27	98.15
MV-AM-softmax [251]	AAAI'20	Classification	3.2M	Attention-56	99.79	98.00*	98.31*	-	-	-	95.63	89.19	95.30	98.00
DUL [26]	CVPR'20	Classification	3.6M	ResNet-64	99.83	98.12*	-	-	94.21	96.84	-	-	<b>98.78</b>	-
DB [22]	CVPR'20	Classification	5.8M	ResNet-50	99.78	96.35*	96.56*	-	-	96.08	92.63	-	97.90	-
CurricularFace [92]	CVPR'20	Classification	5.8M	ResNet-100	99.80	98.71*	98.64*	94.8	96.1	-	<b>96.20</b>	93.13	98.37	<b>98.32</b>
GroupFace [109]	CVPR'20	Classification	5.8M	ResNet-100	<b>99.85</b>	<b>98.74*</b>	<b>98.79*</b>	<b>94.93</b>	<b>96.26</b>	<b>97.8</b>	<b>96.20</b>	<b>93.17</b>	98.63	98.28
FaceNet [189]	CVPR'15	Embedding	400M	GoogleNet-22	99.63	-	-	-	-	95.1	-	-	-	-
VGG Face [172]	BMVC'15	Embedding	2.6M	CNN-36	98.95	<b>64.79</b>	<b>78.32</b>	-	-	<b>97.3</b>	-	-	-	-
N-pair loss [200]	NIPS'16	Embedding	0.5M	CNN-10	98.50	-	-	-	-	-	-	-	-	-
GridFace [343]	ECCV'18	Embedding	10M	GoLeNet-22	<b>99.70</b>	-	-	-	-	95.6	-	-	-	-
DeepID2 [209]	NeurIPS'14	Hybrid	0.3M	CNN-8	99.15	65.21	78.86	-	-	-	-	-	-	-
SparseNet [216]	CVPR'15	Hybrid	0.3M	CNN-15	99.30	-	-	-	-	92.7	-	-	-	-
Center loss [259]	ECCV'16	Hybrid	0.7M	CNN-11	99.28	65.49	80.14	-	-	94.9	-	-	-	-
Ring loss [339]	CVPR'18	Hybrid	3.5M	ResNet-64	99.50	74.93	-	-	-	93.7	-	-	-	-
PRN [106]	ECCV'18	Hybrid	2.8M	ResNet-101	99.76	-	-	84.5	-	96.3	-	-	-	-
RegularFace [336]	CVPR'19	Hybrid	3.1M	ResNet-20	99.61	75.61	91.13	-	-	96.7	-	-	-	-
UniformFace [55]	CVPR'19	Hybrid	3.8M	ResNet-34	99.8	<b>79.98</b>	<b>95.36</b>	-	-	<b>97.7</b>	-	-	-	-
AFRN [105]	ICCV'19	Hybrid	2.8M	ResNet-101	<b>99.85</b>	-	-	<b>88.5</b>	93.0	97.1	<b>96.30</b>	<b>93.48</b>	95.56	<b>96.35</b>
Circle loss [210]	CVPR'20	Hybrid	3.6M	ResNet-34	99.73	<b>97.81*</b>	-	-	<b>93.44</b>	96.38	-	-	<b>96.02</b>	-

verification on these 6,000 pairs. So far, the state-of-the-art accuracy has been saturated on LFW, whereas the total samples in LFW are more than those in the original protocol. Based on this, BLUFR [126] exploits all the face images in LFW for a large-scale unconstrained face recognition evaluation; SLLFW [46] replaces the negative pairs of LFW with more challenging ones. In addition, CFP [191], CPLFW [337], CALFW [338], AgeDB [164] and RFW [246] utilize the similar evaluation metric of LFW to test face recognition with various challenges, such as cross pose, cross age and multiple races. MegaFace [107, 167] and Trillion Pairs [39] focus on the performance at the strict false accept rates (*i.e.*,  $10^{-6}$  and  $10^{-9}$ ) on face verification and identification with million-scale distractors. The above datasets focus on image-to-image face recognition, whereas YTF [262], IJB series [103, 110, 156, 260], and QMUL-SurvFace [234] serve as the evaluation benchmark of video-based face recognition. Especially, IJB-S and QMUL-SurvFace are constructed from real-world surveillance videos, which are much more difficult and realistic than the tasks on still images.

## 5.5 Performance Comparison

Table 9 shows the performance of face representation methods on various test datasets. Among them, CosFace [239] and ArcFace [42] are the two commonly used methods in many applications of face recognition. In addition, with the growing datasets for training and test, the closed-set

Table 10. Summary of the major challenges towards end-to-end deep face recognition.

Challenges		Description
The issues of each element.	Face detection	<ul style="list-style-type: none"> <li>• Trade-off between detection accuracy and efficiency.</li> <li>• Accuracy of the bounding box location.</li> <li>• Detecting faces with a wide range of scale.</li> </ul>
	Face alignment	<ul style="list-style-type: none"> <li>• Annotation ambiguity and granularity.</li> </ul>
	Face representation	<ul style="list-style-type: none"> <li>• limited training data and computational budget.</li> <li>• Surveillance video face recognition.</li> <li>• Noisy label and imbalance data.</li> </ul>
The common issues across the elements.	Facial / image variations	<ul style="list-style-type: none"> <li>• Large pose, extreme expression, occlusion, facial scale.</li> <li>• Motion blur, low illumination, low resolution.</li> </ul>
	Data / label distribution	<ul style="list-style-type: none"> <li>• Limited labeled data, label noise.</li> <li>• Usage of unlabeled data.</li> <li>• Imbalance over scale, identity, race, domain, modality.</li> </ul>
	Computational efficiency	<ul style="list-style-type: none"> <li>• Inference on non-GPU device and edge computing.</li> <li>• Fast training and convergence.</li> </ul>
The issues concerning to the entire system.	Interpretability	<ul style="list-style-type: none"> <li>• Explainable learning and inference.</li> </ul>
	Joint modeling and optimization	<ul style="list-style-type: none"> <li>• End-to-end training and inference.</li> <li>• Unified learning objective.</li> <li>• Mutual promotion.</li> </ul>
	Universal pretraining	<ul style="list-style-type: none"> <li>• Universal pretrained facial representation.</li> </ul>
	Trustworthiness	<ul style="list-style-type: none"> <li>• Robustness, fairness, explainability, security, and privacy.</li> </ul>

classification training on the large-scale datasets enables to approach the open-set face recognition scenario. This could be the reason why the classification based training methods have been widely studied and dominated the state-of-the-art performance in recent years. One can find the publication trend of three supervised training schemes with the increasing scale of public face datasets in the supplementary material.

## 6 DISCUSSION AND CONCLUSION

Deep face recognition still remains a number of issues for each element. In the following, we first analyze the major challenges towards end-to-end deep face recognition and the subcategories of each element. Then, we provide a detailed discussion about the promising future trends for each element and the entire system. Finally, the conclusion of this survey is presented.

### 6.1 Challenge

The top rows of Table 10 elaborate the issues of each element. For face detection, the state-of-the-art methods are eager for trade-off between detection accuracy and efficiency. For example, in many applications, resizing the input image is a common practice of acceleration for detectors, while it harms the recall of tiny faces as well. In the unconstrained condition, human faces with large variation tend to be missed by detectors, whereas the diverse image background often leads to false positives. Besides, detecting faces with a wide range of scale is also a great challenge. As for the face alignment procedure, the facial landmark localization methods are still not robust enough when working with extreme variations, such as severe occlusion, large pose, low illumination. In addition, the annotation ambiguity, such as the landmarks on cheek, is a common problem in datasets. Besides, most of the existing facial landmark datasets provide the annotation of 68 or 106 points. More landmark points enable to depict the abundant facial structure. For the face representation learning, although existing methods achieve high accuracy on various benchmarks, it is still challenging when training data and computational budget are very limited. In addition, surveillance face recognition is a common scenario, where the challenges include various facial

variations, such large poses, motion blur, occlusion, low illumination and resolution, *etc.* Imbalance distribution of training data also brings issues to the face representation learning, such as long-tail distribution over face identities or domains.

The middle rows of Table 10 elaborate the common issues shared between face detection, alignment and representation. We can find that the issues mainly include three aspects, *i.e.*, facial and image variations, data and label distribution, and computational efficiency. For example, in the first aspect, the facial variations include large facial pose, extreme expression, occlusion and facial scale, while the image variations include the objective factors such as motion blur, low illumination and resolution which occur frequently in video face recognition. Another example indicates the need of training efficiency, including fast training and convergence, both of which devote to accelerating the learning of large face representation network (hundreds of layers normally) from weeks to hours; the former generally focuses on the mixed precision training or the distributed framework for large-scale training (over millions of identities), while the latter focuses on improving the supervision, initialization, updating manner, activation, architectures, *etc.* Here, rather than replaying every detail, we leave Table 10 to readers for exploring the common challenges and further improvement. It is worth mentioning that all the elements will benefit from the solutions against these issues, since they are the common issues across the elements.

The bottom of Table 10 indicates the major challenges from the perspective of entire system. For instance, ideally, the three elements should be jointly modeled and optimized with respect to the end-to-end accuracy. On the one hand, such integration provides a possibility to search global optimal solution for the holistic system; on the other hand, the individual elements of the system can benefit from the upstream ones. However, the elements have different learning objectives regarding to their own tasks. How to unify these learning objectives is a challenging and critical issue for the joint optimization. One can find a group of works [41, 76, 256, 263, 314, 332, 334, 340, 343] attempting to integrate face detection and alignment, or face alignment and representation for a joint boost. But face detection is still difficult to be integrated with face representation because they have quite different objectives and implementation mechanisms.

In addition, we are going deeper with Table 11 about the major challenges towards the subcategories of each element. For instance, since the anchor-based face detector needs to pre-define a large number of anchors, the settings of preset anchors need to be carefully tuned for each particular dataset, which limits the generalization ability of face detectors. In contrast, anchor-free face detector needs further exploration for better robustness to false positives and stability in training process.

## 6.2 Future Trend

To address the above challenges, a number of worthwhile research directions need to be explored in the future.

### 6.2.1 Face detection.

- **Generalized anchor settings.** The existing anchor-based methods design the anchor setting from many aspects, such as assignment and matching strategy [120, 125, 145, 221, 322], attributes tuning [32, 321, 347], and sampling strategy [161]. The well-tuned anchors may limit the generalization ability of face detectors. Hence, it is worth to explore a generalized anchor setting that can be used for different application demand.
- **Anchor-free face detection framework.** Anchor-free detectors [116, 224, 346] show flexible designs and more potential in generalization ability for object detection. However, a small number of works [89, 279, 294] have explored the anchor-free mechanism and its advantages for face detection.

Table 11. Summary of the major challenges towards the subcategories for each element.

Element	Subcategory	Challenges Description
Face detection	Multi-stage	• Runtime efficiency.
	Single-stage	• Detecting tiny faces.
	Anchor-based	• Well-tuned anchors.
	Anchor-free	• Training stability and robustness to false positives.
	CPU real-time	• Trade-off between accuracy and efficiency.
	Multi-task learning	• Balance of multi-task training supervision.
	Problem-oriented	• Low-illumination and low-resolution.
Face alignment	Landmark-based – Coordinate regression	• Prediction bias due to poor initialization.
	Landmark-based – Heatmap regression	• High computational cost.
	Landmark-based – 3D model fitting	• Runtime efficiency.
	Landmark-free	• Loss of identity discriminative information.
Face representation	Training supervision – Classification	• Training on imbalance data.
	Training supervision – Feature embedding	• Efficient training on large-scale datasets.
	Training supervision – Hybrid	• Unified training supervision of classification and feature embedding.
	Training supervision – Semi-supervised	• Open-set identities setting.
	Specific Tasks – Cross-age	• Recognizing identities across a wide range of age.
	Specific Tasks – Cross-pose	• Large pose variation.
	Specific Tasks – Racial bias	• Bias reduction.
	Specific Tasks – Cross-modality	• Domain generalization.
	Specific Tasks – Low-shot	• One-shot learning.
	Specific Tasks – Video-based	• Low quality of frames.

