

Going Deeper Into Face Detection: A Survey

Shervin Minaee, Ping Luo, Zhe Lin, Kevin Bowyer

Abstract—Face detection is a crucial first step in many facial recognition and face analysis systems. Early approaches for face detection were mainly based on classifiers built on top of hand-crafted features extracted from local image regions, such as Haar Cascades and Histogram of Oriented Gradients. However, these approaches were not powerful enough to achieve a high accuracy on images of from uncontrolled environments. With the breakthrough work in image classification using deep neural networks in 2012, there has been a huge paradigm shift in face detection. Inspired by the rapid progress of deep learning in computer vision, many deep learning based frameworks have been proposed for face detection over the past few years, achieving significant improvements in accuracy. In this work, we provide a detailed overview of some of the most representative deep learning based face detection methods by grouping them into a few major categories, and present their core architectural designs and accuracies on popular benchmarks. We also describe some of the most popular face detection datasets. Finally, we discuss some current challenges in the field, and suggest potential future research directions.

Index Terms—Face Detection, Face Recognition, Deep Learning, Surveillance Systems, Convolutional Neural Networks.

1 INTRODUCTION

Face detection is an essential early step for tasks such as face recognition, facial attribute classification, face editing, and face tracking, and its performance has a direct impact on the effectiveness of those tasks [1], [2]. Although great improvements have been made in uncontrolled face detection over the past few decades, accurate and efficient face detection in the wild remains an open challenge. This is due to factors such as variations in poses, facial expressions, scale, illumination, image distortion, face occlusion, and other factors. Different from generic object detection, face detection features smaller variations in the aspect ratio, but much larger variations in scale (from several pixels to thousand pixels).

Early face detection efforts were mainly based on the classical approach, in which hand-crafted features were extracted from the image (or from sliding windows on the image) and were fed into a classifier (or ensemble of classifiers) to detect likely face regions. Two landmark classical works for face detection are the Haar Cascades classifier [3] and the Histogram of Oriented Gradients (HOG) followed by SVM [4]. These works represent great improvements on the state-of-the-art at their time. However, face detection accuracy was still limited on challenging images with multiple variation factors such as the ones shown in Fig 1.

With the great success of deep learning in computer vision, researchers have proposed several promising model architectures over the past 6-7 years. Inspired by the cascade of classifiers idea, many of the earlier deep-learning-based models were based on Cascade-CNN architectures. But with the introduction of several novel architectures for general object detection, many more recent deep-learning-based models have shifted toward single-Shot Detection, R-CNN based architectures, feature pyramid network (FPN) models,

and beyond.

Major surveys of face detection research up to circa 2000 include those by Yang et al. [6], Rowley et al.s [7], and Hjeltnäs and Low [8]. Zhang and Zhang survey progress in face detection over roughly the next decade, to about 2010 [9]. Zafeiriou et al [10] survey face detection research over roughly the next five years, to near the beginning of the deep learning wave, around 2015. One of their conclusions is that “even when allowing a relatively large number of false positives (around 1,000), there are still around 15-20% of faces that are not detected.” Our survey picks up where [10] ends, and covers the rapid progress in face detection from the beginning of the deep learning wave through the current time. Table 1 summarizes and compares the existing surveys with our work.

This paper provides a survey of the recent literature in deep-learning-based face detection, including more than fifty such detection methods. It provides a comprehensive review with insights into different aspects of these methods, including the training data, choice of network architectures, loss functions, training strategies, and their key contributions. These works are organized into the following categories, based on their main technical contributions to face detection:

- 1) Cascade-CNN Based Models
- 2) R-CNN and Faster-RCNN Based Models
- 3) Single Shot Detector Models
- 4) Feature Pyramid Network Based Models
- 5) Other models

The remainder of this survey is organized as follows: Section 2 overviews popular Deep Neural Network (DNN) architectures that serve as the backbones of many modern face detection algorithms. Section 3 reviews the most significant state-of-the-art deep learning based face detection models, and their main technical contributions. Section 4 summarizes the most popular benchmarks for face detection, their size and other characteristics. Section 5 lists popular metrics for evaluating deep-learning-based face detection models and also tabulates the performance of models on

- Shervin Minaee is a Machine Learning Lead at Snap Inc.
- Ping Luo is an assistant professor at the University of Hong Kong.
- Zhe Lin is a senior principal scientist at Adobe Research.
- Kevin Bowyer is the Schubmehl-Prein Family Professor of Computer Science and Engineering at the University of Notre Dame.



Fig. 1. Sample images from the “Wider-Face” face detection datasets, showing different variation factors. Courtesy of [5]

TABLE 1
Comparisons of face detection surveys since 1998.

No.	Survey Title	Reference	Year	Venue	Main Content
1	Neural network-based face detection	Rowley et al. [7]	1998	PAMI	A survey of 2-layered neural network for face detection.
2	Face detection: A survey	Hjelmås and Low [8]	2001	CVIU	A survey of traditional feature-based methods such as low-level cues (e.g. edges) and active shape models (e.g. Snakes), as well as image-based methods such as linear subspace methods.
3	Detecting faces in images: a survey	Yang et al. [6]	2002	PAMI	A survey of face detection from a single image, focusing on feature engineering and conventional classifiers such as EigenFace, Naive Bayes, and Support Vector Machine.
4	A survey of recent advances in face detection	Zhang et al. [9]	2010	Technical Report	A survey of Viola-Jones face detector and its variants.
5	A survey on face detection in the wild: past, present and future	Zafeiriou et al. [10]	2015	CVIU	A survey of handcrafted features, boosting and Support Vector Machine, deformable models in face detection before the wave of deep learning.
6	Going deeper into face detection: a survey (this work)	–	2021	–	A survey of the recent advanced deep-learning-based face detection, including more than fifty deep models for face detection.

those datasets. Section 6 discusses the main challenges and opportunities of deep learning-based face detection. Section 7 presents our conclusions.

2 OVERVIEW OF POPULAR DEEP LEARNING ARCHITECTURES

This section provides an overview of prominent DNN architectures used by the computer vision community, including convolutional neural networks, recurrent neural networks and long short-term memory, encoder-decoder and autoencoder models, and generative adversarial networks. Due to space limitations, several other DNN architectures that have been proposed, such as transformers, capsule networks, gated recurrent units, and spatial transformer networks, are not covered.

2.1 Convolutional Neural Networks (CNNs)

CNNs are among the most successful and widely used architectures in the deep learning community, especially for computer vision tasks. CNNs were initially proposed by Fukushima [11] in his seminal paper on the “Neocognitron”, which was based on Hubel and Wiesel’s hierarchical receptive field model of the visual cortex. Subsequently, Waibel

et al. [12] introduced CNNs with weights shared among temporal receptive fields and backpropagation training for phoneme recognition, and LeCun *et al.* [13] developed a practical CNN architecture for document recognition (Fig. 2). CNNs usually include three types of layers: i) convolutional layers, where a kernel (or filter) of weights is convolved to extract features; ii) nonlinear layers, which apply (usually element-wise) an activation function to feature maps, thus enabling the network to model nonlinear functions; and iii) pooling layers, which reduce spatial resolution by replacing small neighborhoods in a feature map with some statistical information about those neighborhoods (mean, max, etc.). The neuronal units in layers are locally connected; that is, each unit receives weighted inputs from a small neighborhood, known as the receptive field, of units in the previous layer. By stacking layers to form multi-resolution pyramids, the higher-level layers learn features from increasingly wider receptive fields. The main computational advantage of CNNs is that all the receptive fields in a layer share weights, resulting in a significantly smaller number of parameters than fully-connected neural networks. Some of the most well known CNN architectures include AlexNet [14], VGGNet [15], and ResNet [16].

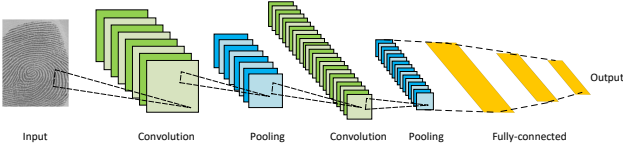


Fig. 2. Architecture of CNNs. From [13].

2.2 R-CNN Based Models

The Regional CNN (R-CNN) and its extensions have proven successful in object detection applications. In particular, the Faster R-CNN [17] architecture (Fig. 3) uses a region proposal network (RPN) that proposes bounding box candidates. The RPN extracts a Region of Interest (RoI), and an RoIPool layer computes features from these proposals to infer the bounding box coordinates and class of the object. Some extensions of R-CNN have been used to address the instance segmentation problem; i.e., the task of simultaneously performing object detection and segmentation.

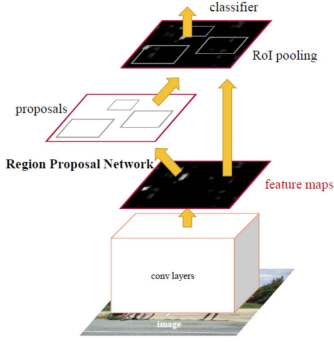


Fig. 3. Faster R-CNN architecture. Each image is processed by convolutional layers and its features are extracted, a sliding window is used in RPN for each location over the feature map, for each location, k ($k = 9$) anchor boxes are used (3 scales of 128, 256 and 512, and 3 aspect ratios of 1:1, 1:2, 2:1) to generate a region proposal; A cls layer outputs $2k$ scores to indicate whether or not there is an object for k boxes; A reg layer outputs $4k$ for the coordinates (box center coordinates, width and height) of k boxes. From [17].

2.3 Single Shot MultiBox Detector

Single Shot Detector (SSD) is a popular deep learning architecture proposed for object detection [18]. It discretizes the output space of bounding boxes into a set of default boxes over different aspect ratios and scales per feature map location. At prediction time, the network generates scores for the presence of each object category in each default box and produces adjustments to the box to better match the object shape. Additionally, the network combines predictions from multiple feature maps with different resolutions to naturally handle objects of various sizes. SSD is simple relative to methods that require object proposals because it completely eliminates proposal generation and subsequent pixel or feature resampling stages and encapsulates all computation in a single network. Figure 4 illustrates the high-level architecture of the original SSD model.

2.4 Feature Pyramid Network (FPN)

Feature pyramids are a basic component in recognition systems for detecting objects at different scales [19]. In this work, Lin et al. exploited the inherent multi-scale, pyramidal hierarchy of deep convolutional networks to construct feature pyramids with marginal extra cost. They developed a top-down architecture with lateral connections for building high-level semantic feature maps at all scales (called FPN), and demonstrated that FPN can bring significant improvements in several vision tasks. Fig 5 shows the high-level architecture of feature pyramid network proposed in [19].

2.5 Generative Adversarial Networks (GANs)

GANs [20] are a newer family of deep learning models. They consist of two networks—a generator and a discriminator (Fig. 7). In the conventional GAN, the generator network G learns a mapping from noise z (with a prior distribution) to a target distribution y , which is similar to the “real” samples. The discriminator network D attempts to distinguish the generated “fake” samples from the real ones. GAN training may be characterized as a minimax game between G and D , where D tries to minimize its classification error in distinguishing fake samples from real ones, hence maximizing a loss function, and G tries to maximize the discriminator network’s error, hence minimizing the loss function. Some of the GAN variants include Convolutional-GANs [21], conditional-GANs [22], Wasserstein-GANs [23], and CycleGAN [24].

