

RefineFace: Refinement Neural Network for High Performance Face Detection

Shifeng Zhang*, Cheng Chi*, Zhen Lei†, *Senior Member, IEEE*, and Stan Z. Li, *Fellow, IEEE*

Abstract—Face detection has achieved significant progress in recent years. However, high performance face detection still remains a very challenging problem, especially when there exists many tiny faces. In this paper, we present a single-shot refinement face detector namely RefineFace to achieve high performance. Specifically, it consists of five 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 regressor. To improve the classification ability for high recall efficiency, STC first filters out most simple negatives from low level detection layers to reduce search space for subsequent classifier, then SML is applied to better distinguish faces from background at various scales and FSM is introduced to let the backbone learn more discriminative features for classification. Besides, RFE is presented to provide more diverse receptive field to better capture faces in some extreme poses. Extensive experiments conducted on WIDER FACE, AFW, PASCAL Face, FDDB, MAFA demonstrate that our method achieves state-of-the-art results and runs at 37.3 FPS with ResNet-18 for VGA-resolution images.

Index Terms—Face detection, refinement network, high performance.

1 INTRODUCTION

FACE detection is a long-standing problem in computer vision with many applications, such as face alignment, face analysis, face recognition and face tracking. Given an image, the goal of face detection is to determine whether there are any faces, and if any, return the bounding box of each face. To detect faces efficiently and accurately, different detection pipelines have been designed after the pioneering work of Viola-Jones [1]. Among them, the single-shot anchor-based approach [2], [3], [4], [5], [6] is the dominant method. It performs face detection based on regular and dense anchors over various locations, scales and aspect ratios. In this framework, the face detection task is decomposed into two sub-tasks: the binary classification and the bounding box regression. The former one aims to classify the preset anchor boxes into face and background, and the latter one is to regress those detected faces to more accurate locations.

With the development of deep convolutional neural networks (CNNs), single-shot anchor-based face detectors have been thoroughly studied and great progress has been made in recent years. In particular, on the very challenging face detection dataset WIDER FACE [7], the average precision (AP) on its Hard subset has been improved from 40.0% to 90.0% by recent algorithms [2], [3], [4], [5], [6] over the past three years. For now, it has become a challenging

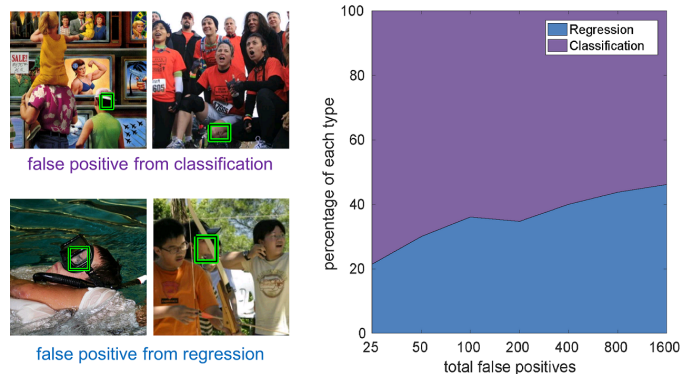


Fig. 1. Illustration of false positives of our baseline face detector on the WIDER FACE validation Hard subset. Left: Example of two error types of false positives. Right: Distribution of two error types of false positives.

problem to further improve these single-shot face detectors' performance, especially when there exists many tiny faces. In our opinion, there remains room for improvement in two aspects: a) location accuracy: accuracy of the bounding box location needs to be improved, *i.e.*, boosting the regression ability; b) recall efficiency: more faces need to be recalled with less false positives, *i.e.*, enhancing the classification ability. To embody these two aspects still can be improved, we utilize the detection analysis tool¹ to analyze the error distribution of our baseline face detector RetinaNet [8] on the WIDER FACE validation Hard subset. As shown in Figure 1, there are two error types of false positives from face detectors: the Regression (LOC) error indicates that a face is detected with a misaligned localization, and the Classification (CLS) error means that a background region

Shifeng Zhang, Zhen Lei and Stan Z. Li are with the Center for Biometrics and Security Research (CBSR), National Laboratory of Pattern Recognition (NLPR), Institute of Automation Chinese Academy of Sciences (CASIA) and University of Chinese Academy of Sciences (UCAS), Beijing, China (e-mail: {shifeng.zhang, zlei, szli}@nlpr.ia.ac.cn).

Cheng Chi is with the Institute of Electronics Chinese Academy of Sciences (IECAS) and University of Chinese Academy of Sciences (UCAS), Beijing, China (e-mail: chicheng15@mails.ucas.ac.cn).

*Equal contribution. † Corresponding author.