### 6.2.2 Face alignment.

- **High robustness and efficiency.** There is a large amount of facial variations in real-world conditions, which requires the alignment methods being robust to various input faces while keeping efficiency as an intermediate step of the system.
- **Dense landmark localization.** The most existing datasets employ 68 or 106 keypoints as annotation configuration. They are enough for face alignment (usually 5 keypoints needed), but insufficient to the complex face analysis tasks, such as facial motion capture. Besides, the dense landmarks will help to locate more accurate alignment-needed keypoints.
- **Video-based landmark localization.** How to make better use of the temporal information is a major challenge for video-based landmark localization. This topic will enable to address the problems in video, such as large poses, motion blur, low illumination and resolution, *etc.*
- **Semi-supervised landmark localization:** The extensive research on landmark localization belongs to the regime of supervised learning, which needs the precise annotated landmarks. However, it is expensive and inefficient to obtain large-scale dataset with the precise annotations. As explored by the pioneering works [52, 53, 81, 182], the semi-supervised routine is a feasible and valuable solution for facial landmark localization.

### 6.2.3 Face representation.

- **Lightweight face recognition:** The large memory and computational cost often makes it impractical to employ heavy-weight networks on mobile or embedded devices. Although many works [30, 56, 154, 265, 266, 270] have studied lightweight face recognition, there is still large room to improve the lightweight models with high efficiency and accuracy.

- **Robustness to variations in video:** It highly requires robust face representation models against varying conditions in surveillance video. The robustness to low image quality and large facial pose is the core demand in many practical applications.
- **Noisy label learning:** Label noise is an inevitable problem when collecting large-scale face dataset. Certain works [39, 42, 234, 329] study how to remove the noisy data, and some others [86, 250, 342] aim at learning noise-robust face representation. But most of them are susceptible to the ability of the initial model, and need to be more flexible in real-world scenarios. It is still an open issue for noisy label learning in face recognition.
- **Cross domain face recognition:** There are many different domain factors in face data, such as facial age, pose, race, imaging modality, and some works [23, 54, 225, 245, 246, 258, 316] have studied the face recognition across a small fraction of them. How to obtain a universal representation for cross domain face recognition is a challenging research topic.
- **Learning with imbalance data:** Representation learning on the long-tail data is long-standing topic in many datasets. With the scarcity of intra-class variations, the subjects with limited training samples are usually neglected. The domain bias caused by imbalance data scale is another problem. It is worth to handle these problems in a unified framework.
- **Learning with unlabeled faces:** There are a large amount of unlabeled face data in practical applications. However, it is excessively expensive to manually annotate them when the dataset keeps growing. Recently, semi-supervised learning and face clustering methods attract increasing attention. How to effectively employ unlabeled data for boosting face recognition is a promising direction.

6.2.4 *Towards the entire system.* There is very little work to solve the major challenges from the perspective of entire system. We present several promising directions of this area in the following.

- **Interpretable deep models:** Although the explainable artificial intelligence, so-called XAI, has been studied for a long time, the explainable deep face recognition is in its infancy [261, 291, 300, 341]. There are two ways to access the interpretability for deep face recognition, *i.e.*, the top-down and bottom-up, respectively. The top-down way resorts to the human prior knowledge for algorithm exploration, since human shows superior ability of face recognition than deep models in many tough conditions. The bottom-up way denotes the exploration from the perspective of face data itself, such as modeling the explainable deep face recognition in spatial and scale dimension.
- **Joint modeling for the holistic system:** Despite the three elements having different optimized objective, it is still worth to exploit the end-to-end trainable deep face recognition, and study how they can be further improved through the jointly learning. Furthermore, beyond the topic of this survey, there is also an open question that how should we develop a single network to perform the end-to-end face recognition.
- **Universal face representation pretraining:** Most studies of face recognition focus on the specific tasks, but overlooking how to learn a pre-trained universal face representation that can be used to facilitate the downstream facial analysis tasks. There is only one work [18] that studies this topic. The findings show that it is promising to obtain significant performance improvement for related facial tasks by employing unsupervised pretraining.
- **Trustworthy face recognition system:** With the wide application, it is important to evaluate and boost the trustworthiness of the recognition system [96]. The pursuit for trustworthy face recognition system is becoming a necessity, which mainly involves several aspects, *i.e.*, robustness, fairness, interpretability, security, and privacy. Further research on these aspects is essential.



### 6.3 Conclusion

In this survey, we review the recent advances of the elements of end-to-end deep face recognition, which consist of face detection, face alignment and face representation. Although there are many surveys about face recognition, they mostly focus on face representation without considering the intrinsic connection from other elements in the pipeline; whereas, this survey is the first one which provides a comprehensive review of the elements of end-to-end deep face recognition. We present a detailed discussion and comparison of many approaches in each element from poly-aspects. Also, we discuss the relationship between the elements and the holistic framework. According to these elaborated contents, we can not only find the suitable methods to establish state-of-the-art face recognition system, but also know which method is quite strong-baseline style for comparison in experiment. Additionally, we analyze the existing challenges and collect certain promising future research directions. We hope this survey could bring helpful thoughts for better understanding of end-to-end face recognition and deeper exploration in a systematic way.

### REFERENCES

- [1] W. Abd-Almageed, Y. Wu, S. Rawls, S. Harel, T. Hassner, I. Masi, J. Choi, J. T. Leksut, J. Kim, P. Natarajan, R. Nevatia, and G. G. Medioni. 2016. Face recognition using deep multi-pose representations. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*. 1–9.
- [2] I. Adjabi, A. Ouahabi, A. Benzaoui, and A. Taleb-Ahmed. 2020. Past, Present, and Future of Face Recognition: A Review. *Electronics* 9, 8 (2020).
- [3] T. Ahonen, A. Hadid, and M. Pietikinen. 2004. Face recognition with local binary patterns. In *Proceedings of the European Conference on Computer Vision*. 469–481.
- [4] Xiang An, Xuhan Zhu, Yanghua Xiao, Lan Wu, Ming Zhang, Yuan Gao, Bin Qin, Debing Zhang, and Yingnan Fu. 2021. Partial FC: Training 10 Million Identities on a Single Machine. *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (2021)*, 1445–1449.
- [5] G. Antipov, M. Baccouche, and J. Dugelay. 2017. Boosting cross-age face verification via generative age normalization. In *Proceedings of the IEEE International Joint Conference on Biometrics*. 191–199.
- [6] G. Antipov, M. Baccouche, and J. Dugelay. 2017. Face aging with conditional generative adversarial networks. In *Proceedings of the IEEE International Conference on Image Processing*. 2089–2093.
- [7] H. R. Arabnia. 2009. A Survey of Face Recognition Techniques. *Journal of Information Processing Systems* 5, 2 (2009), 41–68.
- [8] B. Yang, J. Yan, Z. Lei, and S. Z. Li. 2014. Aggregate channel features for multi-view face detection. In *Proceedings of the IEEE International Joint Conference on Biometrics*. 1–8.
- [9] Y. Bai, Y. Zhang, M. Ding, and B. Ghanem. 2018. Finding tiny faces in the wild with generative adversarial network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 21–30.
- [10] A. Bansal, C. D. Castillo, R. Ranjan, and R. Chellappa. 2017. The Do’s and Don’ts for CNN-Based Face Verification. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 2545–2554.
- [11] A. Bansal, A. Nanduri, C. D. Castillo, R. Ranjan, and R. Chellappa. 2017. UMDFaces: An annotated face dataset for training deep networks. In *Proceedings of the IEEE International Joint Conference on Biometrics*. 464–473.
- [12] P. N. Belhumeur, D. W. Jacobs, D. J. Kriegman, and N. Kumar. 2013. Localizing parts of faces using a consensus of exemplars. *IEEE Trans. Pattern Anal. Mach. Intell.* 35, 12 (2013), 2930–2940.
- [13] P. N. Belhumeur, P. H. Joo, and D. J. Kriegman. 1997. Eigenfaces vs. Fisherfaces: Recognition Using Class Specific Linear Projection. *IEEE Trans. Pattern Anal. Mach. Intell.* 19, 7 (1997), 711–720.
- [14] C. Bhagavatula, C. Zhu, K. Luu, and M. Savvides. 2017. Faster than Real-Time Facial Alignment: A 3D Spatial Transformer Network Approach in Unconstrained Poses. In *Proceedings of the IEEE International Conference on Computer Vision*. 4000–4009.
- [15] V. Blanz and T. Vetter. 2003. Face recognition based on fitting a 3d morphable model. *IEEE Trans. Pattern Anal. Mach. Intell.* 25, 9 (2003), 1063–1074.
- [16] K. W. Bowyer, K. Chang, and P. Flynn. 2006. A survey of approaches and challenges in 3D and multi-modal 3D+ 2D face recognition. *Computer vision and image understanding* 101, 1 (2006), 1–15.
- [17] S. C. Brubaker, J. Wu, J. Sun, M. D. Mullin, and J. M. Rehg. 2008. On the Design of Cascades of Boosted Ensembles for Face Detection. *International Journal of Computer Vision* 77, 1-3 (2008), 65–86.
- [18] Adrian Bulat, Shiyang Cheng, Jing Yang, A. Garbett, Enrique Sanchez, and Georgios Tzimiropoulos. 2021. Pre-training strategies and datasets for facial representation learning. *ArXiv abs/2103.16554* (2021).

- [19] A. Bulat and G. Tzimiropoulos. 2016. Convolutional aggregation of local evidence for large pose face alignment. In *Proceedings of the British Machine Vision Conference*. 86.1–86.12.
- [20] A. Bulat and G. Tzimiropoulos. 2017. How Far are We from Solving the 2D & 3D Face Alignment Problem? (and a Dataset of 230,000 3D Facial Landmarks). In *Proceedings of the IEEE International Conference on Computer Vision*. 1021–1030.
- [21] X. P. Burgos-Artizzu, P. Perona, and P. Dollár. 2013. Robust face landmark estimation under occlusion. In *Proceedings of the IEEE International Conference on Computer Vision*. 1513–1520.
- [22] D. Cao, X. Zhu, X. Huang, J. Guo, and Z. Lei. 2020. Domain Balancing: Face Recognition on Long-Tailed Domains. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5671–5679.
- [23] K. Cao, Y. Rong, C. Li, X. Tang, and C. C. Loy. 2018. Pose-Robust Face Recognition via Deep Residual Equivariant Mapping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5187–5196.
- [24] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman. 2018. VGGFace2: A Dataset for Recognising Faces across Pose and Age. In *Proceedings of the IEEE International Conference on Automatic Face & Gesture Recognition*. 67–74.
- [25] F. Chang, A. T. Tran, T. Hassner, I. Masi, R. Nevatia, and G. Medioni. 2017. FacePoseNet: Making a Case for Landmark-Free Face Alignment. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 1599–1608.
- [26] J. Chang, Z. Lan, C. Cheng, and Y. Wei. 2020. Data Uncertainty Learning in Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5710–5719.
- [27] D. Chen, G. Hua, F. Wen, and J. Sun. 2016. Supervised transformer network for efficient face detection. In *Proceedings of the European Conference on Computer Vision*, Vol. 9909. 122–138.
- [28] J. Chen, R. Ranjan, S. Sankaranarayanan, A. Kumar, C. Chen, V. Patel, C. D. Castillo, and R. Chellappa. 2017. Unconstrained Still/Video-Based Face Verification with Deep Convolutional Neural Networks. *International Journal of Computer Vision* 126 (2017), 272–291.
- [29] L. Chen, H. Su, and Q. Ji. 2019. Face Alignment With Kernel Density Deep Neural Network. In *Proceedings of the IEEE International Conference on Computer Vision*. 6991–7001.
- [30] S. Chen, Y. Liu, X. Gao, and Z. Han. 2018. Mobilefacenet: Efficient cnns for accurate real-time face verification on mobile devices. In *Chinese Conference on Biometric Recognition*. 428–438.
- [31] Y. Cheng, J. Zhao, Z. Wang, Y. Xu, K. Jayashree, S. Shen, and J. Feng. 2017. Know you at one glance: A compact vector representation for low-shot learning. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 1924–1932.
- [32] C. Chi, S. Zhang, J. Xing, Z. Lei, S. Z. Li, and X. Zou. 2019. Selective refinement network for high performance face detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 8231–8238.
- [33] J. Choe, S. Park, K. Kim, J. Park, D. Kim, and H. Shim. 2017. Face Generation for Low-Shot Learning Using Generative Adversarial Networks. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 1940–1948.
- [34] F. Chollet. 2017. Xception: Deep Learning with Depthwise Separable Convolutions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1800–1807.
- [35] A. R. Chowdhury, T. Lin, S. Maji, and E. G. Learned-Miller. 2016. One-to-many face recognition with bilinear CNNs. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*. 1–9.
- [36] T. F. Cootes and C. J. Taylor. 1992. Active shape models—‘smart snakes’. In *Proceedings of the British Machine Vision Conference*. 266–275.
- [37] T. F. Cootes, K. Walker, and C. J. Taylor. 2000. View-based active appearance models. In *Proceedings of the IEEE International Conference on Automatic Face & Gesture Recognition*. 227–232.
- [38] A. Dapogny, M. Cord, and K. Bailly. 2019. DeCaFA: Deep Convolutional Cascade for Face Alignment in the Wild. In *Proceedings of the IEEE International Conference on Computer Vision*. 6892–6900.
- [39] Deepglint. 2020. Trillion Pairs. <http://trillionpairs.deepglint.com/overview>. (Accessed September 15, 2020).
- [40] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 248–255.
- [41] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. 2020. RetinaFace: single-shot multi-Level face localisation in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5203–5212.
- [42] J. Deng, J. Guo, N. Xue, and S. Zafeiriou. 2019. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4690–4699.
- [43] J. Deng, S. Cheng, N. Xue, Y. Zhou, and S. Zafeiriou. 2018. UV-GAN: Adversarial Facial UV Map Completion for Pose-Invariant Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7093–7102.
- [44] J. Deng, G. Trigeorgis, Y. Zhou, and S. Zafeiriou. 2019. Joint Multi-View Face Alignment in the Wild. *Trans. Image Process.* 28, 7 (2019), 3636–3648.