3 OVERVIEW OF DIFFERENT DEEP FACE DETECTION MODELS

There are many approaches proposed for face detection using different deep learning architectures. In order to better summarize the existing works in a more comprehensive way, we grouped these works into a few prominent categories and review the major works of each category in the below sections:

- Cascade-CNN Based Models
- R-CNN Based Models
- Single Shot Detector Models
- Feature Pyramid Network Based Models
- Transformers Based Models
- Other Architectures

Figure 6 provides an illustration summarizing the most popular deep learning based face detection models from 2015 till 2021.

3.1 Cascade-CNN Based Models

Li et al. [25] proposed one of the early deep models for face detection, based on a convolutional neural network cascade. The proposed CNN cascade operates at multiple resolutions, quickly rejects the background regions in the fast low resolution stages, and carefully evaluates a small number of candidates in the last high resolution stage. To improve localization effectiveness, and reduce the number of candidates at later stages, they introduce a CNN-based calibration stage after each of the detection stages in the

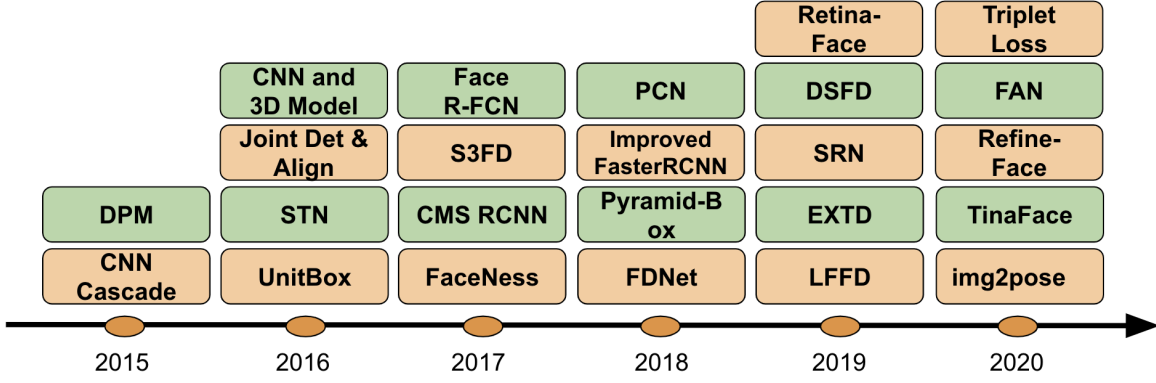


Fig. 6. The timeline of deep-learning-based face detection algorithms from 2015 to 2020.

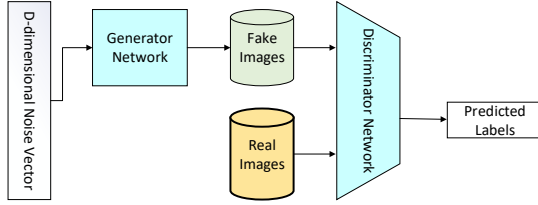


Fig. 7. Architecture of a GAN. Courtesy of Ian Goodfellow.

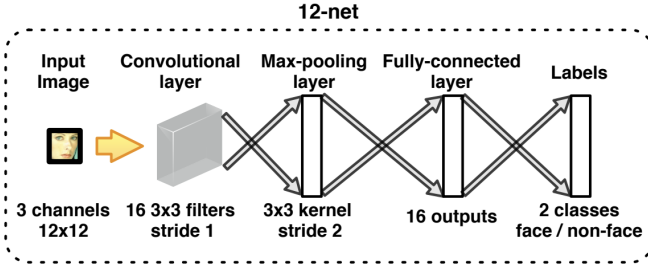


Fig. 8. The architecture of 12-layer Cascade-CNN model. Courtesy of [25].

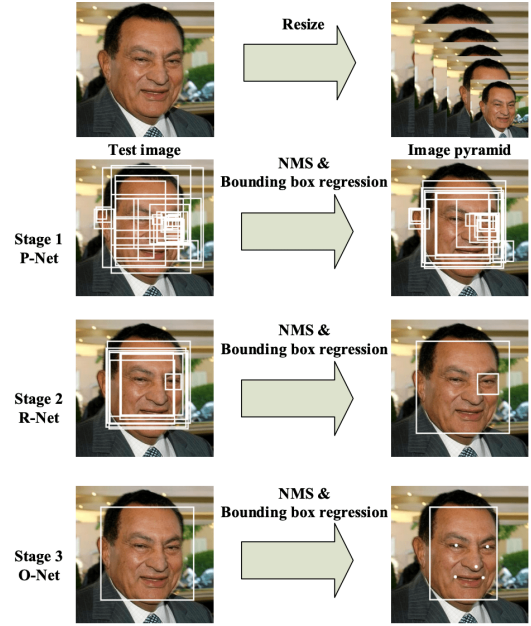


Fig. 9. The architecture of the proposed multi-task cascaded CNN. Courtesy of [26]

In another effort, [35], Wang et al. proposed a region-based fully convolutional network for face detection. They adopt a fully convolutional Residual Network (ResNet) as the backbone network, and exploit several new techniques including position-sensitive average pooling, multi-scale training and testing and on-line hard example mining strategy to improve the detection accuracy. Fig. 13 illustrates the architecture of the proposed R-FCN model.

Some other R-CNN based face detection methods include: "Face detection with different scales based on faster R-CNN" [36], "Face Detection Using Improved Faster RCNN" [37], and "Design of a Deep Face Detector by Mask R-CNN" [38].

3.3 Single Shot Detection Models

Single stage detection (SSD) is another popular and major direction in deep learning based face detection. Unlike two-stage proposal-classification detectors, such as R-CNN models, SSD detects faces in a single stage directly from the early convolutional layers in a classification network.

In [39], Najibi et al. proposed a "Single Stage Headless (SSH)" face detection framework, which achieved state of the art results on all subsets of Wider Faces, FDDB, and Pascal-Faces. Instead of relying on an image pyramid to detect faces with various scales, SSH is scale-invariant by design. It simultaneously detects faces with different scales in a single forward pass of the network, but from different layers. These properties make SSH fast and light-weight. Fig. 14 illustrates the architecture of the proposed SSH model.

Zhang et al. proposed "Single Shot Scale-invariant Face Detector (S³FD)" [40], to address scale variations better with a single deep neural network. Detection of for small faces is an especially common problem with anchor-based detectors. There are three main contributions in this work. First, a scale-equitable face detection framework is proposed to handle different scales of faces well. Second, the recall rate of small faces is improved by a scale compensation anchor-matching strategy. Third, the false positive rate of small faces is reduced via a max-out background label. Fig. 15 illustrates

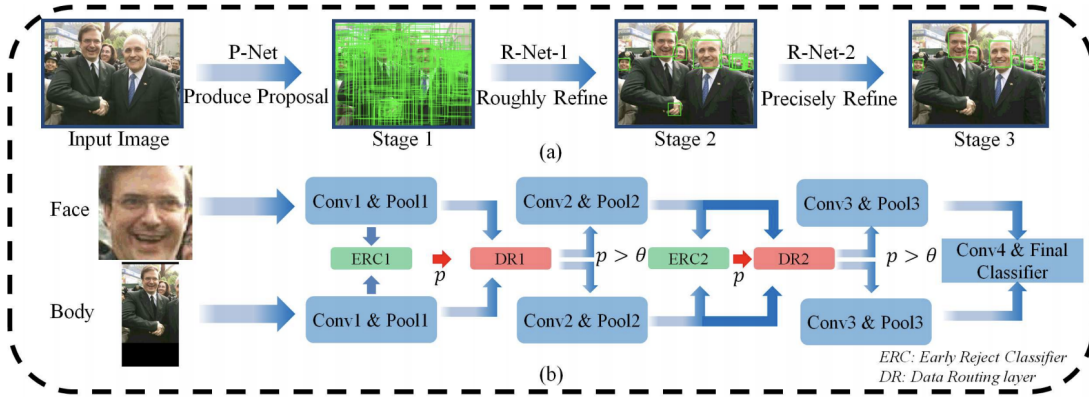


Fig. 10. The architecture of the proposed Inside Cascaded CNN. Courtesy of [27]

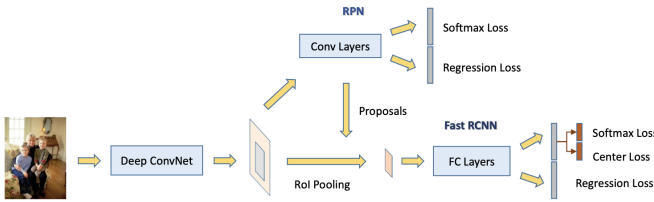


Fig. 11. The architecture of the Face-RCNN Model. Courtesy of [32]

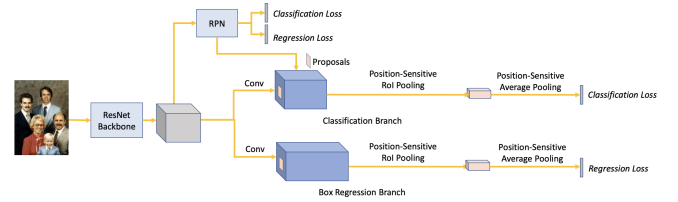


Fig. 13. The architecture of the R-FCN Model. Courtesy of [35]

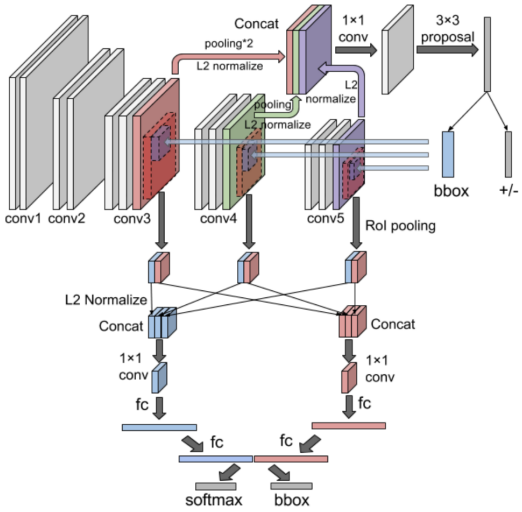


Fig. 12. The architecture of the CMS-RCNN Model. Courtesy of [34]

the architecture of the proposed S3FD model.

Zhang et al. proposed FaceBoxes [41], which is a real-time face detector based on the "Rapidly Digested Convolutional Layers (RDCL)" and the "Multiple Scale Convolutional Layers (MSCL)". RDCL is designed to enable FaceBoxes to achieve real-time speed. The MSCL aims at enriching the receptive fields and discretizing anchors over different layers to handle faces of various scales. They also propose a new anchor densification strategy to make different types of anchors have the same density on the image. This significantly improves the recall rate of small faces. One interesting fact about FaceBoxes is that its speed is invariant to the number of faces.

