

# Anchor Cascade for Efficient Face Detection

Baosheng Yu and Dacheng Tao, *Fellow, IEEE*

**Abstract**—Face detection is essential to facial analysis tasks such as facial reenactment and face recognition. Both cascade face detectors and anchor-based face detectors have translated shining demos into practice and received intensive attention from the community. However, cascade face detectors often suffer from a low detection accuracy, while anchor-based face detectors rely heavily on very large networks pre-trained on large scale image classification datasets such as ImageNet [1], which is not computationally efficient for both training and deployment. In this paper, we devise an efficient anchor-based cascade framework called anchor cascade. To improve the detection accuracy by exploring contextual information, we further propose a context pyramid maxout mechanism for anchor cascade. As a result, anchor cascade can train very efficient face detection models with a high detection accuracy. Specifically, comparing with a popular CNN-based cascade face detector MTCNN [2], our anchor cascade face detector greatly improves the detection accuracy, e.g., from 0.9435 to 0.9704 at 1k false positives on FDDB, while it still runs in comparable speed. Experimental results on two widely used face detection benchmarks, FDDB and WIDER FACE, demonstrate the effectiveness of the proposed framework.

**Index Terms**—Multi-scale Anchors, Cascade Face Detection.

## I. INTRODUCTION

**F**ACE detection, aiming to locate all faces in a given image, is indispensable to a variety of facial analysis tasks, such as facial point detection [3], facial reenactment [4], [5], and face recognition [6], [7], [8]. While the seminal cascade face detector [9] showed impressive performance for real-time near-frontal face detection, face detection in real-world scenarios remains very challenging due to large variations in scale, illumination, expression, pose, and occlusion [10].

The cascade framework is suited to efficient face detection due to its fundamentals, i.e., an extremely unbalanced binary classification problem with a majority of negative samples [9]. A cascade face detector rejects most of easy negative samples in the early stages using simple classifiers. As a result, cascade face detectors are generally extremely fast [11]. To further boost cascade face detector performance in real-world scenarios, robust features are crucial for detecting non-frontal faces with large variations in scale and occlusion [12]. Recently, convolutional neural networks (CNNs) have been effectively applied to image classification [13] due to their great capacity for representation learning [14]. Inspired by this, CNNs have been introduced to traditional cascade face detectors by replacing the original node classifiers, e.g., AdaBoost, with CNNs [15]. However, cascade face detectors

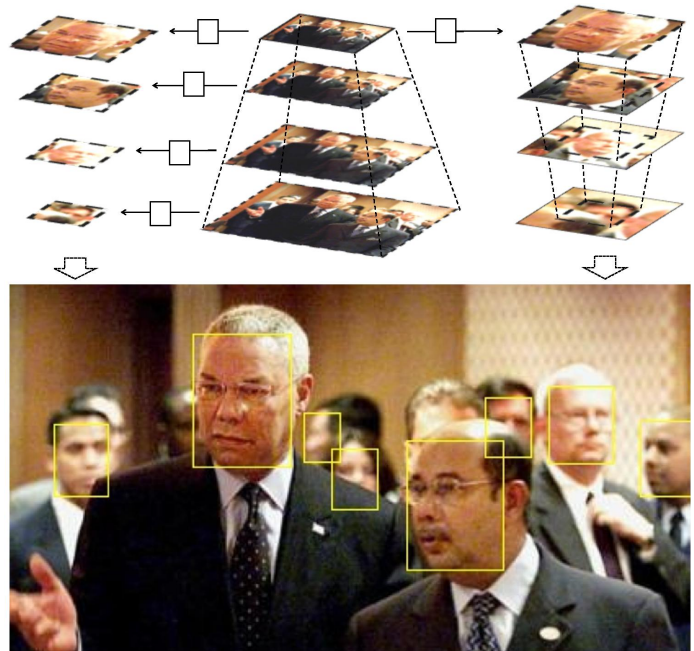


Fig. 1. A comparison between the traditional single-scale detector (left) and the anchor-based detector (right). To detect faces of different scales, a single-scale detector needs to scan an image pyramid slice by slice, while an anchor-based detector only needs to scan the smallest slice of the same image pyramid. By using the anchors, the computational cost are greatly reduced. For example, given the scale factor  $\alpha = 0.7937$  for the image pyramid and a set of 4 anchors, the computational cost of anchor-based detector will reduce to approximate 1/10 comparing with the single-scale detector.

detect faces of different scales using a dense image pyramid and the computations on different slices of the dense image pyramid are not shared. As a result, CNN-based cascade face detectors use only tiny CNNs (e.g., two or three convolutional layers) in early stages, significantly degrading their performances in real-world scenarios [15], [16], [2].

Anchor-based detectors such as Faster R-CNN [17] and single shot multibox detector (SSD) [18] have become the most popular frameworks for generic object detection. With the help of multi-scale anchors, anchor-based detectors can detect objects of different scales simultaneously, and thus work without an image pyramid. We demonstrate a comparison between the traditional single-scale detector and the anchor-based detector in Fig. 1. Anchor-based detectors usually use anchors to cover an extremely wide range of object scales, e.g., from 16x16 pixels to 512x512 pixels. As a result, very large CNNs, e.g., VGGNet [19] and ResNet [20], are usually indispensable for anchor-based detectors to learn scale-invariant features. Furthermore, large CNNs usually need to be pre-trained on a large-scale image classification dataset such as ImageNet [1], which takes intensive computational load and

B. Yu and D. Tao are with the UBTECH Sydney AI Center, The University of Sydney, Darlington, NSW 2008, Australia, and also with the Faculty of Engineering and Information Technologies, School of Information Technologies, The University of Sydney, Darlington, NSW 2008, Australia (e-mail: bayu0826@uni.sydney.edu.au; dacheng.tao@sydney.edu.au).

Manuscript received..