- [45] Jiankang Deng, Yuxiang Zhou, and Stefanos Zafeiriou. 2017. Marginal Loss for Deep Face Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 2006–2014.
- [46] W. Deng, J. Hu, N. Zhang, B. Chen, and J. Guo. 2017. Fine-grained face verification: FGLFW database, baselines, and human-DCMN partnership. *Pattern Recognition* 66 (2017), 63–73.
- [47] Z. Deng, X. Peng, Z. Li, and Y. Qiao. 2019. Mutual Component Convolutional Neural Networks for Heterogeneous Face Recognition. *Trans. Image Process.* 28 (2019), 3102–3114.
- [48] C. Ding and D. Tao. 2015. Robust Face Recognition via Multimodal Deep Face Representation. *IEEE Trans. Multimedia* 17 (2015), 2049–2058.
- [49] C. Ding and D. Tao. 2016. A Comprehensive Survey on Pose-Invariant Face Recognition. *ACM Trans. Intell. Syst. Technol.* 7 (2016), 37:1–37:42.
- [50] C. Ding and D. Tao. 2018. Trunk-Branch Ensemble Convolutional Neural Networks for Video-Based Face Recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* 40 (2018), 1002–1014.
- [51] X. Dong, Y. Yan, W. Ouyang, and Y. Yang. 2018. Style aggregated network for facial landmark detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 379–388.
- [52] X. Dong and Y. Yang. 2019. Teacher Supervises Students How to Learn From Partially Labeled Images for Facial Landmark Detection. In *Proceedings of the IEEE International Conference on Computer Vision*. 783–792.
- [53] X. Dong, S. Yu, X. Weng, S. Wei, Y. Yang, and Y. Sheikh. 2018. Supervision-by-Registration: An Unsupervised Approach to Improve the Precision of Facial Landmark Detectors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 360–368.
- [54] L. Du, H. Hu, and Y. Wu. 2019. Age Factor Removal Network Based on Transfer Learning and Adversarial Learning for Cross-Age Face Recognition. *IEEE Transactions on Circuits and Systems for Video Technology* 30, 9 (2019), 2830 – 2842.
- [55] Y. Duan, J. Lu, and J. Zhou. 2019. UniformFace: Learning Deep Equidistributed Representation for Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3410–3419.
- [56] C. N. Duong, K. G. Quach, N. Le, N. Nguyen, and K. Luu. 2019. Mobiface: A lightweight deep learning face recognition on mobile devices. In *Proceedings of the IEEE International Conference on Biometrics Theory, Applications and Systems*. 1–6.
- [57] M. Everingham, L. V. Gool, C. K. I. Williams, J. Winn, and A. Zisserman. 2010. The Pascal Visual Object Classes (VOC) Challenge. *International Journal of Computer Vision* 88, 2 (2010), p.303–338.
- [58] M. Everingham and J. Winn. 2011. The pascal visual object classes challenge 2012 (voc2012) development kit. *Pattern Analysis, Statistical Modelling and Computational Learning, Tech. Rep* 8 (2011).
- [59] S. S. Farfade, M. J. Saberian, and L. J. Li. 2015. Multi-view face detection using deep convolutional neural networks. In *Proceedings of the 5th ACM on International Conference on Multimedia Retrieval*. 643–650.
- [60] Y. Feng, F. Wu, X. Shao, Y. Wang, and X. Zhou. 2018. Joint 3d face reconstruction and dense alignment with position map regression network. In *Proceedings of the European Conference on Computer Vision*. 534–551.
- [61] Z. Feng, J. Kittler, M. Awais, P. Huber, and X. Wu. 2018. Wing Loss for Robust Facial Landmark Localisation with Convolutional Neural Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2235–2245.
- [62] C. Galea and R. A. Farrugia. 2017. Forensic Face Photo-Sketch Recognition Using a Deep Learning-Based Architecture. *IEEE Singal processing letters* 24 (2017), 1586–1590.
- [63] S. Ge, J. Li, Q. Ye, and Z. Luo. 2017. Detecting masked faces in the wild with lle-cnns. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2682–2690.
- [64] G. Huang, Z. Liu, and K. Q. Weinberger. 2017. Densely Connected Convolutional Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2261–2269.
- [65] S. Gong, Y. Shi, and A. K. Jain. 2019. Low Quality Video Face Recognition: Multi-Mode Aggregation Recurrent Network (MARN). In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 1027–1035.
- [66] S. Gong, Y. Shi, N. D. Kalka, and A. K. Jain. 2019. Video Face Recognition: Component-wise Feature Aggregation Network (C-FAN). In *Proceedings of the International Conference on Biometrics*. 1–8.
- [67] R. Gross, I. A. Matthews, J. F. Cohn, T. Kanade, and S. Baker. 2008. Multi-PIE. In *Proceedings of the IEEE International Conference on Automatic Face & Gesture Recognition*. 1–8.
- [68] P. Grother, R. J. Micheals, and P. J. Phillips. 2003. Face recognition vendor test 2002 performance metrics. In *International Conference on Audio-and Video-based Biometric Person Authentication*. 937–945.
- [69] P. Grother and M. Ngan. 2014. Face recognition vendor test (FRVT):Performance of face identification algorithms. *NIST Interagency report* 8009, 5 (2014), 14.
- [70] J. Guo, Jiankang Deng, Niannan Xue, and S. Zafeiriou. 2018. Stacked Dense U-Nets with Dual Transformers for Robust Face Alignment. In *Proceedings of the British Machine Vision Conference*.

- [71] S. Guo, J. Xu, D. Chen, C. Zhang, X. Wang, and R. Zhao. 2020. Density-Aware Feature Embedding for Face Clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6698–6706.
- [72] Y. Guo and L. Zhang. 2017. One-shot face recognition by promoting underrepresented classes. (2017). arXiv:1707.05574
- [73] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao. 2016. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *Proceedings of the European Conference on Computer Vision*. 87–102.
- [74] C. Han, S. Shan, M. Kan, S. Wu, and X. Chen. 2018. Face recognition with contrastive convolution. In *Proceedings of the European Conference on Computer Vision*. 118–134.
- [75] Z. Hao, Y. Liu, H. Qin, J. Yan, X. Li, and X. Hu. 2017. Scale-aware face detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6186–6195.
- [76] M. Hayat, S. H. Khan, N. Werghi, and R. Goecke. 2017. Joint Registration and Representation Learning for Unconstrained Face Identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1551–1560.
- [77] K. He, X. Zhang, S. Ren, and J. Sun. 2016. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 770–778.
- [78] R. He, J. Cao, L. Song, Z. Sun, and T. Tan. 2020. Adversarial Cross-Spectral Face Completion for NIR-VIS Face Recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* 42 (2020), 1025–1037.
- [79] R. He, X. Wu, Z. Sun, and T. Tan. 2017. Learning Invariant Deep Representation for NIR-VIS Face Recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 9005–9012.
- [80] R. He, X. Wu, Z. Sun, and T. Tan. 2019. Wasserstein CNN: Learning Invariant Features for NIR-VIS Face Recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* 41 (2019), 1761–1773.
- [81] S. Honari, P. Molchanov, S. Tyree, P. Vincent, C. J. Pal, and J. Kautz. 2018. Improving Landmark Localization with Semi-Supervised Learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1546–1555.
- [82] S. Hong, W. Im, J. B. Ryu, and H. Yang. 2017. SSPP-DAN: Deep domain adaptation network for face recognition with single sample per person. In *Proceedings of the IEEE International Conference on Image Processing*. 825–829.
- [83] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. 2017. MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. (2017). arXiv:abs/1704.04861
- [84] J. Hu, L. Shen, and G. Sun. 2018. Squeeze-and-Excitation Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7132–7141.
- [85] P. Hu and D. Ramanan. 2017. Finding tiny faces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1522–1530.
- [86] W. Hu, Y. Huang, F. Zhang, and R. Li. 2019. Noise-Tolerant Paradigm for Training Face Recognition CNNs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11879–11888.
- [87] Y. Hu, X. Wu, B. Yu, R. He, and Z. Sun. 2018. Pose-Guided Photorealistic Face Rotation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8398–8406.
- [88] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. 2007. *Labeled Faces in the Wild: A Database for Studying Face Recognition in Unconstrained Environments*. Technical Report 07-49. University of Massachusetts, Amherst.
- [89] L. Huang, Y. Yang, Y. Deng, and Y. Yu. 2015. Densebox: Unifying landmark localization with end to end object detection. (2015). arXiv:1509.04874
- [90] R. Huang, S. Zhang, T. Li, and R. He. 2017. Beyond Face Rotation: Global and Local Perception GAN for Photorealistic and Identity Preserving Frontal View Synthesis. In *Proceedings of the IEEE International Conference on Computer Vision*. 2458–2467.
- [91] X. Huang, W. Deng, H. Shen, X. Zhang, and J. Ye. 2020. PropagationNet: Propagate Points to Curve to Learn Structure Information. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7265–7274.
- [92] Y. Huang, Y. Wang, Y. Tai, X. Liu, P. Shen, S. Li, J. Li, and F. Huang. 2020. Curricularface: adaptive curriculum learning loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5901–5910.
- [93] H. Wang, Z. Li, X. Ji, and Y. Wang. 2017. Face r-cnn. (2017). arXiv:1706.01061
- [94] F. N. Iandola, M. W. Moskewicz, K. Ashraf, S. Han, W. J. Dally, and K. Keutzer. 2017. SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <1MB model size. (2017). arXiv:1602.07360
- [95] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu. 2015. Spatial transformer networks. In *Advances in neural information processing systems*. 2017–2025.
- [96] Anil K. Jain, Debayan Deb, and Joshua J. Engelsma. 2021. Biometrics: Trust, but Verify. *ArXiv abs/2105.06625* (2021).
- [97] H. Jiang and E. Learned-Miller. 2017. Face detection with the faster r-cnn. In *Proceedings of the IEEE International Conference on Automatic Face & Gesture Recognition*. 650–657.
- [98] H. Jin, S. Zhang, X. U. Zhu, Y. Tang, Z. Lei, and S. Z. Li. 2019. Learning Lightweight Face Detector with Knowledge Distillation. In *Proceedings of the International Conference on Biometrics*. 1–7.