Hu and Ramanam proposed an interesting algorithm designed to address the problem with detecting small faces [42]. This model is able to detect a few hundred of small faces in a single image. One example of such a case is shown in Fig 16. They explored three aspects of the problem of finding small faces: the role of scale invariance, image resolution, and contextual reasoning. They took a different approach from previous works, and trained separate detectors for different scales. To maintain efficiency, detectors are trained in a multi-task fashion: they make use of features extracted from multiple layers of single (deep) feature hierarchy. They show that the context is very important and make use of massively-large receptive fields. The overall architecture of this model is shown in Fig 17.

Most recently, in [43], Zhang et al. proposed a single-shot refinement face detector called RefineFace. This framework is based on RetinaNet [44], with five proposed modules: Selective Two-step Regression (STR), Selective Two-step Classification (STC), Scale-aware Margin Loss (SML), Feature Supervision Module (FSM) and Receptive Field Enhancement (RFE). To enhance the regression ability for high location accuracy, STR coarsely adjusts locations and sizes of anchors from high level detection layers to provide better initialization for subsequent regressors. The overall architecture of this model is shown in Fig 18.

There are other works in this category, e.g. scale-friendly faces [45], YOLO-face: a real-time face detector [46], and "SANet: Smoothed attention network for single stage face detector" [47].

3.4 Feature Pyramid Network Based Models

Another popular family of deep models in face detection is based on feature pyramid networks [19]. Feature pyramid

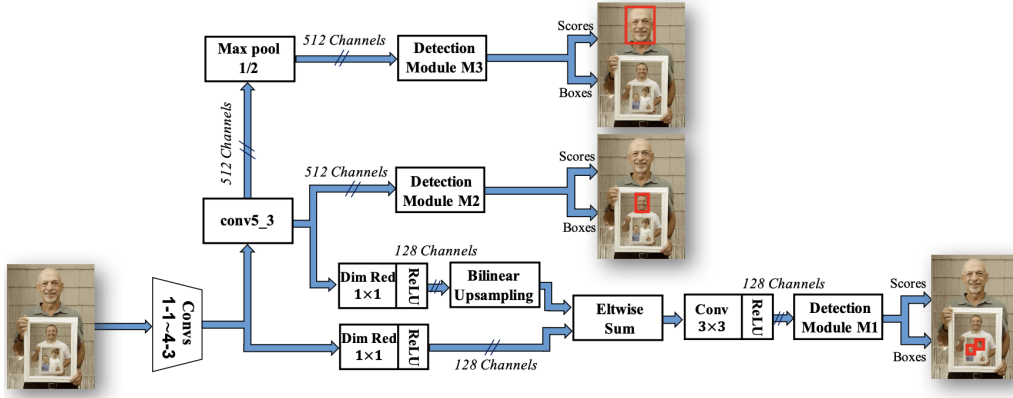


Fig. 14. The architecture of the proposed SSH Model. Courtesy of [39]

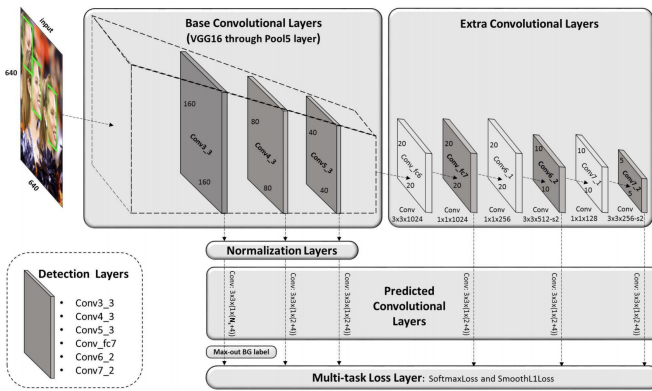


Fig. 15. The architecture of the proposed S3FD Model. Courtesy of [39]

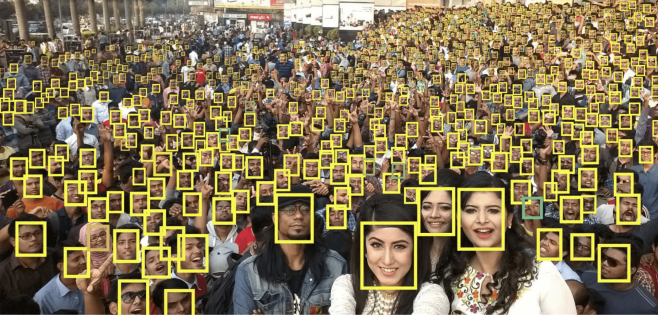


Fig. 16. The detected faces by the tiny face detector model for a sample image. Courtesy of [42]

networks have also been used for object detection, as well as semantic segmentation. A feature pyramid is a neural network structure which combines semantically weak features with semantically strong features using skip-connections.

In [48], Zhang et al. proposed "Feature Agglomeration Networks (FANet)", inspired by the Feature Pyramid Network, for single-stage face detection. The key idea of this framework is to exploit the inherent multi-scale features of a single convolutional neural network by aggregating higher-level semantic feature maps of different scales as contextual cues to augment lower-level feature maps via a hierarchical agglomeration manner at marginal extra computation cost.

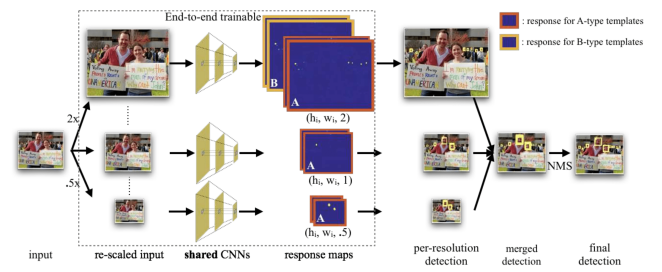


Fig. 17. The overall architecture of the model for detecting tiny faces. Courtesy of [42]

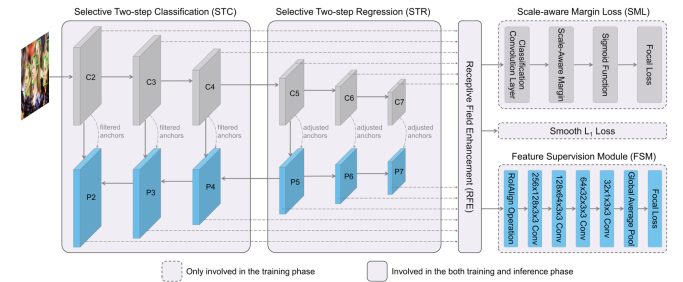


Fig. 18. The overall architecture of the RefineFace model for detecting faces. Courtesy of [43]

They also proposed a Hierarchical Loss to effectively train the FANet model. The overall architecture of FANet model is shown in Fig 25.

In another promising work, Tang et al. proposed Pyramid-Box [49], which is a context-assisted single shot face detector, and tries to address the challenge to detect small, blurred and partially occluded faces in uncontrolled environments. They improved the utilization of contextual information in the following three aspects: First, they designed a novel context anchor to supervise high-level contextual feature learning by a semi-supervised method, which they called PyramidAnchors. Second, they proposed a low-level FPN to combine adequate high-level context semantic feature and low-level facial feature together, which also allows the PyramidBox to predict faces of all scales in a single shot. Third, they introduced a context-sensitive structure to increase the capacity of prediction network to improve the

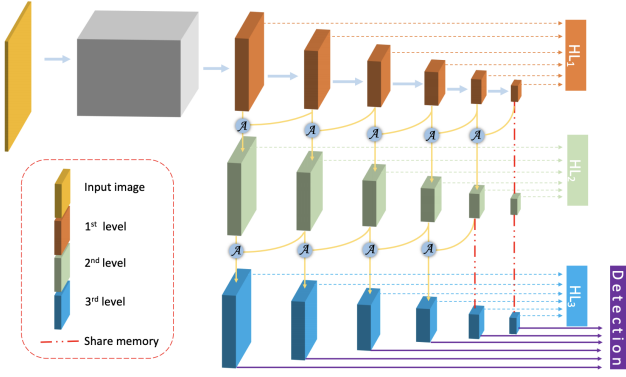


Fig. 19. The overall architecture of FANet model for face detection. Courtesy of [48]

final accuracy of output. The architecture of the proposed PyramidBox model is shown in Fig 20.

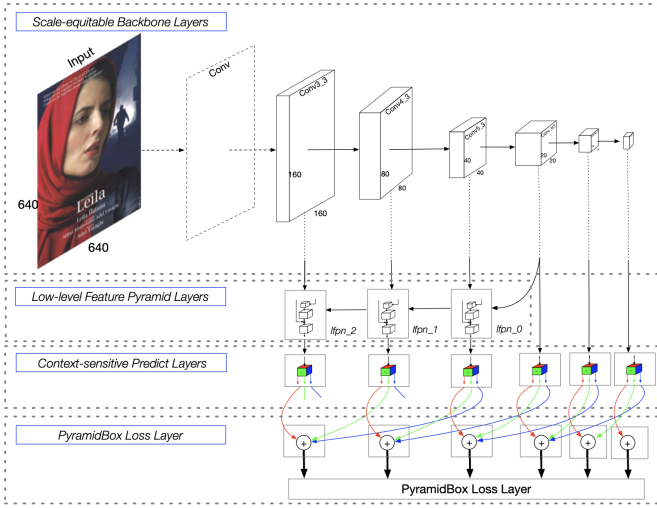


Fig. 20. The overall architecture of PyramidBox model for face detection. Courtesy of [49]

In [50], Chi et al. proposed a new single-shot model for face detection called "Selective Refinement Network (SRN)". SRN consists of two novel modules: the Selective Two-step Classification (STC) module and the Selective Two-step Regression (STR) module. The STC aims to filter out most simple negative anchors from low level detection layers to reduce the search space for the subsequent classifier, while the STR is designed to coarsely adjust the locations and sizes of anchors from high level detection layers to provide better initialization for the subsequent regressor. Furthermore, they also designed a Receptive Field Enhancement (RFE) block to provide a more diverse receptive field, which helps to better capture faces in some extreme poses. The network architecture of the proposed Selective Refinement Network model is shown in Fig 21.

In [51], Li et al. proposed a novel face detection algorithm called "dual shot face detector (DSFD)", with three key contributions: First, they developed a Feature Enhance Module (FEM) for enhancing the original feature maps to extend the single shot detector to dual shot detector. Second, they adopted Progressive Anchor Loss (PAL) computed

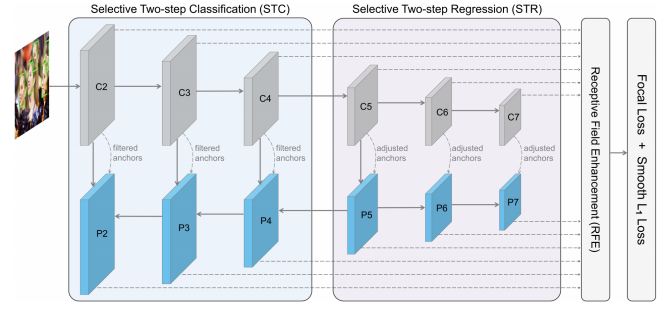


Fig. 21. The model architecture of Selective Refinement Network for face detection. Courtesy of [50]

by two different sets of anchors to effectively facilitate the features. And third, they used an Improved Anchor Matching (IAM) by integrating novel anchor assign strategy into data augmentation to provide better initialization for the regressor. Fig 22 provides the overall architecture of DSFD framework. The "Feature Enhance Module", which is a key component of this model, is shown in Fig 23. In [52], Li et al. proposed PyramidBox++, which is an improved version their earlier PyramidBox framework.