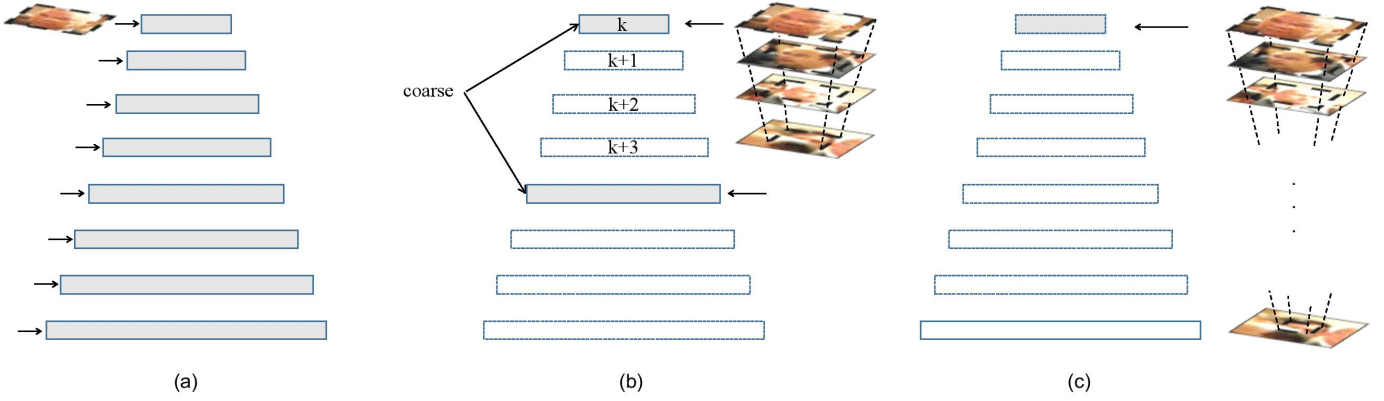


Fig. 2. An example of the trade-off between the multi-scale anchors and the dense image pyramid for dealing with large variations in object scale. (a) A single-scale detector with a dense image pyramid, i.e., traditional cascade face detectors rely on a dense image pyramid to detect faces of different scales. (b) a set of anchors with a coarse image pyramid, which is a computationally efficient. (c) A full set of anchors without an image pyramid, i.e., anchor-based detectors rely on multi-scale anchors to detect faces of all scales.

makes it difficult for us to explore new network architectures for detection. As a special case of generic object detection, face detection has been dominated by anchor-based detectors on detection accuracy [21], [22], [23], [24]. However, anchor-based detectors rely heavily on very large base networks, which have become the bottleneck for efficient face detection in both training and inference phases.

For efficient face detection, it is important to avoid both a dense image pyramid and very large base networks. Two observations raise our concerns: (1) anchor-based detectors use very large CNNs to handle extremely large variations in object scale; and (2) the computational cost for scanning a dense image pyramid has a long-tail property, i.e., the majority of computational cost comes from the largest slice of an image pyramid. Inspired by this, we devise an anchor-based cascade framework, or anchor cascade, by introducing the anchors to cascade framework. More specifically, we try to embed a set of single-position anchors to a coarse image pyramid: (1) we use a set of anchors to handle a small range of object scales, e.g., from 12x12 pixels to 24x24 pixels, and thus avoid very large CNNs; (2) we use a coarse image pyramid to handle extremely large variations in object scale, and thus avoid the dense image pyramid. See a trade-off between a set of anchors and a dense image pyramid in Fig. 2. As a result, anchor cascade significantly reduces the computational cost for dealing with large variations in object scale.

Although anchor cascade is computationally efficient, it still fails to recall some difficult faces such as tiny faces and those in profile. However, this can be addressed by exploring the contextual information, i.e., the areas surrounding the face region [25], [22], [26]. Intuitively, for a fixed detection window, e.g., 24x24 pixels, additional context will always squeeze the face region. On the one hand, additional context such as face contours and ears are crucial for detecting low-resolution faces or faces in profile; On the other hand, face parts such as the eyes, nose, and mouth, help detect partially occluded faces [27]. As a result, there will be a trade-off between adding additional context and preserving face region. Since the optimal context ratio for each candidate window

is not available, we propose a maxout structure based on a diverse set of context templates. More specifically, for each candidate window, we use multiple different context templates, i.e., the context region increases with a fixed scale factor, and a candidate window is rejected only if all context templates fail to recall it. Considering that all context templates forms in a pyramid, we refer to this maxout structure as context pyramid maxout (CPM). As a result, we greatly improve the detection accuracy, especially on difficult faces. Specifically, the context pyramid maxout structure can be efficiently implemented using a parallel-style pipeline (see details in Section III-C and IV).

In summary, we have devised anchor cascade as well as a context pyramid maxout mechanism for efficient face detection: (1) to largely reduce the computational cost and (2) to significantly boost the detection accuracy. The anchor cascade face detector enables to train small models with a high detection accuracy: (1) it outperforms typical CNN-based cascade face detectors, e.g., MTCNN [2], with a large margin, while it still runs in comparable speed; and (2) it achieves comparable detection accuracy comparing with typical anchor-based face detectors, e.g., HR [22], while it uses less than 1/25 parameters and 1/10 inference time. We conduct a number of experiments on two challenging face detection benchmarks, FDDB [28] and WIDER FACE [10], to demonstrate the effectiveness of the proposed anchor cascade framework.

## II. RELATED WORK

**Cascade face detector.** The cascade face detector [9], which uses Haar-like features and AdaBoost to train a cascade of binary classifiers, was the seminal work for real-time detection of near-frontal faces. The cascade face detector has since been improved by using: (1) robust hand-crafted features, e.g., HOG [29], MB-LBP [12], SURF [30], ACF [31], CCF [32], NPD [33], and LDF [34]; (2) asymmetric feature selection approaches, e.g., FFS [11] and SAFS [35]; (3) boosting algorithms, e.g., MILBoost [36], RealBoost [37], and LACBoost [38]; (4) new cascade structures, e.g., boosting chain [39], Closed-Loop [40], and soft-cascade [41]; and (5) multi-task learning, e.g., face/non-face classification, bounding

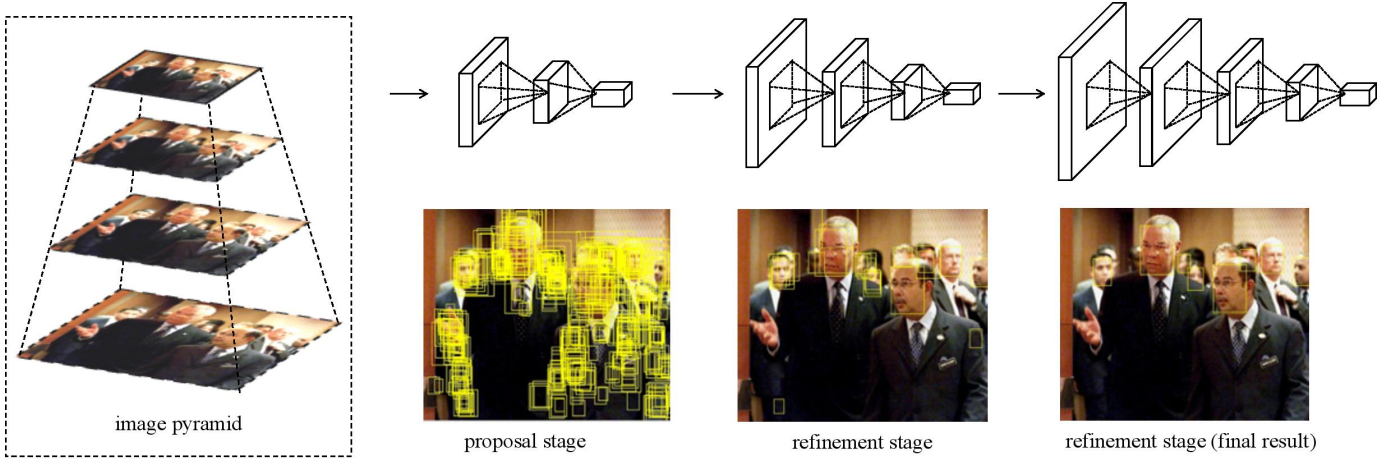


Fig. 3. The main framework of CNN-based cascade face detector. It usually contains one proposal stage and several refinement stages. More specifically, the proposal stage usually is a FCN-based binary classifier, which generates a detection score map for each slice of the dense image pyramid. The refinement stages use larger CNNs to remove hard negative examples.

box regression [42], facial landmark localization [2], and face pose estimation [43], [44].