- [99] X. Jin and X. Tan. 2017. Face alignment in-the-wild: A survey. *Computer Vision and Image Understanding* 162 (2017), 1–22.
- [100] Y. Jing, Q. Liu, and K. Zhang. 2017. Stacked Hourglass Network for Robust Facial Landmark Localisation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 79–87.
- [101] A. Jourabloo and X. Liu. 2016. Large-Pose Face Alignment via CNN-Based Dense 3D Model Fitting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4188–4196.
- [102] A. Jourabloo, M. Ye, X. Liu, and L. Ren. 2017. Pose-Invariant Face Alignment with a Single CNN. In *Proceedings of the IEEE International Conference on Computer Vision*. 3219–3228.
- [103] N. D. Kalka, B. Maze, J. A. Duncan, K. OrConnor, S. Elliott, K. Hebert, J. Bryan, and A. K. Jain. 2018. IJB-S: IARPA Janus Surveillance Video Benchmark. In *Proceedings of the IEEE International Conference on Biometrics Theory, Applications and Systems (BTAS)*. 1–9.
- [104] M. Kan, S. Shan, and X. Chen. 2016. Multi-view deep network for cross-view classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4847–4855.
- [105] B. Kang, Y. Kim, B. Jun, and D. Kim. 2019. Attentional Feature-Pair Relation Networks for Accurate Face Recognition. In *Proceedings of the IEEE International Conference on Computer Vision*. 5471–5480.
- [106] B. Kang, Y. Kim, and D. Kim. 2018. Pairwise relational networks for face recognition. In *Proceedings of the European Conference on Computer Vision*. 628–645.
- [107] I. Kemelmacher-Shlizerman, S. M. Seitz, D. Miller, and E. Brossard. 2016. The megaface benchmark: 1 million faces for recognition at scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4873–4882.
- [108] I. Kemelmacher-Shlizerman, S. Suwajanakorn, and S. Seitz. 2014. Illumination-Aware Age Progression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3334–3341.
- [109] Y. Kim, W. Park, M.I. Roh, and J. Shin. 2020. GroupFace: Learning Latent Groups and Constructing Group-based Representations for Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5621–5630.
- [110] B. Klare, B. Klein, E. Taborsky, A. Blanton, J. Cheney, K. E. Allen, P. Grother, A. Mah, M. Burge, and A. K. Jain. 2015. Pushing the frontiers of unconstrained face detection and recognition: IARPA Janus Benchmark A. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1931–1939.
- [111] M. Kostinger, P. Wohlhart, P. M. Roth, and H. Bischof. 2011. Annotated Facial Landmarks in the Wild: A large-scale, real-world database for facial landmark localization. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 2144–2151.
- [112] A. Krizhevsky, I. Sutskever, and G. E. Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. In *Advances in Neural Information Processing Systems*. 1097–1105.
- [113] A. Kumar and R. Chellappa. 2018. Disentangling 3D Pose in a Dendritic CNN for Unconstrained 2D Face Alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 430–439.
- [114] A.n Kumar, T.K. Marks, W. Mou, Y. Wang, M. Jones, A. Cherian, T. Koike-Akino, X. Liu, and C. Feng. 2020. LUVLi face alignment: estimating landmarks’ location, uncertainty, and visibility likelihood. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8236–8246.
- [115] V. Kushwaha, M. Singh, R. Singh, M. Vatsa, N. K. Ratha, and R. Chellappa. 2018. Disguised Faces in the Wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 1–18.
- [116] H. Law and J. Deng. 2018. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European Conference on Computer Vision*. 734–750.
- [117] V. Le, J. Brandt, Z. L. Lin, L. D. Bourdev, and T. S. Huang. 2012. Interactive Facial Feature Localization. In *Proceedings of the European Conference on Computer Vision*. 679–692.
- [118] J. Lezama, Q. Qiu, and G. Sapiro. 2017. Not Afraid of the Dark: NIR-VIS Face Recognition via Cross-Spectral Hallucination and Low-Rank Embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6807–6816.
- [119] H. Li, Z. Lin, X. Shen, J. Brandt, and G. Hua. 2015. A convolutional neural network cascade for face detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5325–5334.
- [120] J. Li, Y. Wang, C. Wang, Y. Tai, J. Qian, J. Yang, C. Wang, J. Li, and F. Huang. 2019. DSFD: dual shot face detector. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5060–5069.
- [121] S. Li, H. Li, J. Cui, and H. Zha. 2019. Pose-Aware Face Alignment based on CNN and 3DMM. In *Proceedings of the British Machine Vision Conference*. 106.
- [122] S. Z. Li, D. Yi, Z. Lei, and S. Liao. 2013. The CASIA NIR-VIS 2.0 Face Database. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 348–353.
- [123] S. Z. Li, L. Zhu, Z. Zhang, A. Blake, H. Zhang, and H. Shum. 2002. Statistical Learning of Multi-view Face Detection. In *Proceedings of the European Conference on Computer Vision*. 67–81.

- [124] Y. Li, B. Sun, T. Wu, and Y. Wang. 2016. Face detection with end-to-end integration of a convNet and a 3d model. In *Proceedings of the European Conference on Computer Vision*. 420–436.
- [125] Z. Li, X. Tang, J. Han, J. Liu, and R. He. 2019. PyramidBox++: High Performance Detector for Finding Tiny Face. (2019). arXiv:1904.00386
- [126] S. Liao, Z. Lei, D. Yi, and S. Z. Li. 2014. A benchmark study of large-scale unconstrained face recognition. In *Proceedings of the IEEE International Joint Conference on Biometrics*. 1–8.
- [127] T. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2117–2125.
- [128] T. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. 2017. Focal loss for dense object detection. In *Proceedings of the IEEE International Conference on Computer Vision*. 2980–2988.
- [129] W. Lin, J. Chen, C. D. Castillo, and R. Chellappa. 2018. Deep Density Clustering of Unconstrained Faces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8128–8137.
- [130] B. Liu, W. Deng, Y. Zhong, M. Wang, J. Hu, X. Tao, and Y. Huang. 2019. Fair Loss: Margin-Aware Reinforcement Learning for Deep Face Recognition. In *Proceedings of the IEEE International Conference on Computer Vision*. 10051–10060.
- [131] C. Liu and H. Wechsler. 2002. Gabor feature based classification using the enhanced fisher linear discriminant model for face recognition. *Trans. Image Process.* 11 4 (2002), 467–76.
- [132] H. Liu, J. Lu, J. Feng, and J. Zhou. 2018. Two-Stream Transformer Networks for Video-Based Face Alignment. *IEEE Trans. Pattern Anal. Mach. Intell.* 40 (2018), 2546–2554.
- [133] H. Liu, J. Lu, M. Guo, S. Wu, and J. Zhou. 2020. Learning Reasoning-Decision Networks for Robust Face Alignment. *IEEE Trans. Pattern Anal. Mach. Intell.* 42 (2020), 679–693.
- [134] H. Liu, X. Zhu, Z. Lei, and S. Z. Li. 2019. AdaptiveFace: Adaptive Margin and Sampling for Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11947–11956.
- [135] J. Liu, Y. Deng, T. Bai, Z. Wei, and C. Huang. 2015. Targeting Ultimate Accuracy: Face Recognition via Deep Embedding. (2015). arXiv:1506.07310
- [136] W. Liu, D. Anguelov, D. Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C. Berg. 2016. SSD: single shot multiBox detector. In *Proceedings of the European Conference on Computer Vision*. 21–37.
- [137] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song. 2017. Sphreface: Deep hypersphere embedding for face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 212–220.
- [138] W. Liu, Y. Wen, Z. Yu, and M. Yang. 2016. Large-margin softmax loss for convolutional neural networks.. In *ICML*, Vol. 2. 7.
- [139] X. Liu, B. V. K. V. Kumar, C. Yang, Q. Tang, and J. You. 2018. Dependency-Aware Attention Control for Unconstrained Face Recognition with Image Sets. In *Proceedings of the European Conference on Computer Vision*. 548–565.
- [140] X. Liu, L. Song, X. Wu, and T. Tan. 2016. Transferring deep representation for NIR-VIS heterogeneous face recognition. *Proceedings of the International Conference on Biometrics* (2016), 1–8.
- [141] Y. Liu, A. Jourabloo, W. Ren, and X. Liu. 2017. Dense Face Alignment. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 1619–1628.
- [142] Y. Liu, H. Li, and X. Wang. 2017. Rethinking Feature Discrimination and Polymerization for Large-scale Recognition. (2017). arXiv:1710.00870
- [143] Y. Liu, H. Li, J. Yan, F. Wei, X. Wang, and X. Tang. 2017. Recurrent Scale Approximation for Object Detection in CNN. In *Proceedings of the IEEE International Conference on Computer Vision*. 571–579.
- [144] Y. Liu, H. Shen, Y. Si, X. Wang, X. Zhu, H. Shi, Z. Hong, H. Guo, Z. Guo, Y. Chen, B. Li, T. Xi, J. Yu, H. Xie, G. Xie, M. Li, Q. Lu, Z. Wang, S. Lai, Z. Chai, and X. Wei. 2019. Grand Challenge of 106-Point Facial Landmark Localization. In *Proceedings of the IEEE ICME Workshop*. 613–616.
- [145] Y. Liu, X. Tang, J. Han, J. Liu, D. Rui, and X. Wu. 2020. HAMBox: Delving Into Mining High-Quality Anchors on Face Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13043–13051.
- [146] Z. Liu, H. H. CloudMinds, J. Bai, S. Li, and S. L. CloudMinds. 2019. Feature Aggregation Network for Video Face Recognition. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 990–998.
- [147] Z. Liu, X. Zhu, G. Hu, H. Guo, M. Tang, Z. Lei, N. M. Robertson, and J. Wang. 2019. Semantic Alignment: Finding Semantically Consistent Ground-Truth for Facial Landmark Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3462–3471.
- [148] Ping Luo, Zhenyao Zhu, Ziwei Liu, Xiaogang Wang, and Xiaoou Tang. 2016. Face Model Compression by Distilling Knowledge from Neurons. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 3560–3566.
- [149] J. Lv, X. Shao, J. Xing, C. Cheng, and X. Zhou. 2017. A Deep Regression Architecture with Two-Stage Re-initialization for High Performance Facial Landmark Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3691–3700.