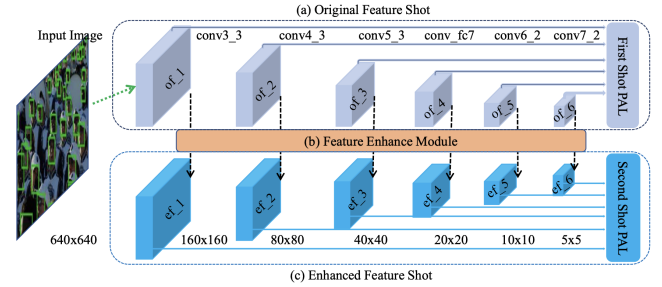


Fig. 22. The high-level architecture of dual shot face detector model. Courtesy of [51]

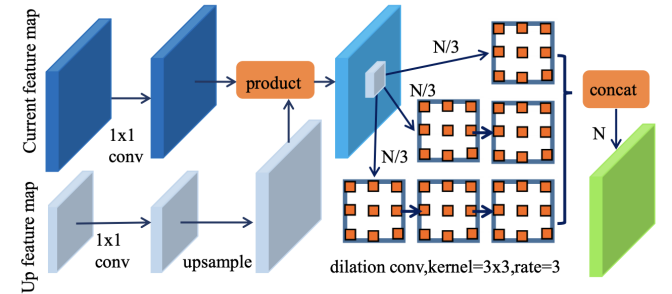


Fig. 23. The details of "Feature Enhance Module", used in DSFD framework. Courtesy of [51]

In [53], Dent et al. proposed a very popular single-stage face detection model of this category, which is called RetinaFace [53]. RetinaFace performs pixel-wise face localisation on various scales of faces by taking advantages of joint extra-supervised and self-supervised multi-task learning. One important contribution of this work is that they manually annotate five facial landmarks on the WIDER FACE dataset and observe significant improvement in hard face detection

with the assistance of this extra supervision signal. Another contribution is that they added a self-supervised mesh decoder branch for predicting a pixel-wise 3D shape face information in parallel with the existing supervised branches. They were able to achieve state of the art performance on several face detection benchmarks. The overall architecture of RetinaFace model is shown in Fig 24.

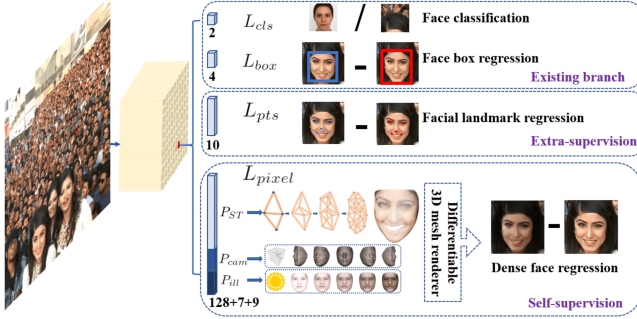


Fig. 24. The architecture of RetinaFace framework for face detection. Courtesy of [53]

In [54], Zhu et al. proposed TinaFace, a simple and strong baseline for face detection, which uses ResNet-50 as the feature extraction part, and 6 level Feature Pyramid Network to extract the multi-scale features of input image, followed by an Inception block to enhance the receptive field. One main purpose of this work was to show that there is no gap between face detection and generic object detection.

In [55], Najibi et al. proposed a novel approach for generating region proposals for performing face detection, called "Floating Anchor Region Proposal Network (FA-RPN)". Instead of classifying anchor boxes using pixel-level features in the convolutional feature map, they generate region proposals using a pooling-based approach. The overall architecture of FA-RPN is shown in Fig 25.

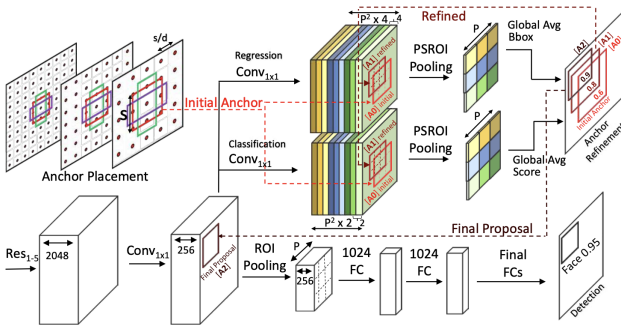


Fig. 25. The architecture of FA-RPN framework for face detection. Courtesy of [55]

There are also other works in this category, including robust and high performance face detector [56], Improved Selective Refinement Network for Face Detection [57], "Fast cascade face detection with pyramid network" [58], "Proposal pyramid networks for fast face detection" [59], etc.

3.5 Other Models

In this section, we cover approaches which do not fall into any of the categories mentioned above, or who

contribution is not in the modeling part but other factors (such as model optimization).

In [60], Zhang et al. developed a novel Automatic and Scalable Face Detector (ASFD), which is based on a combination of neural architecture search techniques as well as a new loss design. They proposed an automatic feature enhance module named Auto-FEM by improved differential architecture search, which allows efficient multi-scale feature fusion and context enhancement. They then used Distance-based Regression and Margin-based Classification (DRMC) multi-task loss to predict accurate bounding boxes and learn highly discriminative deep features. The overall architecture of the proposed AFSD model is shown in Fig 26.

In [61], Bai et al. proposed a framework for finding tiny faces with a generative adversarial network, which addresses the problem of super-resolving and refining jointly. Toward this goal, they developed an algorithm to directly generate a clear high-resolution face from a blurry small one by adopting a generative adversarial network (GAN). The overall architecture of this model is shown in Fig 27.

In [62], Yoo et al. proposed a new multi-scale face detector model containing a tiny number of parameters, called EXT-D. While existing multi-scale face detectors extract feature maps with different scales from a single backbone network, their method generates the feature maps by iteratively reusing a shared lightweight and shallow backbone network. This iterative backbone sharing strategy significantly reduces the number of parameters, and also provides the abstract image semantics captured from the higher stage of the network layers to the lower-level feature map.

In [63], Ranjan et al. proposed a framework called HyperFace, for simultaneous face detection, landmarks localization, pose estimation and gender recognition using deep convolutional neural networks (CNN). The proposed method fuses the intermediate layers of a deep CNN using a separate CNN followed by a multi-task learning algorithm that operates on the fused features. It exploits the synergy among the tasks which boosts up their individual performances. The architecture of this framework is shown in Fig 28.

In [64], Yang et al. proposed Faceness-Net, a deep convolutional neural network (CNN) for face detection leveraging on facial attributes based supervision. They observed a phenomenon that part detectors emerge within CNN trained to classify attributes from uncropped face images, without any explicit part supervision. The observation motivates a new method for finding faces through scoring facial parts responses by their spatial structure and arrangement. Fig 29 illustrates the architecture of the proposed Faceness-Net. The first stage of Faceness-Net applies attribute-aware networks to generate response maps of different facial parts. The maps are subsequently employed to produce face proposals. The second stage of Faceness-Net refines candidate windows generated from first stage using a multi-task convolutional neural network (CNN), where face classification and bounding box regression are jointly optimized.

In [65], Shi et al. proposed a real-time rotation-invariant face detection called "Progressive Calibration Networks". PCN consists of three stages, each of which not only distinguishes the faces from non-faces, but also calibrates the rotation-in-plane orientation of each face candidate to upright progressively. By dividing the calibration process

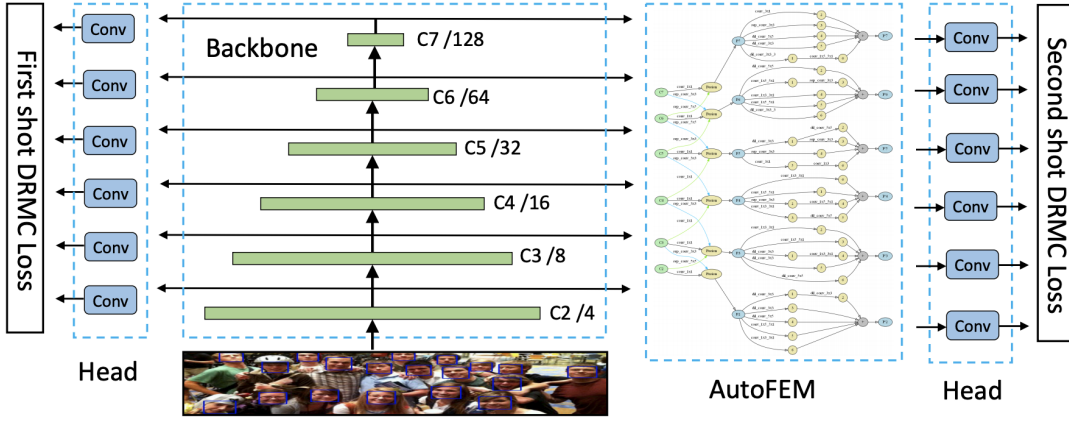


Fig. 26. The high-level idea of the AFSD framework for face detection. Courtesy of [60]

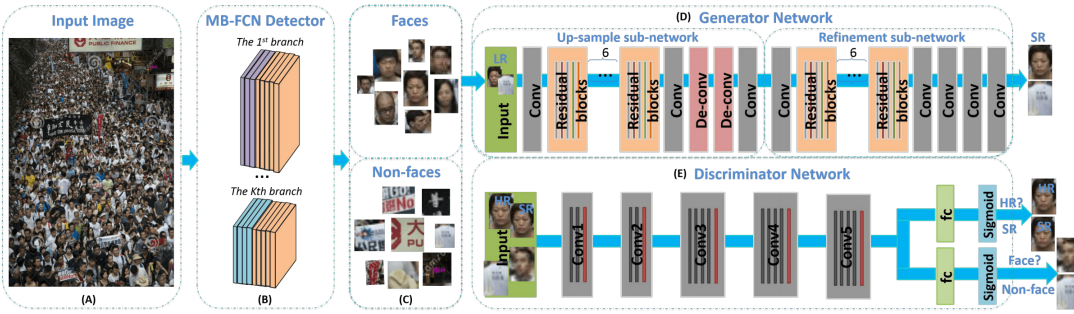


Fig. 27. The architecture of the proposed tiny face detector based on GANs. Courtesy of [61]

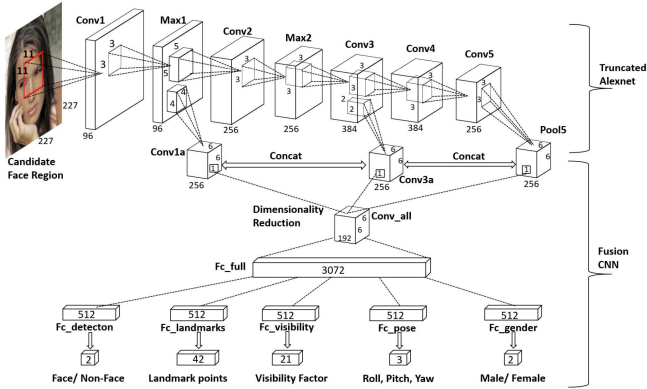


Fig. 28. The architecture of the HyperFace Model. Courtesy of [63]

into several progressive steps and only predicting coarse orientations in early stages, PCN can achieve precise and fast calibration. Figure 30 illustrates the high-level architecture of the PCN framework.