**CNN-based cascade face detector.** CNNs have recently contributed significantly to the progress in image classification [13]. Inspired by this, CNNs have been introduced to cascade face detection [15] and improved by jointly training all stages [16]. Recently, CNN-based cascade face detectors have been further improved by using multi-task learning, i.e., jointly learning face/non-face classification, bounding box regression as well as facial key point detection [2].

**Anchor-based face detector.** Anchor-based detectors, which use multi-scale anchors at each sliding window position to simultaneously predict multiple different candidate regions, were first proposed in Faster R-CNN [17]. After that, SSD [18] tried to assign different anchors to feature maps with different receptive fields. These anchor-based detectors have shown impressive performance in face detection [21]. Recently, anchor-based face detectors have been further improved by addressing (1) the mismatch between the receptive fields and anchor sizes [24], (2) the positions of facial landmarks [45], [46], and the scale distribution histogram of faces [47], [48].

**Contextual information.** Contextual information has turned out to be crucial for object detection [49], [50] and recognition [51], [52], [53], [54]. Recently, [22] analyzed the relationship between templates with different context regions. To explore contextual information in anchor-based face detectors, [25] combined the features from different region of interests (RoIs), while [23] designed a context module using multiple different convolutional filters, e.g., 5x5 and 7x7 filters. Specifically, [55] cascaded the predictions from different RoIs to efficiently boost the performance. To explore contextual information for cascade face detectors, [26] used an additional network to extract features with body information.

### III. MAIN FRAMEWORK

In this section, we first introduce the pipeline of CNN-based cascade face detector. We then describe the long-tail problem

on computational cost and the anchor-based proposal network (APN) used in anchor cascade. After that, we introduce the context pyramid maxout structure, as well as its relationship with the multi-scale anchors. Finally, we briefly discuss the refinement stages, i.e., context-aware refinement networks.

#### A. CNN-based Cascade Face Detector

Cascade face detectors rely on a dense image pyramid to detect faces of different scales, i.e., it slides a detection window and detects single-scale faces on each slice of the dense image pyramid. Inspired by [56], the fully convolutional neural networks [57] (FCN)-based sliding-window approach has been widely used in CNN-based cascade face detectors [16], [2], [26]. More specifically, by replacing last fully connected layer with a 1x1 convolutional layer, the FCN-based sliding window approach shares the computation of the overlapped area between different windows.

**Image pyramid.** Given an image, a dense image pyramid is built by resizing the original image according to the size of detection window. For example, given the size of detection window  $S_w$ , if we want to detect faces of scale  $S$ , the original image then need to be resized by  $S_w/S$ . In practice, we usually construct the image pyramid using a fixed scale factor  $\alpha_I$ , e.g., 0.7937. That is, if the first slice of image pyramid is the original image, the  $k$ -th slice then is obtained by resizing the original image to  $\alpha_I^{k-1}$ .

**Proposal generation.** During test, by replacing the last fully connected layer with a 1x1 convolutional layer, the proposal network can deal with arbitrary input size and a score map will be generated for the input. After that, all candidate windows can be decoded from the score map and a majority of non-face windows are removed by a score threshold as well as the non-maximum suppression (NMS).

**Refinement.** In candidate windows generated by the proposal stage, a number of non-face windows still exist. As a result, refinement stages usually train larger networks by mining hard examples from the proposal stage. To remove



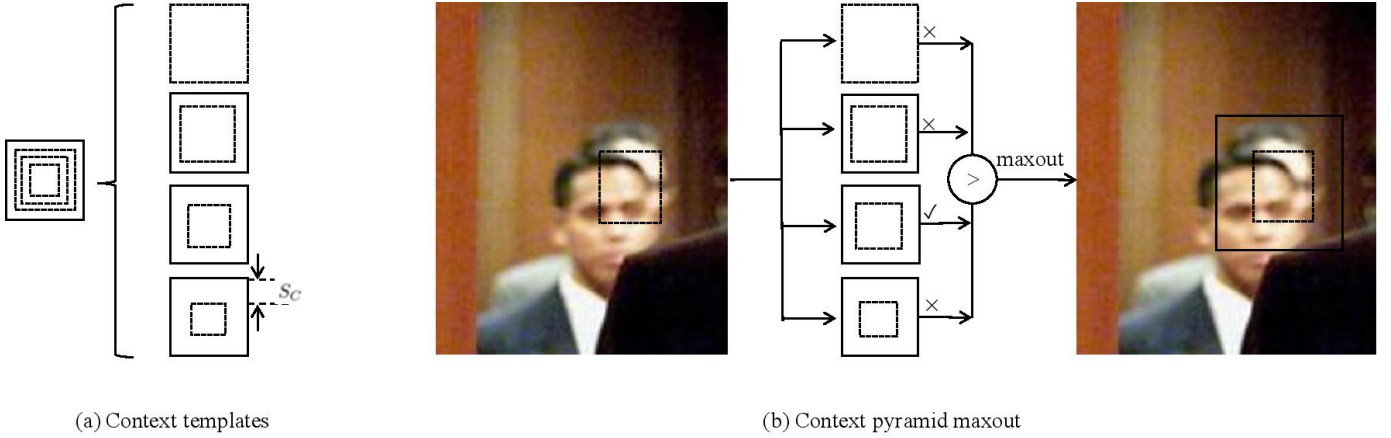


Fig. 6. An example for context pyramid maxout (CPM). In Fig. (a), it is a set of four context templates with different context regions. In Fig. (b), for each candidate window (dashline area), we use multiple context templates and only the template with maximum prediction score are kept for further processing.

long-tail problem, the computational cost will further reduce to less than  $1/n_A$ , i.e.,

$$\frac{C_{n_A}}{\sum_{k=1}^{n_A} C_k} = \frac{\alpha^{n_A}}{\sum_{k=1}^{n_A} \alpha^k} = \frac{\alpha^{n_A-1}(1-\alpha)}{1-\alpha^{n_A+1}}, \quad (2)$$

where  $\alpha = \alpha_A^2$  and  $C_k$  denotes the computational cost on the  $k$ -th slice of an image pyramid. As a result, for  $\alpha_A = 0.7937$  and  $n_A = 4$ , the computational cost will reduce to approximately  $1/9.74$ .