- [150] m. t. Pham and T.J. Cham. 2007. Fast training and selection of Haar features using statistics in boosting-based face detection. In *Proceedings of the IEEE International Conference on Computer Vision*. 1–7.
- [151] N. Ma, X. Zhang, H. Zheng, and J. Sun. 2018. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of the European Conference on Computer Vision*. 116–131.
- [152] R. Manmatha, C.n Wu, A. Smola, and P. Krähenbühl. 2017. Sampling Matters in Deep Embedding Learning. In *Proceedings of the IEEE International Conference on Computer Vision*. 2859–2867.
- [153] B. Martínez, M. F. Valstar, X. Binefa, and M. Pantic. 2013. Local Evidence Aggregation for Regression-Based Facial Point Detection. *IEEE Trans. Pattern Anal. Mach. Intell.* 35 (2013), 1149–1163.
- [154] Y. Martínez-Díaz, L. S. Luevano, H. M. Vazquez, M. Nicolás-Díaz, L. Chang, and M. González-Mendoza. 2019. Shuffle-FaceNet: A Lightweight Face Architecture for Efficient and Highly-Accurate Face Recognition. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 2721–2728.
- [155] I. Masi, S. Rawls, G. G. Medioni, and P. Natarajan. 2016. Pose-Aware Face Recognition in the Wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4838–4846.
- [156] B. Maze, J. C. Adams, J.s A. Duncan, N. D. Kalka, T. Miller, C. Otto, A. K. Jain, W. T. Niggel, J. Anderson, J. Cheney, and Patrick Grother. 2018. IARPA Janus Benchmark - C: Face Dataset and Protocol. In *Proceedings of the International Conference on Biometrics*. 158–165.
- [157] Tao Mei, Bo Yang, Shiqiang Yang, and Xiansheng Hua. 2008. Video collage: presenting a video sequence using a single image. *The Visual Computer* 25 (2008), 39–51.
- [158] D. Merget, M. Rock, and G. Rigoll. 2018. Robust Facial Landmark Detection via a Fully-Convolutional Local-Global Context Network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 781–790.
- [159] X. Miao, X. Zhen, X. Liu, C. Deng, V. Athitsos, and H. Huang. 2018. Direct Shape Regression Networks for End-to-End Face Alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5040–5049.
- [160] Shervin Minaee, Ping Luo, Zhe Lin, and K. Bowyer. 2021. Going Deeper Into Face Detection: A Survey. *ArXiv abs/2103.14983* (2021).
- [161] X. Ming, F. Wei, T. Zhang, D. Chen, and F. Wen. 2019. Group Sampling for Scale Invariant Face Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3446–3456.
- [162] T. Mita, T. Kaneko, and O. Hori. 2005. Joint Haar-like features for face detection. In *Proceedings of the IEEE International Conference on Computer Vision*. 1619–1626.
- [163] P. Mittal, M. Vatsa, and R. Singh. 2015. Composite sketch recognition via deep network - a transfer learning approach. In *Proceedings of International Conference on Biometrics*. 251–256.
- [164] S. Moschoglou, A. Papaioannou, C. Sagonas, J. Deng, I. Kotsia, and S. Zafeiriou. 2017. AgeDB: The First Manually Collected, In-the-Wild Age Database. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 1997–2005.
- [165] M. Najibi, P. Samangouei, R. Chellappa, and L. S. Davis. 2017. SSH: single stage headless face detector. In *Proceedings of the IEEE International Conference on Computer Vision*. 4885–4894.
- [166] M. Najibi, B. Singh, and L. S. Davis. 2019. FA-RPN: Floating Region Proposals for Face Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7723–7732.
- [167] A. Nech and I. Kemelmacher-Shlizerman. 2017. Level Playing Field for Million Scale Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3406–3415.
- [168] A. Newell, K. Yang, and J. Deng. 2016. Stacked Hourglass Networks for Human Pose Estimation. In *Proceedings of the European Conference on Computer Vision*. 483–499.
- [169] T. Ojala, M. Pietikainen, and T. Maenpaa. 2002. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE Trans. Pattern Anal. Mach. Intell.* 24, 7 (2002), 971–987.
- [170] M. Opitz, G. Waltner, G. Poier, H. Possegger, and H. Bischof. 2016. Grid loss: detecting occluded faces. In *Proceedings of the European Conference on Computer Vision*, Vol. 9907. 386–402.
- [171] M. Parchami, S. Bashbaghi, E. Granger, and S. Sayed. 2017. Using deep autoencoders to learn robust domain-invariant representations for still-to-video face recognition. In *Proceedings of the IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*. 1–6.
- [172] O. M. Parkhi, A. Vedaldi, and A. Zisserman. 2015. Deep Face Recognition. In *Proceedings of the British Machine Vision Conference*. 41.1–41.12.
- [173] X. Peng, R. S. Feris, X. Wang, and D. N. Metaxas. 2016. A Recurrent Encoder-Decoder Network for Sequential Face Alignment. In *Proceedings of the European Conference on Computer Vision*. 38–56.
- [174] H. Qin, J. Yan, L. Xiu, and X. Hu. 2016. Joint Training of Cascaded CNN for Face Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3456–3465.
- [175] R. Ranjan, C. D. Castillo, and R. Chellappa. 2017. L2-constrained softmax loss for discriminative face verification. (2017). [arXiv:1703.09507](https://arxiv.org/abs/1703.09507)

- [176] R. Ranjan, V. M. Patel, and R. Chellappa. 2019. HyperFace: A Deep Multi-Task Learning Framework for Face Detection, Landmark Localization, Pose Estimation, and Gender Recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* 41, 1 (2019), 121–135.
- [177] R. Ranjan, S. Sankaranarayanan, A. Bansal, N. Bodla, J. Chen, V. M. Patel, C. D. Castillo, and R. Chellappa. 2018. Deep Learning for Understanding Faces: Machines May Be Just as Good, or Better, than Humans. *IEEE Signal Processing Magazine* 35, 1 (2018), 66–83.
- [178] Y. Rao, J. Lin, J. Lu, and J. Zhou. 2017. Learning Discriminative Aggregation Network for Video-Based Face Recognition. In *Proceedings of the IEEE International Conference on Computer Vision*. 3801–3810.
- [179] Y. Rao, J. Lu, and J. Zhou. 2017. Attention-Aware Deep Reinforcement Learning for Video Face Recognition. In *Proceedings of the IEEE International Conference on Computer Vision*. 3951–3960.
- [180] C. Reale, N. M. Nasrabadi, H. Kwon, and R. Chellappa. 2016. Seeing the Forest from the Trees: A Holistic Approach to Near-Infrared Heterogeneous Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 320–328.
- [181] S. Ren, K. He, R. Girshick, and J. Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*. 91–99.
- [182] J. P. Robinson, Y. Li, N. Zhang, Y. Fu, and S. Tulyakov. 2019. Laplace Landmark Localization. In *Proceedings of the IEEE International Conference on Computer Vision*. 10102–10111.
- [183] A. RoyChowdhury, X. Yu, K. Sohn, E. Learned-Miller, and M. Chandraker. 2020. Improving Face Recognition by Clustering Unlabeled Faces in the Wild. In *Proceedings of the European Conference on Computer Vision*. 119–136.
- [184] C. Sagonas, E. Antonakos, G. Tzimiropoulos, S. Zafeiriou, and M. Pantic. 2016. 300 Faces In-The-Wild Challenge: database and results. *Image Vis. Comput.* 47 (2016), 3–18.
- [185] C. Sagonas, G. Tzimiropoulos, S. Zafeiriou, and M. Pantic. 2013. 300 Faces in-the-Wild Challenge: The First Facial Landmark Localization Challenge. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 397–403.
- [186] M. Sandler, A. G. Howard, M. Zhu, A. Zhmoginov, and L. Chen. 2018. MobileNetV2: Inverted Residuals and Linear Bottlenecks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4510–4520.
- [187] S. Sankaranarayanan, A. Alavi, C. D. Castillo, and R. Chellappa. 2016. Triplet probabilistic embedding for face verification and clustering. In *Proceedings of the IEEE International Conference on Biometrics Theory, Applications and Systems*. 1–8.
- [188] S. Saxena and J. Verbeek. 2016. Heterogeneous Face Recognition with CNNs. In *Proceedings of the European Conference on Computer Vision Workshops*. 483–491.
- [189] F. Schroff, D. Kalenichenko, and J. Philbin. 2015. FaceNet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 815–823.
- [190] M. Schultz and T. Joachims. 2004. Learning a distance metric from relative comparisons. In *Advances in neural information processing systems*. 41–48.
- [191] S. Sengupta, J. Chen, C. Castillo, V. M. Patel, R. Chellappa, and D. W. Jacobs. 2016. Frontal to profile face verification in the wild. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*. 1–9.
- [192] J. Shen, S. Zafeiriou, G. G. Chrysos, J. Kossaifi, G. Tzimiropoulos, and M. Pantic. 2015. The First Facial Landmark Tracking in-the-Wild Challenge: Benchmark and Results. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 1003–1011.
- [193] X. Shi, S. Shan, M. Kan, S. Wu, and X. Chen. 2018. Real-Time rotation-invariant face detection with progressive calibration networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2295–2303.
- [194] Y. Shi, A. K. Jain, and N. D. Kalka. 2019. Probabilistic Face Embeddings. In *Proceedings of the IEEE International Conference on Computer Vision*. 6901–6910.
- [195] Y. Shi, C. Otto, and A. K. Jain. 2018. Face clustering: representation and pairwise constraints. *IEEE Transactions on Information Forensics and Security* 13, 7 (2018), 1626–1640.
- [196] Y. Shi, X. Yu, K. Sohn, M. Chandraker, and A. K. Jain. 2020. Towards Universal Representation Learning for Deep Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6817–6826.
- [197] A. Shrivastava, A. Gupta, and R. Girshick. 2016. Training region-based object detectors with online hard example mining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 761–769.
- [198] K. Simonyan and A. Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. (2015). arXiv:1409.1556
- [199] E. Smirnov, A. Melnikov, S. Novoselov, E. Lubyantsev, and G. Lavrentyeva. 2017. Doppelganger mining for face representation learning. In *Proceedings of the IEEE International Conference on Computer Vision*. 1916–1923.
- [200] K. Sohn. 2016. Improved deep metric learning with multi-class n-pair loss objective. In *Advances in Neural Information Processing Systems*. 1857–1865.



- [201] S. Soltanpour, B. Boufama, and Q. M. J. Wu. 2017. A survey of local feature methods for 3D face recognition. *Pattern Recognition* 72 (2017), 391–406.
- [202] G. Song, Y. Liu, M. Jiang, Y. Wang, J. Yan, and B. Leng. 2018. Beyond Trade-Off: Accelerate FCN-Based Face Detector With Higher Accuracy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7756–7764.
- [203] H. Oh Song, Y. Xiang, S. Jegelka, and S. Savarese. 2016. Deep metric learning via lifted structured feature embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4004–4012.
- [204] J. Song, J. Zhang, L. Gao, X. Liu, and H. Shen. 2018. Dual Conditional GANs for Face Aging and Rejuvenation.. In *International Joint Conference on Artificial Intelligence*. 899–905.
- [205] Lingxiao Song, Man Zhang, Xiang Wu, and Ran He. 2018. Adversarial discriminative heterogeneous face recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.
- [206] Sijie Song, Wei Zhang, Jiaying Liu, and Tao Mei. 2019. Unsupervised Person Image Generation With Semantic Parsing Transformation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2352–2361.
- [207] K. Sun, W. Wu, T. Liu, S. Yang, Q. Wang, Q. Zhou, Z. Ye, and C. Qian. 2019. FAB: A Robust Facial Landmark Detection Framework for Motion-Blurred Videos. In *Proceedings of the IEEE International Conference on Computer Vision*. 5461–5470.
- [208] X. Sun, P. Wu, and S. C. H. Hoi. 2018. Face detection using deep learning: an improved faster rcnn approach. *Neurocomputing* 299 (2018), 42 – 50.
- [209] Y. Sun, Y. Chen, X. Wang, and X. Tang. 2014. Deep learning face representation by joint identification-verification. In *Advances in neural information processing systems*. 1988–1996.
- [210] Y. Sun, C. Cheng, Y. Zhang, C. Zhang, L. Zheng, Z. Wang, and Y. Wei. 2020. Circle Loss: A Unified Perspective of Pair Similarity Optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6398–6407.
- [211] Y. Sun, D. Liang, X. Wang, and X. Tang. 2015. DeepID3: Face Recognition with Very Deep Neural Networks. (2015). arXiv:1502.00873
- [212] Y. Sun, X. Wang, and X. Tang. 2013. Deep Convolutional Network Cascade for Facial Point Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3476–3483.
- [213] Y. Sun, X. Wang, and X. Tang. 2013. Hybrid Deep Learning for Face Verification. In *Proceedings of the IEEE International Conference on Computer Vision*. 1489–1496.
- [214] Y. Sun, X. Wang, and X. Tang. 2014. Deep Learning Face Representation from Predicting 10,000 Classes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1891–1898.
- [215] Y. Sun, X. Wang, and X. Tang. 2015. Deeply learned face representations are sparse, selective, and robust. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 2892–2900.
- [216] Y. Sun, X. Wang, and X. Tang. 2016. Sparsifying neural network connections for face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4856–4864.
- [217] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. 2015. Going deeper with convolutions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1–9.
- [218] O. Tadmor, T. Rosenwein, S. Shalev-Shwartz, Y. Wexler, and A. Shashua. 2016. Learning a Metric Embedding for Face Recognition using the Multibatch Method. In *Advances in neural information processing systems*. 1396–1397.
- [219] Y. Tai, Y. Liang, X. Liu, L. Duan, J. Li, C. Wang, F. Huang, and Y. Chen. 2019. Towards Highly Accurate and Stable Face Alignment for High-Resolution Videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 8893–8900.
- [220] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. 2014. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1701–1708.
- [221] X. Tang, D. K. Du, Z. He, and J. Liu. 2018. Pyramidbox: a context-assisted single shot face detector. In *Proceedings of the European Conference on Computer Vision*. 797–813.
- [222] Zhiqiang Tang, Xi Peng, Shijie Geng, Lingfei Wu, Shaoting Zhang, and Dimitris N. Metaxas. 2018. Quantized Densely Connected U-Nets for Efficient Landmark Localization. In *Proceedings of the European Conference on Computer Vision*. 348–364.
- [223] W. Tian, Z. Wang, H. Shen, W. Deng, Y. Meng, B. Chen, X. Zhang, Y. Zhao, and X. Huang. 2018. Learning Better Features for Face Detection with Feature Fusion and Segmentation Supervision. arXiv:1811.08557
- [224] Z. Tian, C. Shen, H. Chen, and T. He. 2019. FCOS: Fully Convolutional One-Stage Object Detection. In *Proceedings of the IEEE International Conference on Computer Vision*. 9626–9635.
- [225] L. Tran, X. Yin, and X. Liu. 2017. Disentangled Representation Learning GAN for Pose-Invariant Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1283–1292.
- [226] G. Trigeorgis, P. Snape, M. A. Nicolaou, E. Antonakos, and S. Zafeiriou. 2016. Mnemonic Descent Method: A Recurrent Process Applied for End-to-End Face Alignment. In *IEEE Conference on Computer Vision and Pattern Recognition*.