In [66], Cheng et al. proposed an anchor-free and non-maxima-supression-free object detection model, called weakly supervised multi-modal annotation segmentation (WSMA-Seg), which utilizes segmentation models to achieve an accurate and robust object detection without NMS. In addition to this, they proposed a multi-scale pooling segmentation (MSP-Seg) as the underlying segmentation model of WSMA-Seg to achieve a more accurate segmentation and

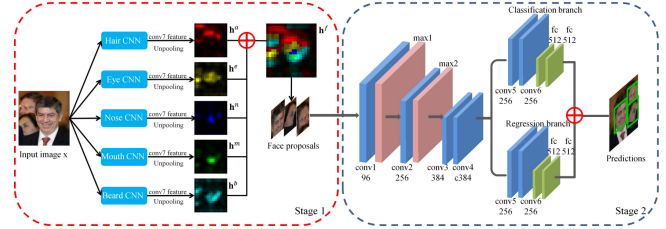


Fig. 29. The architecture of the Faceness-Net framework. Courtesy of [64]

to enhance the detection accuracy of WSMA-Seg. Through experimental results on multiple datasets they show that the proposed WSMA-Seg approach achieves promising results on several detection benchmarks, including widerFace.

In [67], Saha et al. introduced RNNPool, a novel pooling operator based on Recurrent Neural Networks (RNNs), that efficiently aggregates features over large patches of an image and rapidly downsamples activation maps, which they claimed to be more suitable for computer vision tasks, such as face detection. Empirical evaluation indicates that an RNNPool layer can effectively replace multiple blocks in a variety of architectures such as MobileNets, DenseNet when applied to standard vision tasks like face detection.

Some of the other popular models for face detection include: "A Light and Fast Face Detector for Edge Devices" [68], "Accurate face detection for high performance" [69],

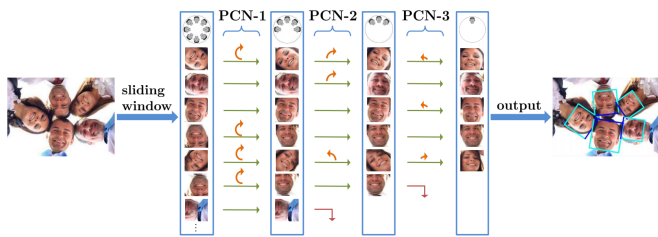


Fig. 30. The architecture of the progressive calibration networks (PCN) for face detection. The PCN progressively calibrates the RIP orientation of each face candidate to upright for better distinguishing faces from non-faces. Specifically, PCN-1 first identifies face candidates and calibrates those facing down to facing up. Then the rotated face candidates are further distinguished and calibrated to an upright range of $[-45, 45]$ in PCN-2, shrinking the RIP ranges by half again. Finally, PCN-3 makes the accurate final decision for each face candidate to determine whether it is a face and predict the precise RIP angle. Courtesy of [65]

UnitBox [70], "HAMBox: Delving Into Mining High-Quality Anchors on Face Detection" [71], "Joint face detection and facial motion re-targeting for multiple faces" [72], "Hierarchical attention for part-aware face detection" [73], "Group Sampling for Scale Invariant Face Detection" [74], "Dafe-fd: Density aware feature enrichment for face detection" [75], "Triple loss for hard face detection" [76], "BlazeFace: Sub-millisecond Neural Face Detection on Mobile GPUs" [77], and Face Detection with End-to-End Integration of a ConvNet and a 3D Model [78].

It is worth providing some high-level comparison between different categories of face detection models. Cascade-CNN based models are designed for efficient high performance face detectors which may be suitable for deployment in edge devices, or camera capture applications. However, these models is hard to bring the state of the art accuracy on challenging cases with low image quality, non-standard poses and lightings. On the other hand, the general object detector-based pipelines such as R-CNN can provide much better accuracy due to the more capacity and learning power of the model architectures used. The single-shot or anchor-free detectors such as SSD, FCOS, etc. provides a good tradeoff between the accuracy and efficiency. For backbone architectures, feature pyramid-based architectures have been shown to be outperforming standard CNN backbones for object detection without much degradation in efficiency. Furthermore, recently transformer-based architectures have shown state of the art results on object detection, and provides a potential promising avenue for future extensions of existing models for boosting the performance of face detection to the next level. Finally other approaches tackling specific challenges in face detection such as GAN-based approaches, or PCN detector specifically designed for better handling of rotated faces, are complementary to existing mainstream approaches.

4 FACE DETECTION BENCHMARK DATASETS

Several datasets have been proposed for face detection over the past decades. Here, we provide an overview of the most widely used datasets in the recent literature. Important factors to consider in comparing datasets include: number

of images and faces, size range of faces, amount of meta-data specified for each face, and range of face/image quality conditions represented.

4.1 FDDB

Face Detection Data Set (also known as FDDB) contains the annotations for 5171 faces in 2845 images taken from the Labelled Faces in the Wild data set [79]. FDDB contains a wide range of difficult elements, including occlusions, difficult poses, and low resolution and out-of-focus faces. The specifications of face regions are provided as elliptical regions. They proposed an evaluations metric based on ROC, coarse and precise score metrics to match between prediction and ground-truth. Three sample images from this dataset are shown in Fig 31.



Fig. 31. Three sample images from FDDB dataset, and the ground-truth face regions. Courtesy of [79]

4.2 WIDER FACE (2016)

The Wider Face dataset [5] was introduced after FDDB and contains 32,203 images and 393,703 labeled faces with a high degree of variability in scale, pose and occlusion, making it a very challenging dataset. The images are selected from the publicly available Wider dataset. The Wider Face dataset is organized based on 61 event classes. For each event class, they randomly selected 40%, 10%, 50% data as training, validation and testing sets accordingly. They adopted the same evaluation metric employed in the PASCAL VOC dataset. Similar to some of the other datasets, they do not release bounding box ground truth for the test images. Users are required to submit final prediction files, which they shall proceed to evaluate. Fig 32 illustrates some of the sample images from this dataset. As we can see there is a high degree of variability in scale, pose, occlusion, expression, appearance and illumination.

4.3 PASCAL Face

PASCAL face dataset contains 1,335 labeled faces in 851 images with large face appearance and pose variations [80]. It is collected from PASCAL person layout test subset. This dataset is small compared to other face detection datasets.

4.4 MALF

Multi-Attribute Labelled Faces (MALF) is the first face detection dataset that supports fine-grained evaluation. MALF consists of 5,250 images and 11,931 faces [81]. Besides bounding box annotations, each face contains other annotations such as: pose deformation level of yaw, pitch and roll (small, medium, large); other facial attributes: gender(female, male, unknown), is-Wearing-Glasses, is-occluded and is-Exaggerated-Expression. Fig 33 shows some of the sample images and annotations from MALF.

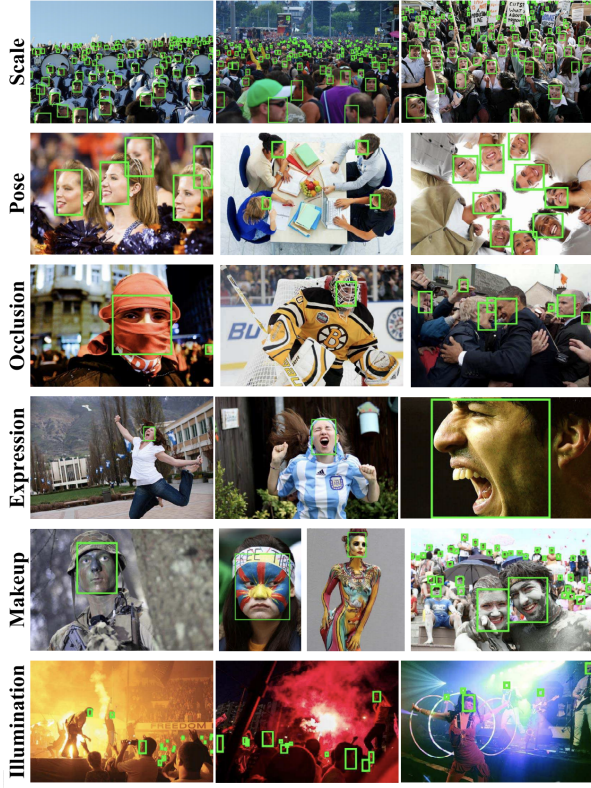


Fig. 32. Some of the sample images from Wider-Face dataset, and the ground-truth face bounding boxes. Courtesy of [5]



Fig. 33. Two sample images and their annotations from MAF dataset. Courtesy of [81]

4.5 UFDD (2018)

Unconstrained Face Detection Dataset (UFDD) [82] contains a total of 6,425 images with 10,897 face-annotations. It involves key degradations or conditions including: rain, snow, haze, lens impediments, blur, illumination variations, and distractors. Fig 34 shows some sample images from this dataset with different variations and the corresponding annotations.

4.6 VGGFace2 Dataset

VGGFace2 dataset contains 3.31 million images of 9131 subjects, with an average of 362.6 images for each subject [83]. Images are downloaded from Google Image Search and

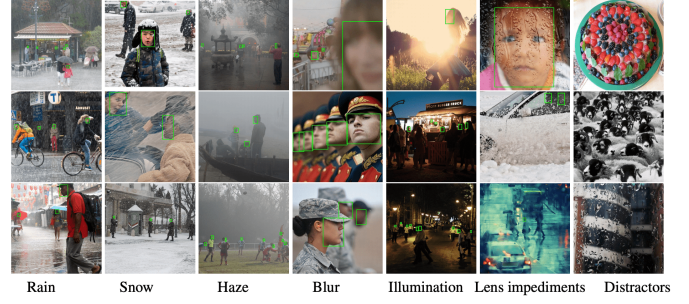


Fig. 34. Sample images from UFDD dataset. Courtesy of [82]

have large variations in pose, age, illumination, ethnicity and profession (e.g. actors, athletes, politicians). Faces in this dataset are detected using the model provided by [26].

5 EXPERIMENTAL PERFORMANCE

In this section, we first introduce some of the popular metrics used to evaluate the performance of face detection models. Then we summarize the quantitative performance of the more promising DL-based face detection models on popular datasets.