Anchors used in APN are inspired by the region proposal network (RPN) used in Faster R-CNN [17], the aim being to share the computations of the dense image pyramid. However, there are several differences between APN and RPN. For simplicity, we refer to these anchors as APN anchors and RPN anchors, respectively. The differences then can be described as follows. First, APN uses a small base network, which does not need to be pre-trained on ImageNet; Second, RPN allocates anchors on a  $W \times H$  feature map, while APN contains only a set of single-position anchors and is trained using image patches as a simple classification network; Third, APN anchors only cover a small range of object scales and thus do not suffer from the problem of a mismatch between anchors and the receptive fields [58], [24].

### C. Context Pyramid Maxout

Contextual information has turned out to be crucial for face detection, especially for detecting challenging faces such as tiny faces [25], [22], [26]. However, for CNN-based cascade face detectors, there is always a trade-off between adding additional context and preserving face regions due to the fixed detection window. See Fig. 6 (a) for a set of context templates with different context regions. For each candidate window, the optimal context region varies due to the large variations in illumination, occlusion, and pose. That is, we do not have prior information on context template selection. To improve the recall rate by using proper context templates, we use a diverse set of context templates with a maxout structure. More specifically, given a candidate window  $x$ , let  $p_i(x)$  denote the

score given by the  $i$ -th context template, the maxout score then is evaluated as follows:

$$p(x) = \max_{i=1,2,\dots,n_C} p_i(x), \quad (3)$$

where  $n_C$  is the number of context templates. By using the maxout score, we will only reject a candidate window if all context templates fail to recall it. A diverse set of context templates can be constructed as follows. Let  $S_C(i)$  denote the padding size of the  $i$ -th context template, we then have

$$S_C(i) = \frac{S_W * (1 - \alpha_C^{i-1})}{2}, \quad (4)$$

where  $\alpha_C \in (0, 1]$  is a scale factor. Furthermore, a set of context templates can be efficiently implemented using the same structure with a set of anchors. See more details about the implementation of context pyramid maxout in Section IV.

### D. Context-Aware Refinement

The motivation of each refinement stage is to reject as many non-face windows as possible and simultaneously retain any face windows. The refinement stage thus is important to boost the performance of a face detector on hard faces. More specifically, our refinement networks share the same structure with APN for training except that (1) refinement networks have more layers and larger input size, e.g.,  $48 \times 48$  pixels for RNet48 and  $96 \times 96$  pixels for RNet96; (2) refinement networks are trained by mining the hard examples from the proposal stage.

During testing, we perform context-aware refinement as follows. Given each candidate window retained by the proposal stage, we also obtain a context template in maxout structure, i.e., the “optimal” template for candidate window  $x$  is

$$i^* \in \arg \max_{i=1,2,\dots,n_C} p_i(x). \quad (5)$$

After that, we use the  $i^*$ -th context template of the refinement network to classify the candidate window  $x$ . Considering all stages in our cascade framework are based on anchors, we refer to the proposed framework as anchor-based cascade or anchor cascade.

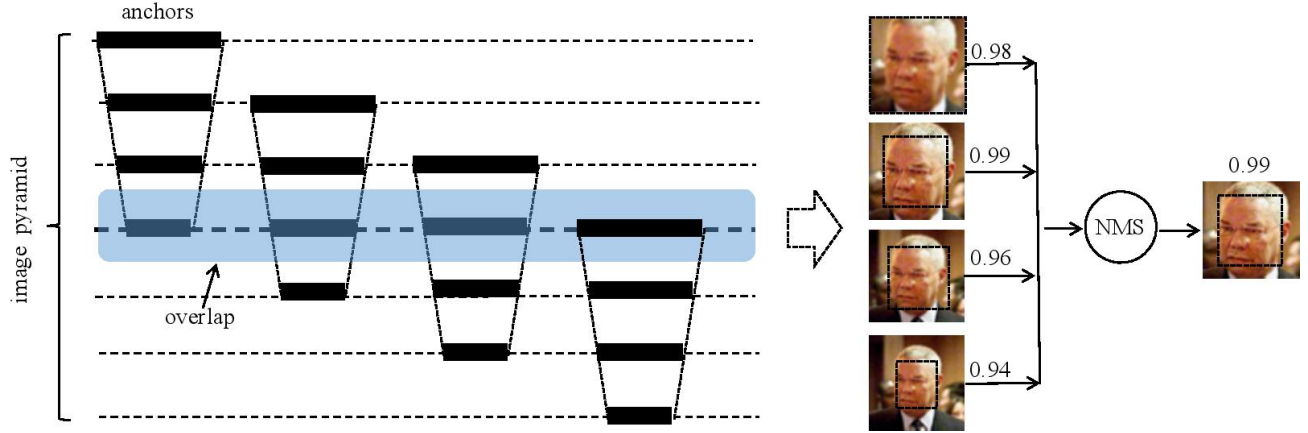


Fig. 7. An efficient implementation of context pyramid maxout. Given an APN with  $n_A$  anchors, a context pyramid on the  $k$ -th slice of the image pyramid is constructed as follows: the  $i$ -th context templates is the  $i$ -th anchor of APN, i.e., APN actually works on the  $(k - i + 1)$ -th slice of the image pyramid. As a result, a set of context templates exist on the  $k$ -th slice of the image pyramid constructed by APN working on the  $k$ -th,  $(k - 1)$ -th,  $\dots$ , slices of the image pyramid; that is, the overlap between APN anchors forms a set of context templates in a parallel-style pipeline. The maxout structure can be trivially implemented by using NMS with a large threshold, e.g., 0.9.

#### IV. IMPLEMENTATION DETAILS

In this section, we first describe the implementation details of context pyramid maxout. We then introduce the multi-task loss as well as the network architectures used in each stage of anchor cascade face detector. Finally, we discuss data preparation as well as hyper parameters for training.

##### A. Implementation of CPM

To share the computations between different context templates, we find that a set of anchors share the same input in APN, and each anchor thus contains different context information. Inspired by this, our context templates share the same scale factor with APN, i.e.,  $\alpha_C = \alpha_A$ , and a set of anchors thus form a diverse set of context templates (see Fig. 6 (a) for an example). More specifically, the context templates induced by the multi-scale anchors run in a parallel pipeline style. We demonstrate the parallel pipeline of context pyramid maxout in Fig. 7. As a result, the maximum number of context templates  $n_C$  is limited by the number of anchors in APN, i.e.,

$$n_C = n_A - i + 1. \quad (6)$$

A shortcoming for such an implementation is that more context templates always degrade the computational efficiency of APN. That is, there will be a trade-off between using more context templates and keeping computational efficiency.