4177–4187.

- [227] M. Turk and A. Pentland. 1991. Eigenfaces for Recognition. *Journal of Cognitive Neuroscience* 3, 1 (1991), 71–86.
- [228] T. Zheng, W. Deng, and J. Hu. 2017. Age Estimation Guided Convolutional Neural Network for Age-Invariant Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 503–511.
- [229] Jain V and E. Learned-Miller. 2010. *Fddb: A Benchmark for Face Detection in Unconstrained Settings*. Technical Report UM-CS-2010-009.
- [230] Paul Viola and Michael Jones. 2001. Rapid object detection using a boosted cascade of simple features. In *Proceedings of the conference on computer vision and pattern recognition*, Vol. 1. I–I.
- [231] Zhao w, R. Chellappa, P. J. Phillips, and A. Rosenfeld. 2003. Face recognition: A literature survey. *ACM computing surveys* 35, 4 (2003), 399–458.
- [232] S. Wan, Z. Chen, T. Zhang, B. Zhang, and K. Wong. 2016. Bootstrapping face detection with hard negative examples. (2016). arXiv:1608.02236
- [233] W. Wan, Y. Gao, and H. J. Lee. 2019. Transfer deep feature learning for face sketch recognition. *Neural Computing and Applications* (2019), 1–10.
- [234] F. Wang, L. Chen, C. Li, S. Huang, Y. Chen, C. Qian, and C. C. Loy. 2018. The Devil of Face Recognition is in the Noise. In *Proceedings of the European Conference on Computer Vision*. 765–780.
- [235] F. Wang, J. Cheng, W. Liu, and H. Liu. 2018. Additive margin softmax for face verification. *IEEE Signal processing letters* 25, 7 (2018), 926–930.
- [236] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang. 2017. Residual Attention Network for Image Classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6450–6458.
- [237] F. Wang, X. Xiang, J. Cheng, and A. L. Yuille. 2017. Normface: l<sub>2</sub> hypersphere embedding for face verification. In *Proceedings of the 25th ACM international conference on Multimedia*. 1041–1049.
- [238] H. Wang, D. Gong, Z. Li, and W. Liu. 2019. Decorrelated Adversarial Learning for Age-Invariant Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3522–3531.
- [239] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu. 2018. Cosface: Large margin cosine loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5265–5274.
- [240] J. Wang, k. Sun, T. Cheng, B. Jiang, C. Deng, Y. Zhao, D. Liu, Y. Mu, M. Tan, X. Wang, W. Liu, and B. Xiao. 2021. Deep high-resolution representation learning for visual recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* 43 (2021), 3349–3364.
- [241] J. Wang, Y. Yuan, B. Li, G. Yu, and S. Jian. 2018. Sface: An efficient network for face detection in large scale variations. (2018). arXiv:1804.06559
- [242] J. Wang, Y. Yuan, and G. Yu. 2017. Face attention network: an effective face detector for the occluded faces. (2017). arXiv:1711.07246
- [243] L. Wang, V. Sindagi, and V. M. Patel. 2018. High-Quality Facial Photo-Sketch Synthesis Using Multi-Adversarial Networks. In *Proceedings of the IEEE International Conference on Automatic Face & Gesture Recognition*. 83–90.
- [244] M. Wang and W. Deng. 2018. Deep Face Recognition: A Survey. *Neurocomputing* 312 (2018), 135–153.
- [245] M. Wang and W. Deng. 2020. Mitigating Bias in Face Recognition Using Skewness-Aware Reinforcement Learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9322–9331.
- [246] M. Wang, W. Deng, J. Hu, X. Tao, and Y. Huang. 2019. Racial Faces in the Wild: Reducing Racial Bias by Information Maximization Adaptation Network. In *Proceedings of the IEEE International Conference on Computer Vision*. 692–702.
- [247] N. Wang, X. Gao, D. Tao, H. Yang, and X. Li. 2018. Facial feature point detection: A comprehensive survey. *Neurocomputing* 275 (2018), 50–65.
- [248] W. Wang, Z. Cui, Y. Yan, J. Feng, S. Yan, X. Shu, and N. Sebe. 2016. Recurrent Face Aging. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2378–2386.
- [249] X. Wang, L. Bo, and L. Fuxin. 2019. Adaptive Wing Loss for Robust Face Alignment via Heatmap Regression. In *Proceedings of the IEEE International Conference on Computer Vision*. 6970–6980.
- [250] X. Wang, S. Wang, J. Wang, H. Shi, and T. Mei. 2019. Co-Mining: Deep Face Recognition With Noisy Labels. In *Proceedings of the IEEE International Conference on Computer Vision*. 9358–9367.
- [251] Xiaobo Wang, Shifeng Zhang, Shuo Wang, Tianyu Fu, Hailin Shi, and Tao Mei. 2020. Mis-classified vector guided softmax loss for face recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 12241–12248.
- [252] Y. Wang, D. Gong, Z. Zhou, X. Ji, H. Wang, Z. Li, W. Liu, and T. Zhang. 2018. Orthogonal Deep Features Decomposition for Age-Invariant Face Recognition. In *Proceedings of the European Conference on Computer Vision*. 738–753.
- [253] Y. Wang, X. Ji, Z. Zhou, H. Wang, and Z. Li. 2017. Detecting faces using region-based fully convolutional networks. (2017). arXiv:1709.05256
- [254] Z. Wang, X. Tang, W. Luo, and S. Gao. 2018. Face Aging with Identity-Preserved Conditional Generative Adversarial Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7939–7947.

- [255] Z. Wang, L. Zheng, Y. Li, and S. Wang. 2019. Linkage based face clustering via graph convolution network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1117–1125.
- [256] H. Wei, P. Lu, and Y. Wei. 2020. Balanced Alignment for Face Recognition: A Joint Learning Approach. (2020). arXiv:2003.10168
- [257] K. Q. Weinberger and L. K. Saul. 2006. Distance Metric Learning for Large Margin Nearest Neighbor Classification. In *Advances in neural information processing systems*. 1473–1480.
- [258] Y. Wen, Z. Li, and Y. Qiao. 2016. Latent Factor Guided Convolutional Neural Networks for Age-Invariant Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4893–4901.
- [259] Y. Wen, K. Zhang, Z. Li, and Y. Qiao. 2016. A discriminative feature learning approach for deep face recognition. In *Proceedings of the European Conference on Computer Vision*. 499–515.
- [260] C. Whitelam, E. Taborsky, A. Blanton, B. Maze, J. C. Adams, T. Miller, N. D. Kalka, A. K. Jain, J. A. Duncan, K. E. Allen, J. Cheney, and P. Grother. 2017. IARPA Janus Benchmark-B Face Dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 592–600.
- [261] Jonathan R Williford, Brandon B May, and Jeffrey Byrne. 2020. Explainable Face Recognition. In *Proceedings of the European Conference on Computer Vision*. 248–263.
- [262] L. Wolf, T. Hassner, and I. Maoz. 2011. Face recognition in unconstrained videos with matched background similarity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 529–534.
- [263] W. Wu, M. Kan, X. Liu, Y. Yang, S. Shan, and X. Chen. 2017. Recursive Spatial Transformer (ReST) for Alignment-Free Face Recognition. In *Proceedings of the IEEE International Conference on Computer Vision*. 3792–3800.
- [264] W. Wu, C. Qian, S. Yang, Q. Wang, Y. Cai, and Q. Zhou. 2018. Look at Boundary: A Boundary-Aware Face Alignment Algorithm. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2129–2138.
- [265] X. Wu, R. He, and Z. Sun. 2015. A Lightened CNN for Deep Face Representation. (2015). arXiv:1511.02683
- [266] X. Wu, R. He, Z. Sun, and T. Tan. 2018. A Light CNN for Deep Face Representation With Noisy Labels. *IEEE Transactions on Information Forensics and Security* 13 (2018), 2884–2896.
- [267] X. Wu, H. Huang, V. M. Patel, R. He, and Z. Sun. 2019. Disentangled Variational Representation for Heterogeneous Face Recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 9005–9012.
- [268] Y. Wu, T. Hassner, K. Kim, G. Medioni, and P. Natarajan. 2018. Facial Landmark Detection with Tweaked Convolutional Neural Networks. *IEEE Trans. Pattern Anal. Mach. Intell.* 40, 12 (2018), 3067–3074.
- [269] Y. Wu, H. Liu, and Y. Fu. 2017. Low-Shot Face Recognition with Hybrid Classifiers. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 1933–1939.
- [270] Y. Wu, Y. Wu, R. Gong, Y. Lv, K. Chen, D. Liang, X. Hu, X. Liu, and J. Yan. 2020. Rotation Consistent Margin Loss for Efficient Low-Bit Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6866–6876.
- [271] X. Fan, R. Liu, K. Huan, Y. Feng, and Z. Luo. 2018. Self-Reinforced Cascaded Regression for Face Alignment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.
- [272] S. Xiao, J. Feng, L. Liu, X. Nie, W. Wang, S. Yan, and A. Kassim. 2017. Recurrent 3D-2D Dual Learning for Large-Pose Facial Landmark Detection. In *Proceedings of the IEEE International Conference on Computer Vision*. 1642–1651.
- [273] S. Xiao, J. Feng, J. Xing, H. Lai, S. Yan, and A. A. Kassim. 2016. Robust Facial Landmark Detection via Recurrent Attentive-Refinement Networks. In *Proceedings of the European Conference on Computer Vision*. 57–72.
- [274] W. Xie, L. Shen, and A. Zisserman. 2018. Comparator Networks. In *Proceedings of the European Conference on Computer Vision*. 782–797.
- [275] Yilin Xiong, Zijian Zhou, Yuhao Dou, and Zhizhong Su. 2020. Gaussian vector: An efficient solution for facial landmark detection. In *Proceedings of the Asian Conference on Computer Vision*. 70–87.
- [276] C. Xu, Q. Liu, and M. Ye. 2017. Age invariant face recognition and retrieval by coupled auto-encoder networks. *Neurocomputing* 222 (2017), 62–71.
- [277] X. Xu and I. A. Kakadiaris. 2017. Joint head pose estimation and face alignment framework using global and local cnn features. In *Proceedings of the IEEE International Conference on Automatic Face & Gesture Recognition*. 642–649.
- [278] Xiqing Xu, Qiang Meng, Yunxiao Qin, Jianzhu Guo, Chenxu Zhao, Feng Zhou, and Zhen Lei. 2021. Searching for Alignment in Face Recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 3065–3073.
- [279] Y. Xu, W. Yan, H. Sun, G. Yang, and J. Luo. 2019. CenterFace: Joint Face Detection and Alignment Using Face as Point. arXiv:1911.03599
- [280] J. Yan, X. Zhang, Z. Lei, and S. Z. Li. 2014. Face detection by structural models. *Image and Vision Computing* 32, 10 (2014), 790 – 799.
- [281] B. Yang, J. Yan, Z. Lei, and S. Z. Li. 2015. Fine-grained evaluation on face detection in the wild. In *Proceedings of the IEEE International Conference on Automatic Face & Gesture Recognition*, Vol. 1. 1–7.
- [282] Jing Yang, Adrian Bulat, and Georgios Tzimiropoulos. 2020. Fan-face: a simple orthogonal improvement to deep face recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 12621–12628.