5.1 Evaluation Metrics

There are a few metrics which are widely used for evaluating the performance of face detection models, including: Average Precision (AP), Precision-Recall (PR) Curve, and receiver operating characteristic (ROC) curve. We give brief definitions of these metrics below.

5.1.1 Precision-Recall Curve

Before getting into the description of precision-recall curve, we need to give the mathematical definition of precision and recall metrics. These metrics are defined as Eq 1:

$$\text{Precision} = \frac{TP}{TP + FP}; \quad \text{Recall} = \frac{TP}{TP + FN}, \quad (1)$$

Now assuming that the model output is a score in $[0,1]$, we can compare the output with a threshold to get the class label. Therefore, for each threshold value we can find the corresponding precision and recall values. **Precision-Recall curve** shows the precision as a function of recall, for all different threshold values.

5.1.2 Average Precision (AP)

As mentioned above, a model's precision refers to precision at a particular decision threshold. Average precision calculates the average precision at all such possible thresholds, which is also similar to the area under the precision-recall curve. It is a useful metric to compare how well model's predictions are, without considering any specific decision threshold.

5.1.3 ROC Curve

The receiver operating characteristic (ROC) curve is plot which shows the performance of a model as a function of its cut-off threshold (similar to precision-recall curve). It essentially shows the true positive rate (TPR) against the false positive rate (FPR) for various threshold values. In general,

TABLE 2

The performance of face detection models on different versions of Wider-Face dataset.

Method	Easy	Medium	Hard
Faceness [64]	71.3	63.4	34.5
Multiscale Cascade CNN [5]	69.1	66.4	42.4
CMS-RCNN [34]	90.2	87.4	64.3
LFFD [68]	89.6	86.5	77.0
img2pose [84]	90.0	89.1	83.9
S3FD [40]	92.8	91.3	84.4
EXTD [62]	91.2	90.3	85.0
FACE R-FCN [35]	94.3	93.1	87.6
SRN [50]	95.9	94.8	89.6
FDNet [37]	95.0	93.9	89.6
DSFD [51]	96.0	95.3	90.0
PyramidBox [49]	95.6	94.6	90.0
AlnoFace [69]	96.5	95.7	91.2
RetinaFace [53]	-	-	91.4
TinaFace [54]	-	-	92.4

TABLE 3

The performance of face detection models on the FDDB dataset.

Method	AP
CascadeCNN [25]	85.7
Joint-Cascade [28]	86.3
HyperFace [63]	90.1
Faceness [64]	90.3
DP2MFD [50]	90.3
UnitBox [70]	95.1
FaceBoxes [41]	96.0
Faster R-CNN [17]	96.1
DPSSD [85]	96.1
LFFD [68]	97.3
S3FD [40]	98.3
PyramidBox [49]	98.7
SRN [50]	98.8
FACE R-FCN [35]	99.0
DSFD [51]	99.1

the lower the cut-off threshold on positive class, the more samples predicted as positive class, i.e. higher true positive rate (recall) and also higher false positive rate (corresponding to the right side of this curve). Therefore, there is a trade-off between how high the recall could be versus how much we want to bound the error (FPR).

5.2 Quantitative Performance of DL-Based Face Detection Models

In this section we tabulate the performance of several of the previously discussed algorithms on popular face detection benchmarks. Table 2 provides the performance of some of the prominent works on Wider-Face dataset. It is worth mentioning that there are three versions available for the Wider-Face test set. We report the performance of each model on all three versions, when they are available. As we can see, there has been a huge progress on the performance of deep models on this dataset. Tables 3 and 4 provide the performance of some of the deep learning based face detection models on FDDB, and PASCAL Face datasets, respectively. As it can be seen, the performances on these two datasets are typically higher than those on Wider-Face.

As we can see the performance of the recent models on the FDDB and PASCAL face datasets is higher than 99%, which

TABLE 4

The performance of face detection models on PASCAL Face dataset.

Method	AP
Headhunter	89.63
DPM [86]	90.29
Faceness [64]	92.11
STN [30]	94.10
HyperFace [63]	96.20
FaceBoxes [41]	96.30
S3FD [40]	98.49
Anchor-based [87]	99.00
SRN [50]	99.09

means that even with a much more powerful and novel model we will not be able to see a significant quantitative gain over the previous models. Hence, there is a real need for developing of larger and harder datasets of face detection.

In addition to the average precision of the models listed above, we provide the Precision-Recall curves of those models on WiderFace dataset (in Fig 35), and their Receiver operating characteristic (ROC) curves on FDDB benchmark (in Fig 36).

We also provide a comparison of some of the prominent face detection models, in terms of F1-score, in Fig 37

6 CHALLENGES AND OPPORTUNITIES

With no doubt, face detection has seen a great improvement over the last few years, thanks to new deep learning based models. However, several challenges lie ahead. We will next discuss some of the promising research directions that we believe will help in further advancing face detection.

6.1 Detection Robustness on Tiny Faces

Several researchers have already focused specific efforts on detecting faces that appear in tiny regions in images. However, there is still much room for improving the accuracy of these models for detecting tiny faces in general, and the more challenging cases of tiny faces with heavy occlusion, or facial disguise.

6.2 Face Occlusion

Occlusion is one of the most critical challenges in face detection. Various factors can lead to occlusion, such as: accessories (Goggles, Cap), medical masks, beard/moustache, facial disguise, or blockage by another person or object. Although many of the current face detection models handle occlusion to some extent, it is still very challenging to detect faces with heavy occlusion. There are two avenues to improve face detection models with occlusion. First, to develop new large-scale datasets of faces with occlusion. Second, to develop model architectures which are better curated toward detecting faces with occlusion.

6.3 Accurate Lightweight Models

Face detection is an important step in many real-world problems and use-cases, and therefore it is becoming deployed in various mobile applications. Hence, it is important for these models to be lightweight and memory efficient, in addition

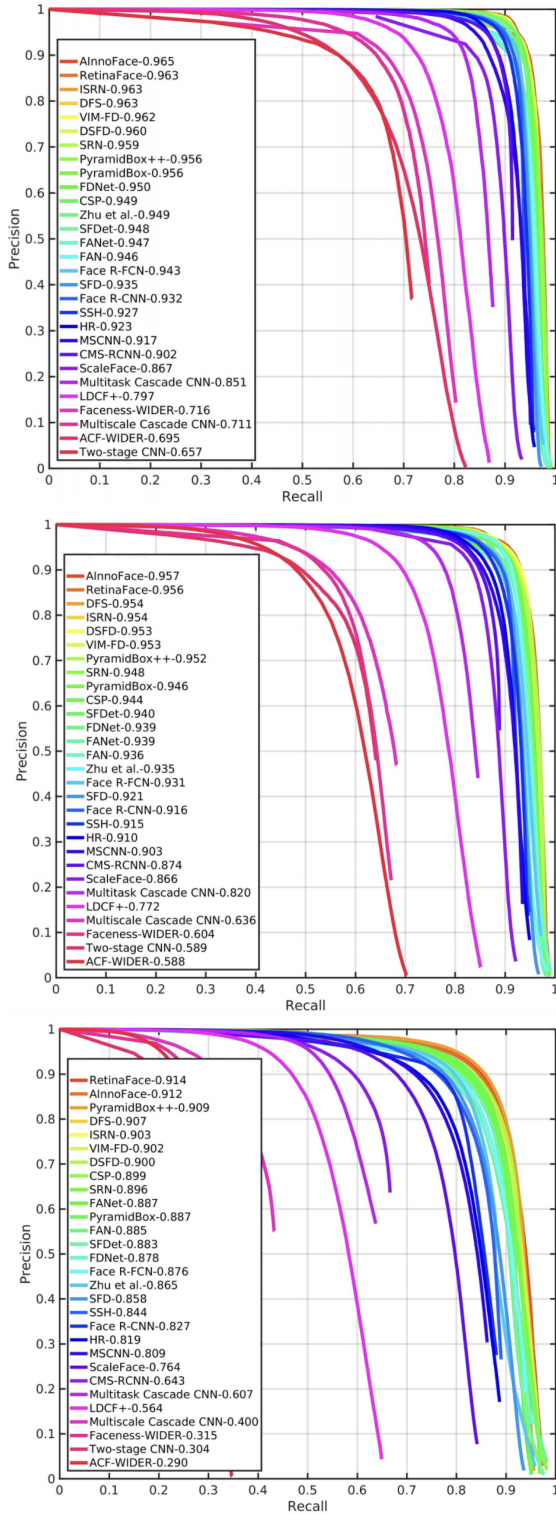


Fig. 35. The top, middle, and bottom figures present the Precision-Recall curves of promising face-detection models on WiderFace Easy/Medium/Hard test sets, respectively. Courtesy of [43].

to obtaining high accuracy. This can be achieved either by investigating novel lightweight model architectures in the first place, or by applying model compression and neural architecture search techniques to already-developed models.

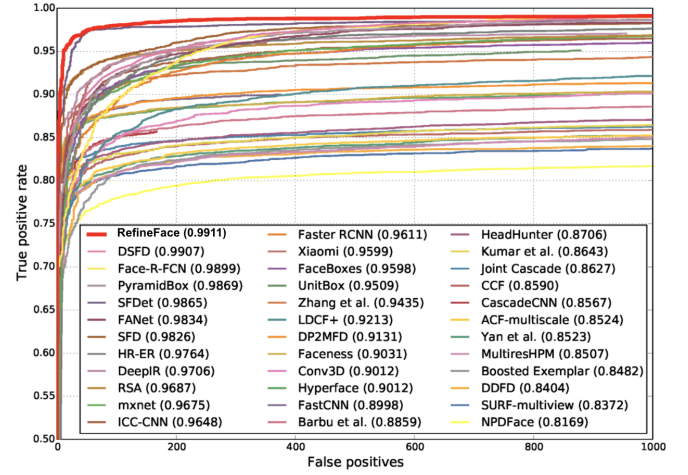


Fig. 36. The Receiver operating characteristic (ROC) curves of face detection models on Fddb benchmark. Courtesy of [43]

6.4 Few-Shot Face Detection

Most of the current face detection models are trained on very large-scale datasets of annotated faces, and there have not been many deep-learning based works for face detection from few labeled samples. This could be specially useful for new use-cases for which a large-scale dataset is not available, such as animation face detection, or detecting all faces in the crowd of images (for which there usually exist more than 100 faces per image).

6.5 Interpretable Deep Models

While DL-based models have achieved promising performance on challenging face detection benchmarks, there remain open questions about these models. For example, what exactly are deep models learning? What is a minimal neural architecture that can achieve a certain accuracy on a given dataset? How should we interpret the features learned by these models? Although there are some techniques available to visualize the learned convolutional kernels of these models, a comprehensive study of the underlying dynamics of these models is lacking. A better understanding of the theoretical aspects of these models can enable us to develop better models curated for face detection.

6.6 Face Detection Bias Reduction

The topic of “bias”, or accuracy disparity across demographic groups, is a hot issue currently. Researchers are actively engaged on this topic for face recognition [88], [89], [90], [91] and for analytics such as gender-from-face [92]. However, we know of no work that has focused specifically on possible bias in face detection. Experimental study of this topic presents some non-trivial challenges. The meta-data for an experimental dataset should not be created using a face detection algorithm, or the dataset may inherit any bias in the algorithm. Also, factors that are likely to affect accuracy in general should be balanced across demographic groups in the dataset, so that observed disparities can be inferred to be due to demographic factors.