##### B. Multi-task Loss

Face/non-face classification and bounding box regression are two important tasks in object detection. It has turned out that joint learning classification and regression tasks improves object detectors [17], [2], [59], [60]. Inspired by this, we train our anchor cascade in a similar multi-task learning framework [2]. For face/non-face classification, we use the softmax loss, i.e., the softmax activation function with cross entropy loss. For bounding box regression, we adopt the Euclidean loss based on the parameterizations of four coordinates as

follows. Let  $x_a^1, y_a^1, x_a^2, y_a^2$  denote the coordinates of the top-left corner and bottom-right corner of the candidate window, and let  $x_g^1, y_g^1, x_g^2, y_g^2$  denote the ground truth. The regression parameters then are defined as follows:

$$\delta_x^1 = (x_g^1 - x_a^1)/w, \quad (7)$$

$$\delta_y^1 = (y_g^1 - y_a^1)/h, \quad (8)$$

$$\delta_x^2 = (x_g^2 - x_a^2)/w, \quad (9)$$

$$\delta_y^2 = (y_g^2 - y_a^2)/h, \quad (10)$$

where  $w$  and  $h$  are the width and height of the candidate window, respectively. Specifically, all context templates share the same regression parameters, i.e., the regression parameters are evaluated without considering the padding size. The final loss function uses a loss weight  $\lambda$  to balance the face/non-face classification loss  $L_{cls}$  and bounding box regression loss  $L_{reg}$ , i.e.,

$$L_{total} = L_{cls} + \lambda * L_{reg}. \quad (11)$$

##### C. Network Architectures

For simplicity, we use plain CNNs, i.e., stacked convolutional layers, in both the proposal and refinement stages of the anchor cascade face detector. Unlike [15], [2], [16], we use convolutional layers with stride 2 to replace all max-pooling layers. As a result, all networks used in anchor cascade face detector contain only 3x3 convolutional layers followed by a PReLU [61] activation function, except the last two layers (see a typical configuration of network architectures in Figure 8).

##### D. Data Preparation and Training

We collect training and validation data from the training set of WIDER FACE [10]. Inspired by [15], [2], we randomly generate positive (IoU > 0.65) and negative (IoU < 0.3) samples by keeping a ratio pos:neg=1:3. For bounding box regression, we randomly collect semi-positive (IoU > 0.4) samples with the same ratio as positive samples. For the

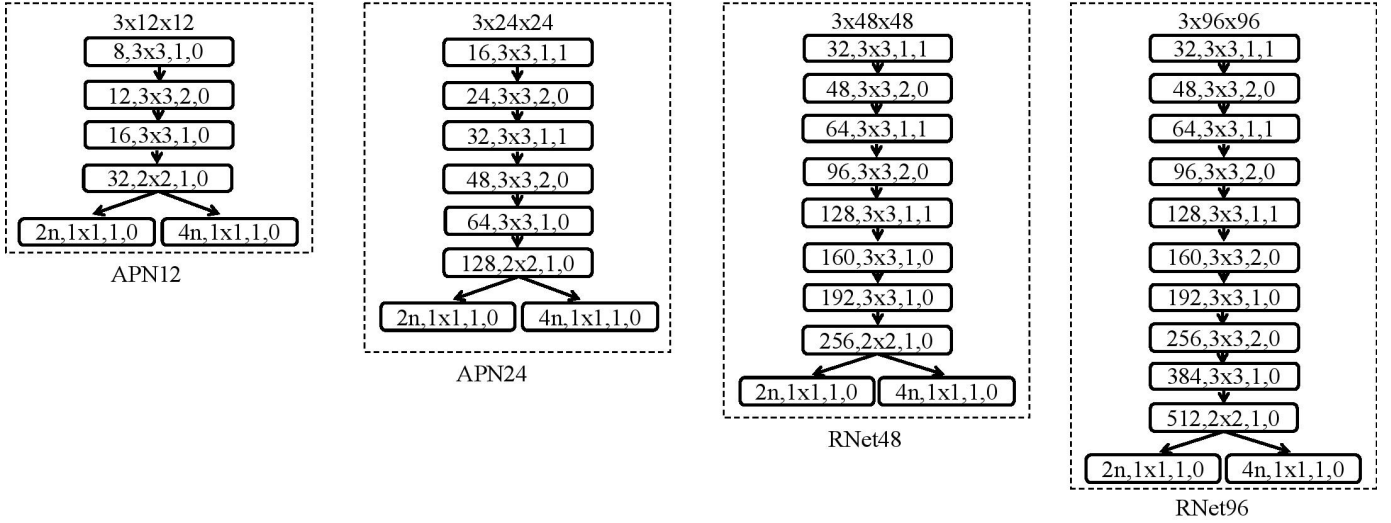


Fig. 8. A typical configuration of network architectures used in anchor cascade. The text in convolutional layers refers to the number of channels, size of kernel, stride, and padding size, respectively. We use four anchors  $n_A = 4$  in our current implementation, i.e.,  $n = 4$ .

proposal stage, we use 0.8M positive samples. For the two refinement stages, we use 1M and 1.1M positive samples, respectively. More specifically, in the first refinement stage, we use 3.0M negative samples (1.0M from the proposal stage; 2.0M hard negative samples). In the second refinement stage, we use 3.3M negative samples (1.1M from the first refinement stage; 2.2M hard negative samples).

The proposed anchor cascade is implemented using a customized Caffe [62]. We use  $\lambda = 1.0$  for the second refinement stage and  $\lambda = 0.5$  for other stages. All models are trained using stochastic gradient descent (SGD) with momentum 0.9 and weight decay  $2e-5$ . For the proposal stage, we start with a learning rate 0.1 for 4 epochs and then divide by 5 for every 2 epochs, with 12 epochs in total. For both refinement stages, we use a learning rate of 0.01 for 5 epochs and 0.001 for 3 epochs. We use the batch size 480, 360 and 240 for APN24, RNet48, and RNet96, respectively. Specifically, we use a dropout of 0.5 for the second refinement stage, i.e., RNet96.

## V. EXPERIMENT

In this section, we first briefly introduce two widely used benchmarks in face detection, FDDB [28] and WIDER FACE [10]. We then conduct a number of experiments to demonstrate the effectiveness of the proposed anchor cascade framework.

### A. Datasets and Benchmarks

**FDDB** contains 2,845 images with 5,171 faces. For FDDB evaluations, we follow the “unrestricted training” setting, i.e., we use all ten folds of the dataset as validation data without performing 10-fold cross-validation.

**WIDER FACE** contains 32,203 images with 393,793 faces, of which 40% are used for training, 10% for validation, and 50% for testing. According to the detection rate, the validation data are divided into three classes: “easy”, “medium”, and “hard”. We train all our models using the training set and test on the validation set.