- [283] J. Yang, P. Ren, D. Zhang, D. Chen, F. Wen, H. Li, and G. Hua. 2017. Neural Aggregation Network for Video Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5216–5225.
- [284] L. Yang, D. Chen, X. Zhan, R. Zhao, C. C. Loy, and D. Lin. 2020. Learning to Cluster Faces via Confidence and Connectivity Estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13369–13378.
- [285] L. Yang, X. Zhan, D. Chen, J. Yan, C. C. Loy, and D. Lin. 2019. Learning to cluster faces on an affinity graph. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2298–2306.
- [286] Ming-Hsuan Yang, D. Kriegman, and N. Ahuja. 2002. Detecting Faces in Images: A Survey. *IEEE Trans. Pattern Anal. Mach. Intell.* 24 (2002), 34–58.
- [287] S. Yang, P. Luo, C. C. Loy, and X. Tang. 2015. From facial parts responses to face detection: a deep learning approach. In *Proceedings of the IEEE International Conference on Computer Vision*. 3676–3684.
- [288] S. Yang, P. Luo, C. C. Loy, and X. Tang. 2016. WIDER FACE: A Face Detection Benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5525–5533.
- [289] D. Yi, Z. Lei, S. Liao, and S. Z. Li. 2014. Learning Face Representation from Scratch. (2014). arXiv:1411.7923
- [290] Z. Yi, H. Zhang, P. Tan, and M. Gong. 2017. DualGAN: Unsupervised Dual Learning for Image-to-Image Translation. In *Proceedings of the IEEE International Conference on Computer Vision*. 2868–2876.
- [291] Bangjie Yin, Luan Tran, Haoxiang Li, Xiaohui Shen, and Xiaoming Liu. 2019. Towards interpretable face recognition. In *Proceedings of the IEEE International Conference on Computer Vision*. 9348–9357.
- [292] X. Yin, X. Yu, K. Sohn, X. Liu, and M. Chandraker. 2019. Feature Transfer Learning for Face Recognition with Under-Represented Data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5704–5713.
- [293] H. Yu, Y. Fan, K. Chen, H. Yan, X. Lu, J. Liu, and D. Xie. 2019. Unknown Identity Rejection Loss: Utilizing Unlabeled Data for Face Recognition. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 2662–2669.
- [294] J. Yu, Y. Jiang, Z. Wang, Z. Cao, and T. Huang. 2016. Unitbox: An advanced object detection network. In *Proceedings of the 24th ACM international conference on Multimedia*. 516–520.
- [295] Y. Yu, G. Song, M. Zhang, J. Liu, Y. Zhou, and J. Yan. 2019. Towards flops-constrained face recognition. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 2698–2702.
- [296] L. Yue, X. Miao, P. Wang, B. Zhang, X. Zhen, and X. Cao. 2018. Attentional Alignment Networks. In *Proceedings of the British Machine Vision Conference*, Vol. 2. 6–13.
- [297] z. Liu, p. Luo, x. Wang, and X. Tang. 2015. Deep learning face attributes in the wild. In *Proceedings of the IEEE International Conference on Computer Vision*. 3730–3738.
- [298] S. Zafeiriou, G. Trigeorgis, G. Chrysos, J. Deng, and Ji. Shen. 2017. The Menpo Facial Landmark Localisation Challenge: A Step Towards the Solution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition W*. 2116–2125.
- [299] S. Zafeiriou, C. Zhang, and Z. Zhang. 2015. A survey on face detection in the wild: Past, present and future. *Computer Vision and Image Understanding* 138 (2015), 1–24.
- [300] Timothy Zee, Geeta Gali, and Ifeoma Nwogu. 2019. Enhancing human face recognition with an interpretable neural network. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 514–522.
- [301] D. Zeng, H. Liu, F. Zhao, S. Ge, W. Shen, and Z. Zhang. 2019. Proposal pyramid networks for fast face detection. *Information Sciences* 495 (2019), 136 – 149.
- [302] X. Zhan, Z. Liu, J. Yan, D. Lin, and C. C. Loy. 2018. Consensus-Driven Propagation in Massive Unlabeled Data for Face Recognition. In *Proceedings of the European Conference on Computer Vision*. 568–583.
- [303] B. Zhang, J. Li, Y. Wang, Y. Tai, C. Wang, J. Li, F. Huang, Y. Xia, W. Pei, and R. Ji. 2020. ASFD: Automatic and Scalable Face Detector. arXiv:2003.11228
- [304] C. Zhang, X. Xu, and D. Tu. 2018. Face detection using improved faster rcnn. (2018). arXiv:1802.02142
- [305] C. Zhang and Z. Zhang. 2014. Improving multiview face detection with multi-task deep convolutional neural networks. In *IEEE Winter Conference on Applications of Computer Vision*. 1036–1041.
- [306] D. Zhang, L. Lin, T. Chen, X. Wu, W. Tan, and E. Izquierdo. 2017. Content-Adaptive Sketch Portrait Generation by Decompositional Representation Learning. *Trans. Image Process.* 26 (2017), 328–339.
- [307] F. Zhang. 2006. Face recognition from a single image per person: A survey. *Pattern Recognition* 39, 9 (2006), 1725–1745.
- [308] F. Zhang, X. Fan, G. Ai, J. Song, Y. Qin, and J. Wu. 2019. Accurate Face Detection for High Performance. (2019). arXiv:1905.01585
- [309] G. Zhang, H. Han, S. Shan, X. Song, and X. Chen. 2018. Face Alignment across Large Pose via MT-CNN Based 3D Shape Reconstruction. In *Proceedings of the IEEE International Conference on Automatic Face & Gesture Recognition*. 210–217.
- [310] H. Zhang, B. S. Riggan, S. Hu, N. J. Short, and V. M. Patel. 2019. Synthesis of High-Quality Visible Faces from Polarimetric Thermal Faces using Generative Adversarial Networks. *International Journal of Computer Vision* 127

- (2019), 845–862.
- [311] J. Zhang and H. Hu. 2019. Stacked Hourglass Network Joint with Salient Region Attention Refinement for Face Alignment. In *Proceedings of the IEEE International Conference on Automatic Face & Gesture Recognition*. 1–7.
- [312] J. Zhang, S. Shan, M. Kan, and X. Chen. 2014. Coarse-to-fine auto-encoder networks (cfan) for real-time face alignment. In *Proceedings of the European Conference on Computer Vision*. 1–16.
- [313] J. Zhong, X. Wu, S. C. Hoi, and J. Zhu. 2020. Feature agglomeration networks for single stage face detection. *Neurocomputing* 380 (2020), 180–189.
- [314] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao. 2016. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Signal Processing Letters* 23, 10 (2016), 1499–1503.
- [315] L. Zhang, L. Lin, X. Wu, S. Ding, and L. Zhang. 2015. End-to-End Photo-Sketch Generation via Fully Convolutional Representation Learning. In *Proceedings of the ACM on International Conference on Multimedia Retrieval*. 627–634.
- [316] M. Zhang, Y. Li, N. Wang, Y. Chi, and X. Gao. 2020. Cascaded Face Sketch Synthesis Under Various Illuminations. *IEEE Transactions on Image Processing* 29 (2020), 1507–1521.
- [317] M. Zhang, R. Wang, X. Gao, J. Li, and D. Tao. 2019. Dual-Transfer Face Sketch–Photo Synthesis. *Trans. Image Process.* 28 (2019), 642–657.
- [318] S. Zhang, C. Chi, Z. Lei, and S. Z. Li. 2020. Refineface: Refinement neural network for high performance face detection. *IEEE Trans. Pattern Anal. Mach. Intell.* (2020).
- [319] S. Zhang, L. Wen, X. Bian, Z. Lei, and S. Z. Li. 2018. Single-Shot Refinement Neural Network for Object Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4203–4212.
- [320] S. Zhang, R. Zhu, X. Wang, H. Shi, F. Fu, S. Wang, T. Mei, and S. Z. Li. 2019. Improved Selective Refinement Network for Face Detection. (2019). arXiv:1901.06651
- [321] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Z. Li. 2017. FaceBoxes: a CPU real-time face detector with high accuracy. In *Proceedings of the IEEE International Joint Conference on Biometrics*. 1–9.
- [322] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Z. Li. 2017. S<sup>3</sup>FD: single shot scale-Invariant face detector. In *Proceedings of the IEEE International Conference on Computer Vision*. 192–201.
- [323] S. Zhang, X. Zhu, Z. Lei, X. Wang, and Stan Z Li. 2018. Detecting face with densely connected face proposal network. *Neurocomputing* 284 (2018), 119–127.
- [324] W. Zhang, X. Wang, and X. Tang. 2011. Coupled information-theoretic encoding for face photo-sketch recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 513–520.
- [325] X. Zhang, Z. Fang, Y. Wen, Z. Li, and Y. Qiao. 2017. Range Loss for Deep Face Recognition with Long-Tailed Training Data. In *Proceedings of the IEEE International Conference on Computer Vision*. 5419–5428.
- [326] X. Zhang, R. Zhao, Y. Qiao, X. Wang, and H. Li. 2019. AdaCos: Adaptively Scaling Cosine Logits for Effectively Learning Deep Face Representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10823–10832.
- [327] X. Zhang, R. Zhao, J. Yan, M. Gao, Y. Qiao, X. Wang, and H. Li. 2019. P2SGrad: Refined Gradients for Optimizing Deep Face Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9898–9906.
- [328] X. Zhang, X. Zhou, M. Lin, and J. Sun. 2018. ShuffleNet: An Extremely Efficient Convolutional Neural Network for Mobile Devices. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6848–6856.
- [329] Y. Zhang, W. Deng, M. Wang, J. Hu, X. Li, D. Zhao, and D. Wen. 2020. Global-Local GCN: Large-Scale Label Noise Cleansing for Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7728–7737.
- [330] Y. Zhang, X. Xu, and X. Liu. 2019. Robust and High Performance Face Detector. arXiv:1901.02350
- [331] Z. Zhang, L. Ping, C. L. Chen, and X. Tang. 2014. Facial Landmark Detection by Deep Multi-task Learning. In *Proceedings of the European Conference on Computer Vision*. 94–108.
- [332] H. Zhao, X. Ying, Y. Shi, X. Tong, J. Wen, and H. Zha. 2020. RDCFace: Radial Distortion Correction for Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7721–7730.
- [333] J. Zhao, Y. Cheng, Y. P. Cheng, Y. Yang, H. Lan, F. Zhao, L. Xiong, Y. Xu, J. Li, S. Pranata, S. Shen, J. Xing, H. Liu, S. Yan, and J. Feng. 2019. Look Across Elapse: Disentangled Representation Learning and Photorealistic Cross-Age Face Synthesis for Age-Invariant Face Recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 9251–9258.
- [334] J. Zhao, Y. Cheng, Y. Xu, L. Xiong, J. Li, F. Zhao, K. Jayashree, S. Pranata, S. Shen, J. Xing, S. Yan, and J. Feng. 2018. Towards Pose Invariant Face Recognition in the Wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2207–2216.
- [335] J. Zhao, L. Xiong, J. Karlekar, J. Li, F. Zhao, Z. Wang, S. Pranata, S. Shen, S. Yan, and J. Feng. 2017. Dual-Agent GANs for Photorealistic and Identity Preserving Profile Face Synthesis. In *Advances in neural information processing systems*. 65–75.