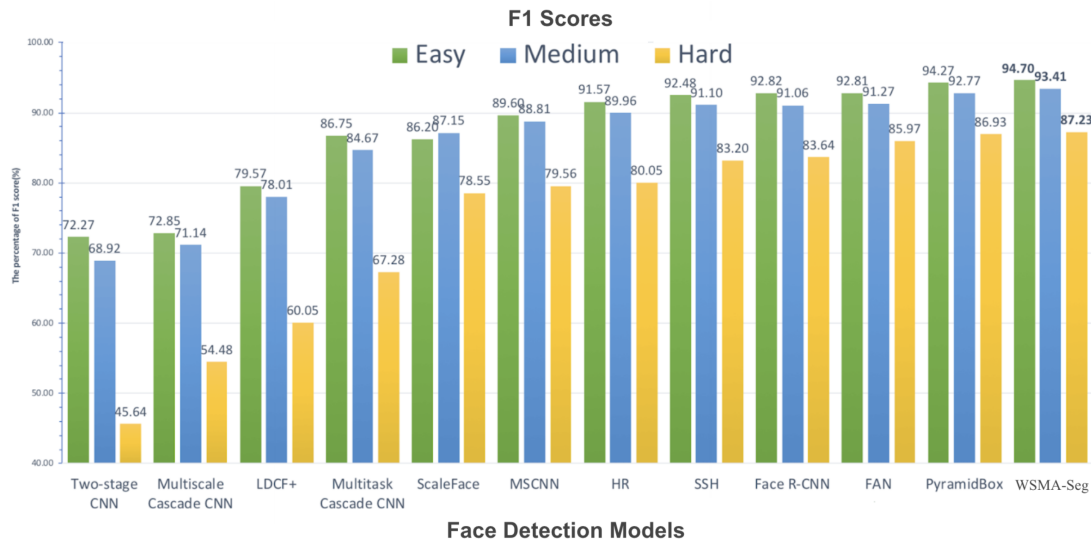


Fig. 37. The F1 scores of face detections models on the WiderFace benchmark. Courtesy of [66]

7 SUMMARY AND CONCLUSIONS

We have surveyed recent face detection methods based on deep learning models, which have achieved promising results on various face detection benchmarks and enabled the usage of these models on real-world applications. We categorized these models into several architectural groups: Cascaded CNN models, R-CNN based models, SSD-based models, and FPN-based models, and identified their main technical contributions. We also provided an overview of some of the popular face detection benchmarks, such as Wider-Face, FDDB, and PASCAL Face. We then summarized the quantitative performance of these models on these popular benchmarks. Finally, we discussed some of the open challenges and promising directions for deep-learning-based face detection in the coming years.

ACKNOWLEDGMENTS

The authors would like to thank Aleksei Stoliar for his comments and suggestions regarding this work.

REFERENCES

- [1] M. Wang and W. Deng, "Deep face recognition: A survey," *arXiv preprint arXiv:1804.06655*, 2018.
- [2] S. Minaee, A. Abdolrashidi, H. Su, M. Bennamoun, and D. Zhang, "Biometric recognition using deep learning: A survey," *arXiv preprint arXiv:1912.00271*, 2019.
- [3] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *Proceedings of the 2001 IEEE computer society conference on computer vision and pattern recognition*. CVPR 2001, vol. 1. IEEE, 2001, pp. I-I.
- [4] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR'05)*, vol. 1. Ieee, 2005, pp. 886-893.
- [5] S. Yang, P. Luo, C.-C. Loy, and X. Tang, "Wider face: A face detection benchmark," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5525-5533.
- [6] M.-H. Yang, D. J. Kriegman, and N. Ahuja, "Detecting faces in images: a survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 1, pp. 34-58, 2002.
- [7] H. A. Rowley, S. Baluja, and T. Kanade, "Neural network-based face detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 1, pp. 23-38, 1998.
- [8] E. Hjelmås and B. K. Low, "Face detection: A survey," *Computer Vision and Image Understanding*, vol. 83, no. 3, pp. 236-274, 2001.
- [9] C. Zhang and Z. Zhang, "A survey of recent advances in face detection," *MSR-TR-2010-66*, 2010.
- [10] S. Zafeiriou, C. Zhang, and Z. Zhang, "A survey on face detection in the wild: Past, present and future," *Computer Vision and Image Understanding*, vol. 138, pp. 1-24, 2015.
- [11] K. Fukushima, "Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position," *Biological cybernetics*, vol. 36, no. 4, pp. 193-202, 1980.
- [12] A. Waibel, T. Hanazawa, G. Hinton, K. Shikano, and K. J. Lang, "Phoneme recognition using time-delay neural networks," *IEEE transactions on acoustics, speech, and signal processing*, vol. 37, no. 3, pp. 328-339, 1989.
- [13] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner *et al.*, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278-2324, 1998.
- [14] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097-1105.
- [15] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [16] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770-778.
- [17] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *arXiv preprint arXiv:1506.01497*, 2015.
- [18] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21-37.
- [19] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117-2125.
- [20] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672-2680.
- [21] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.
- [22] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.
- [23] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein generative

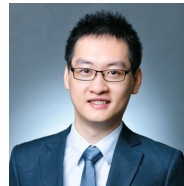
- adversarial networks," in *International conference on machine learning*. PMLR, 2017, pp. 214–223.
- [24] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [25] H. Li, Z. Lin, X. Shen, J. Brandt, and G. Hua, "A convolutional neural network cascade for face detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 5325–5334.
- [26] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016.
- [27] K. Zhang, Z. Zhang, H. Wang, Z. Li, Y. Qiao, and W. Liu, "Detecting faces using inside cascaded contextual cnn," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3171–3179.
- [28] H. Qin, J. Yan, X. Li, and X. Hu, "Joint training of cascaded cnn for face detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 3456–3465.
- [29] R. Qi, R.-S. Jia, Q.-C. Mao, H.-M. Sun, and L.-Q. Zuo, "Face detection method based on cascaded convolutional networks," *IEEE Access*, vol. 7, pp. 110740–110748, 2019.
- [30] D. Chen, G. Hua, F. Wen, and J. Sun, "Supervised transformer network for efficient face detection," in *European Conference on Computer Vision*. Springer, 2016, pp. 122–138.
- [31] H. Jiang and E. Learned-Miller, "Face detection with the faster r-cnn," in *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)*. IEEE, 2017, pp. 650–657.
- [32] H. Wang, Z. Li, X. Ji, and Y. Wang, "Face r-cnn," *arXiv preprint arXiv:1706.01061*, 2017.
- [33] X. Sun, P. Wu, and S. C. Hoi, "Face detection using deep learning: An improved faster rcnn approach," *Neurocomputing*, vol. 299, pp. 42–50, 2018.
- [34] C. Zhu, Y. Zheng, K. Luu, and M. Savvides, "Cms-rcnn: contextual multi-scale region-based cnn for unconstrained face detection," in *Deep learning for biometrics*. Springer, 2017, pp. 57–79.
- [35] Y. Wang, X. Ji, Z. Zhou, H. Wang, and Z. Li, "Detecting faces using region-based fully convolutional networks," *arXiv preprint arXiv:1709.05256*, 2017.
- [36] W. Wu, Y. Yin, X. Wang, and D. Xu, "Face detection with different scales based on faster r-cnn," *IEEE transactions on cybernetics*, vol. 49, no. 11, pp. 4017–4028, 2018.
- [37] C. Zhang, X. Xu, and D. Tu, "Face detection using improved faster rcnn," *arXiv preprint arXiv:1802.02142*, 2018.
- [38] O. Cakiroglu, C. Ozer, and B. Günsel, "Design of a deep face detector by mask r-cnn," in *2019 27th Signal Processing and Communications Applications Conference (SIU)*. IEEE, 2019, pp. 1–4.
- [39] M. Najibi, P. Samangouei, R. Chellappa, and L. S. Davis, "Ssh: Single stage headless face detector," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 4875–4884.
- [40] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Z. Li, "S3fd: Single shot scale-invariant face detector," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 192–201.
- [41] —, "Faceboxes: A cpu real-time face detector with high accuracy," in *2017 IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 2017, pp. 1–9.
- [42] P. Hu and D. Ramanan, "Finding tiny faces," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 951–959.
- [43] S. Zhang, C. Chi, Z. Lei, and S. Z. Li, "Refineface: Refinement neural network for high performance face detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [44] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [45] S. Yang, Y. Xiong, C. C. Loy, and X. Tang, "Face detection through scale-friendly deep convolutional networks," *arXiv preprint arXiv:1706.02863*, 2017.
- [46] W. Chen, H. Huang, S. Peng, C. Zhou, and C. Zhang, "Yolo-face: a real-time face detector," *The Visual Computer*, pp. 1–9, 2020.
- [47] L. Shi, X. Xu, and I. A. Kakadiaris, "Sanet: Smoothed attention network for single stage face detector," in *2019 International Conference on Biometrics (ICB)*. IEEE, 2019, pp. 1–7.
- [48] J. Zhang, X. Wu, S. C. Hoi, and J. Zhu, "Feature agglomeration networks for single stage face detection," *Neurocomputing*, vol. 380, pp. 180–189, 2020.
- [49] X. Tang, D. K. Du, Z. He, and J. Liu, "Pyramidbox: A context-assisted single shot face detector," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 797–813.
- [50] C. Chi, S. Zhang, J. Xing, Z. Lei, S. Z. Li, and X. Zou, "Selective refinement network for high performance face detection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 8231–8238.
- [51] J. Li, Y. Wang, C. Wang, Y. Tai, J. Qian, J. Yang, C. Wang, J. Li, and F. Huang, "Dsf: dual shot face detector," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5060–5069.
- [52] Z. Li, X. Tang, J. Han, J. Liu, and R. He, "Pyramidbox++: High performance detector for finding tiny face," *arXiv preprint arXiv:1904.00386*, 2019.
- [53] J. Deng, J. Guo, Y. Zhou, J. Yu, I. Kotsia, and S. Zafeiriou, "Retinaface: Single-stage dense face localisation in the wild," *arXiv preprint arXiv:1905.00641*, 2019.
- [54] Y. Zhu, H. Cai, S. Zhang, C. Wang, and Y. Xiong, "Tinaface: Strong but simple baseline for face detection," *arXiv preprint arXiv:2011.13183*, 2020.
- [55] M. Najibi, B. Singh, and L. S. Davis, "Fa-rpn: Floating region proposals for face detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7723–7732.
- [56] Y. Zhang, X. Xu, and X. Liu, "Robust and high performance face detector," *arXiv preprint arXiv:1901.02350*, 2019.
- [57] S. Zhang, R. Zhu, X. Wang, H. Shi, T. Fu, S. Wang, T. Mei, and S. Z. Li, "Improved selective refinement network for face detection," *arXiv preprint arXiv:1901.06651*, 2019.
- [58] D. Zeng, F. Zhao, S. Ge, and W. Shen, "Fast cascade face detection with pyramid network," *Pattern Recognition Letters*, vol. 119, pp. 180–186, 2019.
- [59] D. Zeng, H. Liu, F. Zhao, S. Ge, W. Shen, and Z. Zhang, "Proposal pyramid networks for fast face detection," *Information Sciences*, vol. 495, pp. 136–149, 2019.
- [60] B. Zhang, J. Li, Y. Wang, Y. Tai, C. Wang, J. Li, F. Huang, Y. Xia, W. Pei, and R. Ji, "Asfd: Automatic and scalable face detector," *arXiv preprint arXiv:2003.11228*, 2020.
- [61] Y. Bai, Y. Zhang, M. Ding, and B. Ghanem, "Finding tiny faces in the wild with generative adversarial network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 21–30.
- [62] Y. Yoo, D. Han, and S. Yun, "Ext: Extremely tiny face detector via iterative filter reuse," *arXiv preprint arXiv:1906.06579*, 2019.
- [63] R. Ranjan, V. M. Patel, and R. Chellappa, "Hyperface: A deep multi-task learning framework for face detection, landmark localization, pose estimation, and gender recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 1, pp. 121–135, 2017.
- [64] S. Yang, P. Luo, C. C. Loy, and X. Tang, "Faceness-net: Face detection through deep facial part responses," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 8, pp. 1845–1859, 2017.
- [65] X. Shi, S. Shan, M. Kan, S. Wu, and X. Chen, "Real-time rotation-invariant face detection with progressive calibration networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2295–2303.
- [66] Z. Cheng, Y. Wu, Z. Xu, T. Lukasiewicz, and W. Wang, "Segmentation is all you need," *arXiv preprint arXiv:1904.13300*, 2019.
- [67] O. Saha, A. Kusupati, H. V. Simhadri, M. Varma, and P. Jain, "Rnnpool: Efficient non-linear pooling for ram constrained inference," *arXiv preprint arXiv:2002.11921*, 2020.
- [68] Y. He, D. Xu, L. Wu, M. Jian, S. Xiang, and C. Pan, "Lffd: A light and fast face detector for edge devices," *arXiv preprint arXiv:1904.10633*, 2019.
- [69] F. Zhang, X. Fan, G. Ai, J. Song, Y. Qin, and J. Wu, "Accurate face detection for high performance," *arXiv preprint arXiv:1905.01585*, 2019.
- [70] J. Yu, Y. Jiang, Z. Wang, Z. Cao, and T. Huang, "Unitbox: An advanced object detection network," in *Proceedings of the 24th ACM international conference on Multimedia*, 2016, pp. 516–520.
- [71] Y. Liu, X. Tang, J. Han, J. Liu, D. Rui, and X. Wu, "Hambox: Delving into mining high-quality anchors on face detection," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2020, pp. 13 043–13 051.
- [72] B. Chaudhuri, N. Vedapant, and B. Wang, "Joint face detection and facial motion retargeting for multiple faces," in *Proceedings of*