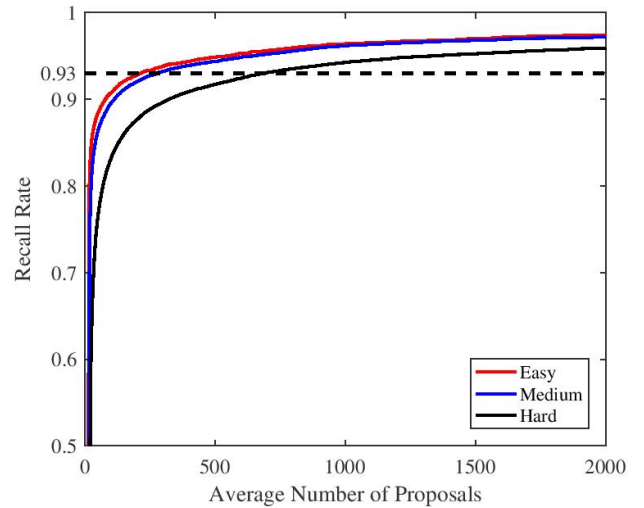


Fig. 9. The recall rate of APN24 on WIDER FACE. Specifically, we use min-face  $S_{min} = 10$  and  $n_C = 3$  in this experiment.

### B. Effectiveness for Proposal Generation

Proposal generation is very important for efficient face detection, especially for cascade face detectors. On the one hand, refinement stages are based on generated proposals and they will never recall any faces outside generated proposals; On the other hand, refinement stages usually contain larger networks for hard examples and thus are not efficient for thousands of proposals. As a result, it is crucial for efficient face detection to generate high quality proposals.

In this experiment, we train several APNs using input size  $12 \times 12$  and  $24 \times 24$  pixels, i.e., APN12 and APN24, respectively. To demonstrate effectiveness of APNs for proposal generation, we compare the recall rate of APNs under different number of proposals. More specifically, in Table I, we evaluate the recall rate on FDDB using 50 and 100 proposals per



Fig. 10. Qualitative results on Fddb and best view in color.

TABLE I

THE RECALL RATE OF APN24 ON Fddb. WE USE MIN-FACE  $S_{min} = 16, 20, 24$ , AS Fddb DATASET CONTAINS FACE AROUND  $20 \times 20$  PIXELS. SPECIFICALLY, APN24(S) SHARES THE SAME STRUCTURE WITH APN24, EXCEPT THE NUMBER OF CHANNELS IN EACH CONVOLUTIONAL LAYER, I.E., 8, 16, 24, 32, 64, 128, RESPECTIVELY.

Network	$n_C$	$S_{min}$	R@50	R@100	Model	Speed	Network	$n_C$	$S_{min}$	R@50	R@100	Model	Speed
PNet [2]	1	16	0.9604	0.9712	28KB	25.17 FPS	APN24	3	16	0.9920	0.9948	360KB	15.73 FPS
PNet [2]	1	20	0.9557	0.9700	28KB	33.75 FPS	APN24	3	20	0.9898	0.9919	360KB	23.16 FPS
PNet [2]	1	24	0.9540	0.9673	28KB	43.15 FPS	APN24	3	24	0.9851	0.9876	360KB	30.13 FPS
APN24	1	16	0.9872	0.9910	360KB	22.91 FPS	APN12	2	16	0.9598	0.9685	69KB	43.31 FPS
APN24	1	20	0.9860	0.9882	360KB	32.63 FPS	APN12	2	20	0.9484	0.9553	69KB	55.78 FPS
APN24	1	24	0.9825	0.9856	360KB	43.39 FPS	APN12	2	24	0.9354	0.9408	69KB	72.11 FPS
APN24	2	16	0.9907	0.9946	360KB	20.82 FPS	APN24(S)	2	16	0.9760	0.9885	164KB	38.33 FPS
APN24	2	20	0.9888	0.9921	360KB	29.68 FPS	APN24(S)	2	20	0.9604	0.9743	164KB	51.41 FPS
APN24	2	24	0.9859	0.9901	360KB	38.02 FPS	APN24(S)	2	24	0.9564	0.9624	164KB	65.81 FPS

image (in average). We find that (1) APN24 achieves very high recall rate comparing with PNet used in MTCNN, e.g., 0.9825 vs 0.9540; (2) APN24 runs in comparable speed with PNet, even though it is 12 times larger than PNet. Furthermore, with an average of 250 proposals per image, APN24 recalls almost all faces in Fddb, i.e., more than 99.9% faces. To further demonstrate the efficiency of APNs, we train two small models, i.e., APN12 and APN24(S). We find that both APN12 and APN24(S) run very fast, and APN24(S) outperforms PNet with a clear margin. A possible reason for the weakness

of APN12 is that the smallest anchor in APN12, i.e.,  $6 \times 6$  pixels, is too small to carry useful information for challenge faces. In Figure 9, we demonstrate the recall rate on WIDER FACE, which contains a large number of tiny faces. We see that APN24 recalls more than 93% faces with less than 1k proposals per image on “hard” set of WIDER FACE. By comparison, the default face detector used in WIDER FACE, i.e., Faceness-Net [63], achieves a recall rate 93.1% with 10k proposals.





Fig. 11. Qualitative results on WIDER FACE. Please zoom in to see tiny detections.

### C. Two-stage Anchor Cascade

APNs achieve very high recall rate, while there are still a large number of false positives in generated proposals. To further remove hard false positives, we use a refinement network RNet48 and demonstrate its effectiveness for efficient face detection by comparing with (1) a popular CNN-based cascade face detector MTCNN [2]; and (2) a typical anchor-based face detector HR [22]. In Table II, we see that (1) the two-stage anchor cascade face detector is comparable with MTCNN in both model size and running speed, while it greatly outperforms MTCNN in detection accuracy; (2) it achieves comparable performance with HR on Fddb using less than 1/25 parameters and 1/10 inference time. As a result, with

TABLE II  
THE PERFORMANCE OF THE TWO-STAGE ANCHOR CASCADE ON Fddb. SPECIFICALLY, WE USE MIN-FACE  $S_{min} = 18$  FOR BOTH APN24(S) AND APN24.

Method	$n_C$	FP=1k	Model	Speed
MTCNN [2]	-	0.9435	2.10MB	25.09 FPS
HR-ResNet101 [22]	-	0.9698	103.7MB	2.12 FPS
APN24(S)+RNet48	1	0.9416	3.56MB	40.79 FPS
APN24+RNet48	1	0.9704	3.76MB	24.74 FPS

anchor cascade, we further bridge the gap between CNN-based cascade face detectors and anchor-based face detectors.