- [336] K. Zhao, K. Xu, and M. Cheng. 2019. RegularFace: Deep Face Recognition via Exclusive Regularization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1136–1144.
- [337] T. Zheng and W. Deng. 2018. Cross-pose lfw: A database for studying crosspose face recognition in unconstrained environments. *Beijing University of Posts and Telecommunications, Tech. Rep* (2018), 18–01.
- [338] T. Zheng, W. Deng, and J. Hu. 2017. Cross-age lfw: A database for studying cross-age face recognition in unconstrained environments. (2017). arXiv:1708.08197
- [339] Y. Zheng, D. K. Pal, and M. Savvides. 2018. Ring loss: Convex feature normalization for face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5089–5097.
- [340] Y. Zhong, J. Chen, and B. Huang. 2017. Toward End-to-End Face Recognition Through Alignment Learning. *IEEE Singal processing letters* 24, 8 (2017), 1213–1217.
- [341] Yaoyao Zhong and Weihong Deng. 2018. Deep Difference Analysis in Similar-looking Face recognition. In *Proceedings of the International Conference on Pattern Recognition (ICPR)*. IEEE, 3353–3358.
- [342] Y. Zhong, W. Deng, M. Wang, J. Hu, J. Peng, X. Tao, and Y. Huang. 2019. Unequal-Training for Deep Face Recognition With Long-Tailed Noisy Data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7804–7813.
- [343] E. Zhou, Z. Cao, and J. Sun. 2018. Gridface: Face rectification via learning local homography transformations. In *Proceedings of the European Conference on Computer Vision*. 3–19.
- [344] E. Zhou, H. Fan, Z. Cao, Y. Jiang, and Q. Yin. 2013. Extensive Facial Landmark Localization with Coarse-to-Fine Convolutional Network Cascade. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 386–391.
- [345] S. K. Zhou and D. Comaniciu. 2007. Shape regression machine. In *Biennial International Conference on Information Processing in Medical Imaging*. 13–25.
- [346] C. Zhu, Y. He, and M. Savvides. 2019. Feature Selective Anchor-Free Module for Single-Shot Object Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 840–849.
- [347] C. Zhu, R. Tao, K. Luu, and M. Savvides. 2018. Seeing small faces from robust anchor’s perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5127–5136.
- [348] C. Zhu, Y. Zheng, K. Luu, and M. Savvides. 2017. Cms-rcnn: contextual multi-scale region-based cnn for unconstrained face detection. In *Deep learning for biometrics*. 57–79.
- [349] J. Zhu, T. Park, P. Isola, and A. A. Efros. 2017. Unpaired Image-to-Image Translation Using Cycle-Consistent Adversarial Networks. In *Proceedings of the IEEE International Conference on Computer Vision*. 2242–2251.
- [350] M. Zhu, F. Shi, M. Zheng, and M. Sadiq. 2019. Robust facial landmark detection via occlusion-adaptive deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3486–3496.
- [351] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li. 2016. Face Alignment Across Large Poses: A 3D Solution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 146–155.
- [352] X. Zhu, H. Liu, Z. Lei, H. Shi, F. Yang, D. Yi, G. Qi, and S. Z. Li. 2019. Large-scale bisample learning on id versus spot face recognition. *International Journal of Computer Vision* 127, 6-7 (2019), 684–700.
- [353] X. Zhu and D. Ramanan. 2012. Face detection, pose estimation, and landmark localization in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2879–2886.
- [354] Zheng Zhu, Guan Huang, Jiankang Deng, Yun Ye, Junjie Huang, Xinze Chen, Jiagang Zhu, Tian Yang, Jiwen Lu, Dalong Du, and Jie Zhou. 2021. WebFace260M: A Benchmark Unveiling the Power of Million-Scale Deep Face Recognition. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2021)*, 10487–10497.
- [355] C. Zhuang, S. Zhang, X. Zhu, Z. Lei, J. Wang, and S. Z. Li. 2019. FLDet: A CPU Real-time Joint Face and Landmark Detector. In *Proceedings of the International Conference on Biometrics*. 1–8.
- [356] X. Zou, J. Kittler, and K. Messer. 2007. Illumination Invariant Face Recognition: A Survey. In *Proceedings of the IEEE International Conference on Biometrics: Theory, Applications, and Systems*. 1–8.
- [357] X. Zou, S. Zhong, L. Yan, X. Zhao, J. Zhou, and Y. Wu. 2019. Learning Robust Facial Landmark Detection via Hierarchical Structured Ensemble. In *Proceedings of the IEEE International Conference on Computer Vision*. 141–150.
- [358] O. Çeliktutan, S. Ulukaya, and B. Sankur. 2013. A comparative study of face landmarking techniques. *EURASIP Journal on Image and Video Processing* 2013 (2013), 1–27.

## A REPRESENTATIVE SURVEYS OF FACE RECOGNITION

A number of face recognition surveys have been published in the past twenty years. We summarize them in Table 12.

Table 12. Representative surveys of face recognition

Title	Year	Description
Face Recognition: A Literature Survey [231]	2003	Traditional image- and video-based methods in face recognition. Not covering deep face recognition.
Face Recognition from a Single Image per Person: A Survey [307]	2006	The methods to address the single sample problem in face recognition, not covering deep face recognition.
A Survey of Approaches and Challenges in 3D and Multi-modal 3D+2D Face Recognition [16]	2006	A survey of 3D and multi-modal face recognition, not covering deep face recognition.
Illumination Invariant Face Recognition: A Survey [356]	2007	Focus on illumination-invariant face recognition task, not covering deep face recognition.
A Survey of Face Recognition Techniques [7]	2009	Traditional face recognition methods on different modal face data, not covering deep face recognition.
A Comprehensive Survey on Pose-Invariant Face Recognition [49]	2016	Focus on pose-invariant face recognition task.
A Survey of Local Feature Methods for 3D Face Recognition [201]	2017	A review of feature extraction based methods for 3D face recognition.
Deep Learning for Understanding Faces [177]	2018	Provide a brief overview of the end-to-end deep face recognition, not covering the recent works.
Deep Face Recognition: A Survey [244]	2018	Focus on the deep face representation learning.
Past, Present, and Future of Face Recognition: A Review [2]	2020	A review of 2D and 3D face recognition, not covering end-to-end deep face recognition.

## B FACE DETECTION

### B.1 Single-stage and multi-stage face detectors

Fig. 9 illustrates the difference between single-stage and multi-stage face detectors. For comparison, the single-stage face detector accomplishes the detection processing directly from the feature maps, whereas the multi-stage face detector adopts a proposal stage to generate candidates and one or more stages to refine these candidates.

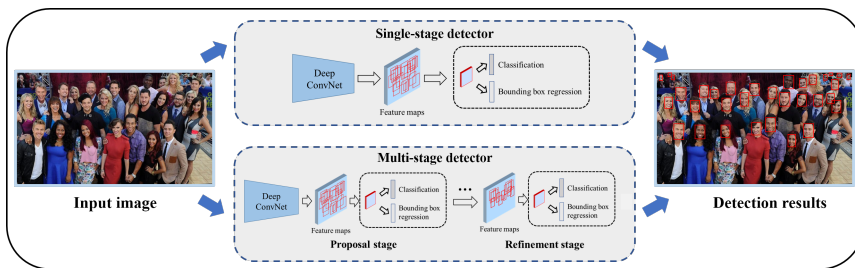


Fig. 9. The illustration of single-stage and multi-stage face detectors. The single-stage detector accomplishes the face detection directly from the feature maps, whereas the multi-stage detector adopts a proposal stage to generate candidates and one or more stages to refine these candidates.

### B.2 Performance comparison of CPU real-time face detection methods

Table 13 shows the running efficiency of CPU real-time face detection methods, among which the lightweight backbone [41, 279], rapidly digested convolutional layer [321, 323], knowledge distillation [98] and region-of-interest (RoI) convolution [27] are the common practices.

Table 13. Running efficiency of CPU real-time face detectors. “Accuracy (%)” denotes the true positive rate at 1000 false positives on Fddb.

Method	Publication	CPU-model	Speed (FPS)	Accuracy (%)
Faceboxes [321]	IJCB’17	E5-2660v3@2.60GHz	20	96.0
STN [27]	ECCV’16	i7-4770K	30	-
DCFPN [323]	Neurocomputing’18	2.60GHz	30	-
FBI [98]	ICB’19	E5-2660v3@2.60GHz	20	96.8
PCN [193]	CVPR’18	3.40GHz	29	-
PPN [301]	Information Sciences’19	i5	60	-
RetinaFace [41]	CVPR’19	i7-6700K	60	-
CenterFace [279]	arXiv’19	i7-6700@2.60GHz	30	98.0

## C FACE ALIGNMENT

### C.1 Hourglass network for facial landmark localization

Hourglass [168] is a bottom-up and top-down architecture, playing an important role in the deep stack of bottleneck blocks along with intermediate supervision. Fig. 10 is an illustration of stacked hourglass network.

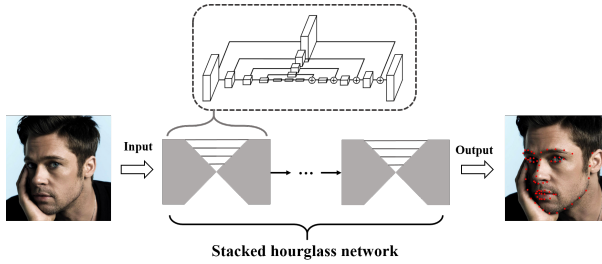


Fig. 10. An illustration of stacked hourglass network [168] for facial landmark localization. In each hourglass structure, the width (*i.e.*, feature channels) is consistent, and the boxes represent the residual modules.

### C.2 3D model fitting for facial landmark localization

As illustrated in Fig. 11, some 3D model fitting based methods employ cascaded regression manner with a dense 3D Morphable Model (3DMM) [15] to estimate the 3D face shape.

### C.3 Landmark-free face alignment

Landmark-free face alignment methods integrate the alignment transformation processing into DCNNs and output aligned face without relying on facial landmarks (Fig. 12).

## D FACE REPRESENTATION

Fig. 13 shows the pipeline of face representation training phase and test phase. In the training phase, two types of training supervision are widely used, *i.e.*, classification and feature embedding. As for test phase, there are two major tasks, *i.e.*, face verification or face identification.

In addition, as shown Fig. 14, we can observe that the publications of classification based training supervision exceed those of the feature embedding and hybrid methods with the growing scale of available face data. The reason is that the closed-set classification training on the large-scale datasets enables to approach open-set face recognition scenario.



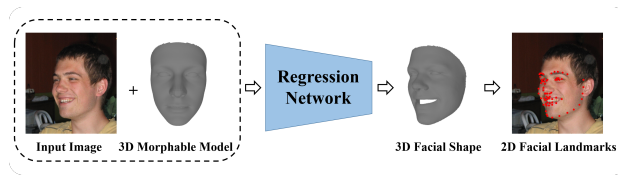


Fig. 11. The process of 3D model fitting for face alignment. A dense 3D Morphable Model is used to model a 2D face to 3D mesh. The regression network estimates the parameters of 3D shape and projection matrix, and then the 3D shape is projected onto the image plane to obtain the 2D landmarks.

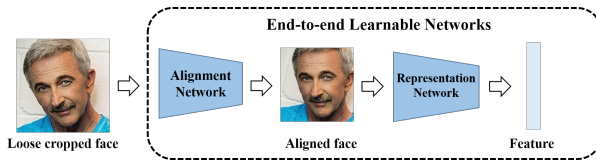


Fig. 12. An illustration of integrated framework that accomplishes landmark-free face alignment and representation computation.

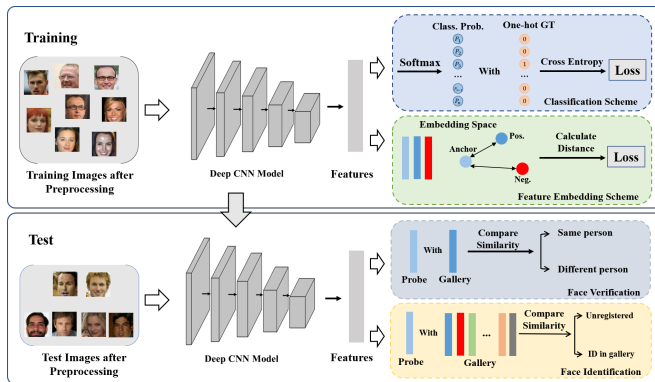


Fig. 13. The pipeline of face representation training phase and test phase. In the training phase, two schemes, *i.e.*, classification and feature embedding, are often used for learning face representation. In the test phase, face verification and face identification are the major tasks.

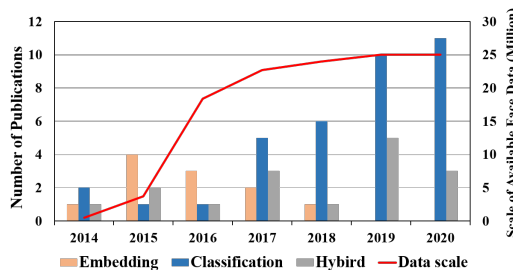


Fig. 14. The publication trend of three supervised face representation learning schemes with the growing scale of available face datasets from 2014 - 2020.