the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 9719–9728.

- [73] S. Wu, M. Kan, S. Shan, and X. Chen, “Hierarchical attention for part-aware face detection,” *International Journal of Computer Vision*, vol. 127, no. 6, pp. 560–578, 2019.
- [74] X. Ming, F. Wei, T. Zhang, D. Chen, and F. Wen, “Group sampling for scale invariant face detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3446–3456.
- [75] V. A. Sindagi and V. Patel, “Dafe-fd: Density aware feature enrichment for face detection,” in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2019, pp. 2185–2195.
- [76] Z. Fang, J. Ren, S. Marshall, H. Zhao, Z. Wang, K. Huang, and B. Xiao, “Triple loss for hard face detection,” *Neurocomputing*, vol. 398, pp. 20–30, 2020.
- [77] V. Bazarevsky, Y. Kartynnik, A. Vakunov, K. Raveendran, and M. Grundmann, “Blazeface: Sub-millisecond neural face detection on mobile gpus,” *arXiv preprint arXiv:1907.05047*, 2019.
- [78] Y. Li, B. Sun, T. Wu, and Y. Wang, “Face detection with end-to-end integration of a convnet and a 3d model,” in *European Conference on Computer Vision*. Springer, 2016, pp. 420–436.
- [79] V. Jain and E. Learned-Miller, “Fddb: A benchmark for face detection in unconstrained settings,” UMass Amherst technical report, Tech. Rep., 2010.
- [80] J. Yan, X. Zhang, Z. Lei, and S. Z. Li, “Face detection by structural models,” *Image and Vision Computing*, vol. 32, no. 10, pp. 790–799, 2014.
- [81] B. Yang, J. Yan, Z. Lei, and S. Z. Li, “Fine-grained evaluation on face detection in the wild,” in *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, vol. 1. IEEE, 2015, pp. 1–7.
- [82] H. Nada, V. A. Sindagi, H. Zhang, and V. M. Patel, “Pushing the limits of unconstrained face detection: a challenge dataset and baseline results,” in *2018 IEEE 9th International Conference on Biometrics Theory, Applications and Systems (BTAS)*. IEEE, 2018, pp. 1–10.
- [83] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, “Vggface2: A dataset for recognising faces across pose and age,” in *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*. IEEE, 2018, pp. 67–74.
- [84] V. Albiero, X. Chen, X. Yin, G. Pang, and T. Hassner, “img2pose: Face alignment and detection via 6dof, face pose estimation,” *arXiv preprint arXiv:2012.07791*, 2020.
- [85] R. Ranjan, A. Bansal, J. Zheng, H. Xu, J. Gleason, B. Lu, A. Nanduri, J.-C. Chen, C. D. Castillo, and R. Chellappa, “A fast and accurate system for face detection, identification, and verification,” *IEEE Transactions on Biometrics, Behavior, and Identity Science*, vol. 1, no. 2, pp. 82–96, 2019.
- [86] S. Liao, A. K. Jain, and S. Z. Li, “A fast and accurate unconstrained face detector,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 2, pp. 211–223, 2015.
- [87] Z. Zhang, W. Shen, S. Qiao, Y. Wang, B. Wang, and A. Yuille, “Robust face detection via learning small faces on hard images,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020, pp. 1361–1370.
- [88] V. Albiero, K. S. Krishnapriya, K. Vangara, K. Zhang, M. C. King, and K. W. Bowyer, “Analysis of gender inequality in face recognition accuracy,” in *2020 IEEE Winter Applications of Computer Vision Workshops (WACVW)*, 2020, pp. 81–89.
- [89] J. G. Cavazos, P. J. Phillips, C. D. Castillo, and A. J. O’Toole, “Accuracy comparison across face recognition algorithms: Where are we on measuring race bias?” *IEEE Transactions on Biometrics, Behavior, and Identity Science*, vol. 3, no. 1, pp. 101–111, 2021.
- [90] K. S. Krishnapriya, V. Albiero, K. Vangara, M. C. King, and K. W. Bowyer, “Issues related to face recognition accuracy varying based on race and skin tone,” *IEEE Transactions on Technology and Society*, vol. 1, no. 1, pp. 8–20, 2020.
- [91] P. Dhar, J. Gleason, H. Sour, C. D. Castillo, and R. Chellappa, “An adversarial learning algorithm for mitigating gender bias in face recognition,” *arXiv preprint arXiv:2006.07845*, 2020.
- [92] P. Drodowski, C. Rathgeb, A. Dantcheva, N. Damer, and C. Busch, “Demographic bias in biometrics: A survey on an emerging challenge,” *IEEE Transactions on Technology and Society*, vol. 1, no. 2, pp. 89–103, 2020.



Shervin Minaee is a machine learning lead in the computer vision team at Snapchat, Inc. He received his PhD in Electrical Engineering and Computer Science from New York University, in 2018. His research interests include computer vision, image segmentation, biometric recognition, and applied deep learning. He has published more than 40 papers and patents during his PhD. He previously worked as a research scientist at Samsung Research, AT&T Labs, Huawei Labs, and as a data scientist at Expedia group. He has been a reviewer for more than 20 computer vision related journals from IEEE, ACM, Elsevier, and Springer. He has won several awards, including the best research presentation at Samsung Research America in 2017 and the Verizon Open Innovation Challenge Award in 2016.



Ping Luo is an Assistant Professor in the department of computer science, The University of Hong Kong (HKU). He received his PhD degree in 2014 from Information Engineering, the Chinese University of Hong Kong (CUHK), supervised by Prof. Xiaou Tang and Prof. Xiaogang Wang. He was a Postdoctoral Fellow in CUHK from 2014 to 2016. He joined SenseTime Research as a Principal Research Scientist from 2017 to 2018. His research interests are machine learning and computer vision. He has published 70+ peer-reviewed articles in top-tier conferences and journals such as TPAMI, IJCV, ICML, ICLR, CVPR, and NIPS. His work has high impact with 7,000 citations according to Google Scholar. He has won a number of competitions and awards such as the first runner up in 2014 ImageNet ILSVRC Challenge, the first place in 2017 DAVIS Challenge on Video Object Segmentation, Gold medal in 2017 Youtube 8M Video Classification Challenge, the first place in 2018 Drivable Area Segmentation Challenge for Autonomous Driving, 2011 HK PhD Fellow Award, and 2013 Microsoft Research Fellow Award (ten PhDs in Asia).



Zhe Lin is Senior Principal Scientist in Creative Intelligence Lab, Adobe Research. He received his Ph.D. degree in Electrical and Computer Engineering from University of Maryland at College Park in May 2009. Prior to that, he obtained his M.S. degree in Electrical Engineering and Computer Science from Korea Advanced Institute of Science and Technology in August 2004, and B.Eng. degree in Automation from University of Science and Technology of China. He has been a member of Adobe Research since May 2009. His

research interests include computer vision, image processing, machine learning, deep learning, artificial intelligence. He has served as a reviewer for many computer vision conferences and journals since 2009, and recently served as an Area Chair for WACV 2018, CVPR 2019, ICCV 2019, CVPR 2020, ECCV 2020, ACM Multimedia 2020.



Kevin W. Bowyer (Fellow, IEEE) received the Ph.D. degree in computer science from Duke University, Durham, NC, USA. He is the Schubmehl-Prein Family Professor of computer science and engineering with the University of Notre Dame, Notre Dame, IN, USA. Prof. Bowyer received the Technical Achievement Award from the IEEE Computer Society, with the citation “for pioneering contributions to the science and engineering of biometrics.” He served as the Editor-in-Chief for the IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE. He is currently serving as the Editor-in-Chief for the IEEE TRANSACTIONS ON BIOMETRICS, BEHAVIOR AND IDENTITY SCIENCE. In 2019, he was elected as a fellow of the American Association for the Advancement of Science (AAAS).

He is currently serving as the Editor-in-Chief for the IEEE TRANSACTIONS ON BIOMETRICS, BEHAVIOR AND IDENTITY SCIENCE. In 2019, he was elected as a fellow of the American Association for the Advancement of Science (AAAS).