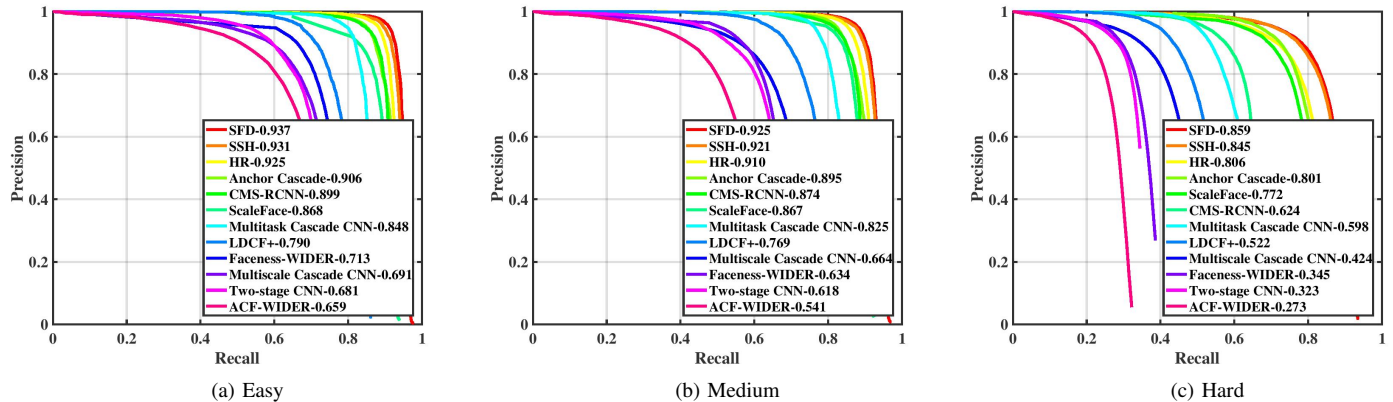


Fig. 12. The performance of three-stage anchor cascade on WIDER FACE.

TABLE III  
THE PERFORMANCE OF ANCHOR CASCADE ON WIDER FACE. WE USE  
MIN-FACE  $S_{min} = 8$ .

Method	$n_C$	Easy	Medium	Hard
MTCNN [2]	-	0.848	0.825	0.598
HR-ResNet101 [22]	-	0.925	0.910	0.806
APN24+RNet48	3	0.883	0.879	0.761
APN24+RNet48+RNet96	3	0.906	0.895	0.801

#### D. Three-stage Anchor Cascade

Although the two-stage anchor cascade face detector achieves a good trade-off between detection accuracy and inference speed, we find that most of false positives come from detecting tiny faces. To further boost the performance, we use one more refinement stage RNet96 as the three-stage anchor cascade face detector and demonstrate its performance in Table III. We show that the three-stage anchor cascade face detector achieves comparable detection accuracy with HR [22] on hard set of WIDER FACE. In Fig. 13, we also find that it achieves a detection rate 98.37% at 1k false positives on FDDB, which outperforms both SSH [23] (98.1%) and SFD<sup>1</sup> [24] (98.26%). However, for detecting tiny faces, there is still a clear margin between anchor cascade face detectors and the state-of-art anchor-based detectors such as SSH [23] and SFD [24], as shown in Fig. 12. A possible reason is that the anchor cascade face detectors use limited contextual information for efficient face detection, while it is more effective for tiny faces to explore more contextual information, such as body part [25], [26].

## VI. CONCLUSION

In this paper, we propose an anchor cascade framework for efficient face detection by exploring the multi-scale anchors in CNN-based cascade face detectors. To further improve the recall rate, we devise a context pyramid maxout mechanism in harmony with the anchor cascade framework. By using anchor cascade face detector, we further bridge the gap between

<sup>1</sup>They add 238 unlabelled faces as ground truth.

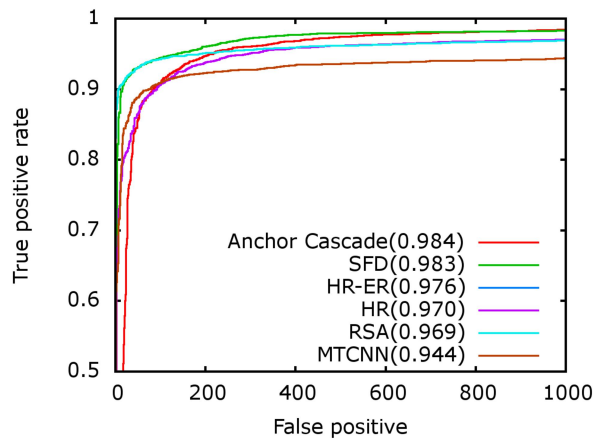


Fig. 13. The performance of three-stage anchor cascade on FDDB.

anchor-based face detectors and CNN-based cascade face detectors. Specifically, our anchor cascade face detector is comparable with typical CNN-based cascade face detectors, e.g., MTCNN, in both model size and running speed, while the detection accuracy has been greatly improved, e.g., from 0.9435 to 0.9704 at 1k false positives on FDDB. Experimental results on FDDB and WIDER FACE demonstrate the effectiveness of the proposed anchor cascade framework for efficient face detection.

## ACKNOWLEDGMENT

The authors would like to thank...

## REFERENCES

- [1] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, “Imagenet large scale visual recognition challenge,” *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [2] 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.
- [3] Y. Sun, X. Wang, and X. Tang, “Deep convolutional network cascade for facial point detection,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2013.

- [4] J. Thies, M. Zollhöfer, M. Nießner, L. Valgaerts, M. Stamminger, and C. Theobalt, "Real-time expression transfer for facial reenactment." *ACM Trans. Graph.*, vol. 34, no. 6, pp. 183–1, 2015.
- [5] J. Thies, M. Zollhofer, M. Stamminger, C. Theobalt, and M. Nießner, "Face2face: Real-time face capture and reenactment of rgb videos," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 2387–2395.
- [6] Y. Sun, X. Wang, and X. Tang, "Deep learning face representation from predicting 10,000 classes," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 1891–1898.
- [7] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 815–823.
- [8] O. M. Parkhi, A. Vedaldi, A. Zisserman et al., "Deep face recognition." in *BMVC*, vol. 1, no. 3, 2015, p. 6.
- [9] P. Viola and M. J. Jones, "Robust real-time face detection," *International Journal of Computer Vision*, vol. 57, no. 2, pp. 137–154, 2004.
- [10] S. Yang, P. Luo, C.-C. Loy, and X. Tang, "Wider face: A face detection benchmark," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 5525–5533.
- [11] J. Wu, S. C. Brubaker, M. D. Mullin, and J. M. Rehg, "Fast asymmetric learning for cascade face detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 30, no. 3, pp. 369–382, 2008.
- [12] L. Zhang, R. Chu, S. Xiang, S. Liao, and S. Z. Li, "Face detection based on multi-block lbp representation," in *International Conference on Biometrics*. Springer, 2007, pp. 11–18.
- [13] 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.
- [14] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [15] H. Li, Z. Lin, X. Shen, J. Brandt, and G. Hua, "A convolutional neural network cascade for face detection," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 5325–5334.
- [16] H. Qin, J. Yan, X. Li, and X. Hu, "Joint training of cascaded cnn for face detection," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 3456–3465.
- [17] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems*, 2015, pp. 91–99.
- [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] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.
- [20] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [21] H. Jiang and E. Learned-Miller, "Face detection with the faster r-cnn," in *Automatic Face & Gesture Recognition (FG), 2017 12th IEEE International Conference on*. IEEE, 2017, pp. 650–657.
- [22] P. Hu and D. Ramanan, "Finding tiny faces," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [23] M. Najibi, P. Samangouei, R. Chellappa, and L. S. Davis, "Ssh: Single stage headless face detector," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [24] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Z. Li, "S3fd: Single shot scale-invariant face detector," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [25] 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.
- [26] K. Zhang, Z. Zhang, H. Wang, Z. Li, Y. Qiao, and W. Liu, "Detecting faces using inside cascaded contextual cnn," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [27] S. Yang, P. Luo, C. C. Loy, and X. Tang, "Faceness-net: Face detection through deep facial part responses," *arXiv preprint arXiv:1701.08393*, 2017.
- [28] V. Jain and E. Learned-Miller, "Fddb: A benchmark for face detection in unconstrained settings," University of Massachusetts, Amherst, Tech. Rep. UM-CS-2010-009, 2010.
- [29] Q. Zhu, M.-C. Yeh, K.-T. Cheng, and S. Avidan, "Fast human detection using a cascade of histograms of oriented gradients," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2. IEEE, 2006, pp. 1491–1498.
- [30] J. Li and Y. Zhang, "Learning surf cascade for fast and accurate object detection," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2013.
- [31] B. Yang, J. Yan, Z. Lei, and S. Z. Li, "Aggregate channel features for multi-view face detection," in *IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 2014, pp. 1–8.
- [32] —, "Convolutional channel features," in *Computer Vision (ICCV), 2015 IEEE International Conference on*. IEEE, 2015, pp. 82–90.
- [33] 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, 2016.
- [34] Y. Yuan, W. Zou, Y. Zhao, X. Wang, X. Hu, and N. Komodakis, "A robust and efficient approach to license plate detection," *IEEE Transactions on Image Processing*, vol. 26, no. 3, pp. 1102–1114, 2017.
- [35] B. Yu, M. Fang, D. Tao, and J. Yin, "Submodular asymmetric feature selection in cascade object detection," in *AAAI*, 2016, pp. 1387–1393.
- [36] C. Zhang, J. C. Platt, and P. A. Viola, "Multiple instance boosting for object detection," in *Advances in Neural Information Processing Systems*, 2006, pp. 1417–1424.
- [37] C. Huang, H. Ai, Y. Li, and S. Lao, "High-performance rotation invariant multiview face detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 29, no. 4, pp. 671–686, 2007.
- [38] C. Shen, P. Wang, S. Paisitkriangkrai, and A. van den Hengel, "Training effective node classifiers for cascade classification," *International Journal of Computer Vision*, vol. 103, no. 3, pp. 326–347, 2013.
- [39] R. Xiao, L. Zhu, and H.-J. Zhang, "Boosting chain learning for object detection," in *The IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2003, pp. 709–715.
- [40] L. Galteri, L. Seidenari, M. Bertini, and A. Del Bimbo, "Spatio-temporal closed-loop object detection," *IEEE Transactions on Image Processing*, vol. 26, no. 3, pp. 1253–1263, 2017.
- [41] L. Bourdev and J. Brandt, "Robust object detection via soft cascade," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2. IEEE, 2005, pp. 236–243.
- [42] D. Chen, S. Ren, Y. Wei, X. Cao, and J. Sun, "Joint cascade face detection and alignment," in *European Conference on Computer Vision*. Springer, 2014, pp. 109–122.
- [43] X. Zhu and D. Ramanan, "Face detection, pose estimation, and landmark localization in the wild," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2012, pp. 2879–2886.
- [44] C. Zhang and Z. Zhang, "Improving multiview face detection with multi-task deep convolutional neural networks," in *Applications of Computer Vision (WACV), 2014 IEEE Winter Conference on*. IEEE, 2014, pp. 1036–1041.
- [45] 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.
- [46] 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.
- [47] Y. Liu, H. Li, J. Yan, F. Wei, X. Wang, and X. Tang, "Recurrent scale approximation for object detection in cnn," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [48] Z. Hao, Y. Liu, H. Qin, J. Yan, X. Li, and X. Hu, "Scale-aware face detection," *arXiv preprint arXiv:1706.09876*, 2017.
- [49] S. K. Divvala, D. Hoiem, J. H. Hays, A. A. Efros, and M. Hebert, "An empirical study of context in object detection," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2009, pp. 1271–1278.
- [50] Y. Li, S. Tang, M. Lin, Y. Zhang, J. Li, and S. Yan, "Implicit negative sub-categorization and sink diversion for object detection," *IEEE Transactions on Image Processing*, vol. 27, no. 4, pp. 1561–1574, 2018.
- [51] A. Oliva and A. Torralba, "The role of context in object recognition," *Trends in Cognitive Sciences*, vol. 11, no. 12, pp. 520–527, 2007.
- [52] J. Garcia, N. Martinel, A. Gardel, I. Bravo, G. L. Foresti, and C. Micheloni, "Discriminant context information analysis for post-ranking person re-identification," *IEEE Transactions on Image Processing*, vol. 26, no. 4, pp. 1650–1665, 2017.
- [53] X. Song, S. Jiang, and L. Herranz, "Multi-scale multi-feature context modeling for scene recognition in the semantic manifold," *IEEE Transactions on Image Processing*, vol. 26, no. 6, pp. 2721–2735, 2017.
- [54] H. Lee and H. Kwon, "Going deeper with contextual cnn for hyper-spectral image classification," *IEEE Transactions on Image Processing*, vol. 26, no. 10, pp. 4843–4855, 2017.
- [55] W. Ouyang, K. Wang, X. Zhu, and X. Wang, "Chained cascade network for object detection," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.

- [56] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. Lecun, "Overfeat: Integrated recognition, localization and detection using convolutional networks," in *International Conference on Learning Representations (ICLR), CBLIS, April 2014*, 2014.
- [57] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 3431–3440.
- [58] Z. Cai, Q. Fan, R. S. Feris, and N. Vasconcelos, "A unified multi-scale deep convolutional neural network for fast object detection," in *European Conference on Computer Vision*. Springer, 2016, pp. 354–370.
- [59] J. Yu, Y. Jiang, Z. Wang, Z. Cao, and T. Huang, "Unitbox: An advanced object detection network," in *Proceedings of the 2016 ACM on Multimedia Conference*. ACM, 2016, pp. 516–520.
- [60] 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," *arXiv preprint arXiv:1603.01249*, 2016.
- [61] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1026–1034.
- [62] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, 2014, pp. 675–678.
- [63] S. Yang, P. Luo, C.-C. Loy, and X. Tang, "From facial parts responses to face detection: A deep learning approach," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 3676–3684.



**Baosheng Yu** Biography text here.



**Dacheng Tao** Biography text here.