1. <http://web.engr.illinois.edu/~dhoiem/projects/detectionAnalysis>

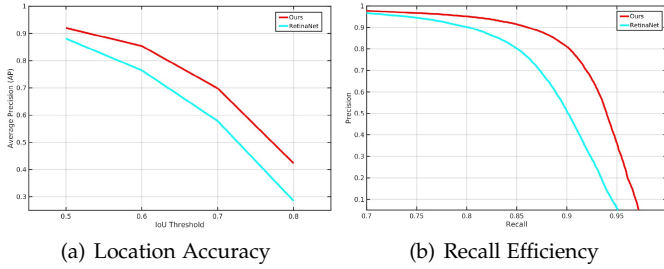


Fig. 2. (a) As the IoU threshold increases, the AP of RetinaNet drops dramatically. Our method improves its location accuracy by boosting the regression ability. (b) RetinaNet produces about 50% false positives when the recall rate is 90% and it also misses about 5% faces. Our method improves its recall efficiency by enhancing the classification ability.

is mistakenly detected as a face. Apart from false positives, false negatives indicate that a face fails to be detected, which also belong to the CLS error. These two error modes are elaborated as follows.

The LOC error is triggered by the lack of strong regression ability and its manifestation is plenty of false detections with inaccurate localization. If the regression ability of the face detector can be enhanced, these false positives from the LOC error will be reduced. Actually, the location accuracy in the face detection task has attracted much more attention of researchers in recent years. Although current evaluation criteria of most face detection datasets [9], [10], [11] do not focus on the location accuracy, the WIDER FACE Challenge² adopts MS COCO [12] evaluation criterion, which puts more emphasis on bounding box location accuracy. To visualize the location accuracy issue, we use different IoU thresholds to evaluate our baseline face detector RetinaNet [8] on the WIDER FACE dataset. As shown in Figure 2(a), as the IoU threshold increases, the AP drops dramatically, indicating that the accuracy of the bounding box location needs to be improved. To this end, Gidaris et al. [13] propose iterative regression during inference to improve the accuracy. Cascade R-CNN [14] addresses this issue by cascading R-CNN with different IoU thresholds. RefineDet [15] applies two-step regression to single-shot detector. However, blindly adding multi-step regression to the face detection task is often counterproductive, which needs more exploration.

The CLS error is caused by the non-robust classification ability and its manifestation is lots of false alarms and missing faces. Specifically, the average precision (AP) of current face detection algorithms is already very high, but the recall efficiency is not high enough. As shown in Figure 2(b) of our baseline, its precision is only about 50% (half of detections are false alarms) when the recall rate is equal to 90%, and its highest recall rate is only about 95% (the remaining 5% faces are still not detected). Reflected on the shape of the Precision-Recall curve, it has not extended far enough to the right as well as not steep enough. In our opinion, the reasons behind the recall efficiency issue in the CLS error of single-shot detectors are 1) class imbalance: plenty of small anchors need to be tiled to detect tiny faces, causing the extreme class imbalance problem; 2) scale problem: different scales of anchors have different degrees of classification

difficulty, and the classification of smaller anchors are more difficult; 3) feature misalignment: different anchors at the same location are classified based on the same misaligned features. They are the culprit leading to the CLS error. If we can improve the recall efficiency via enhancing the classification ability, more faces can be correctly detected from the complex background, making fewer CLS errors and higher average precision (AP). Therefore, it is worth further study to improve the classification ability of the face detector.

To reduce the aforementioned two errors, we propose five improvements to enhance the regression and classification ability of the high performance face detector and present a new state-of-the-art method namely RefineFace. Specifically, to improve the regression ability, we apply the STR to coarsely adjust the locations and sizes of anchors from high level detection layers to provide better initialization for the subsequent regressor. To enhance the classification ability, we first use the STC to filter out most simple negatives from low level detection layers to reduce the search space for the subsequent classifier, then employ the SML to better distinguish faces from background at various scales and the FSM to let the backbone network learn more discriminative features for classification. Besides, we design the RFE to provide more diverse receptive field to better capture faces in some extreme poses. We conduct extensive experiments on the WIDER FACE, AFW, PASCAL Face, FDDB and MAFA benchmark datasets and achieve the state-of-the-art results with 37.3 FPS for the VGA-resolution image. The main contributions of this paper can be summarized below.

- Designing a STR module to coarsely adjust the locations and sizes of anchors from high level layers to provide better initialization for the subsequent regressor.
- Presenting a STC module to filter out most simple negative samples from low level layers to reduce the classification search space.
- Introducing a SML module to better distinguish faces from background across different scales.
- Proposing a FSM module to learn more discriminative features for the classification task.
- Constructing a RFE module to provide more diverse receptive fields for detecting extreme-pose faces.
- Achieving state-of-the-art performances on AFW, PASCAL face, FDDB, MAFA and WIDER FACE datasets.

Preliminary results of this work have been published in [2]. The current work has been improved and extended from the conference version in several important aspects. (1) We introduce a Scale-aware Margin Loss (SML) function to better distinguish faces from the complex background across different scales. (2) We design a Feature Supervision Module (FSM) to learn more discriminative features for classification. (3) We noticeably improve the accuracy of the detector in our previous work without introducing any additional overhead during the inference phase. (4) All sections are rewritten with more details, more references and more experiments to have a more elaborate presentation.

2. <http://wider-challenge.org>

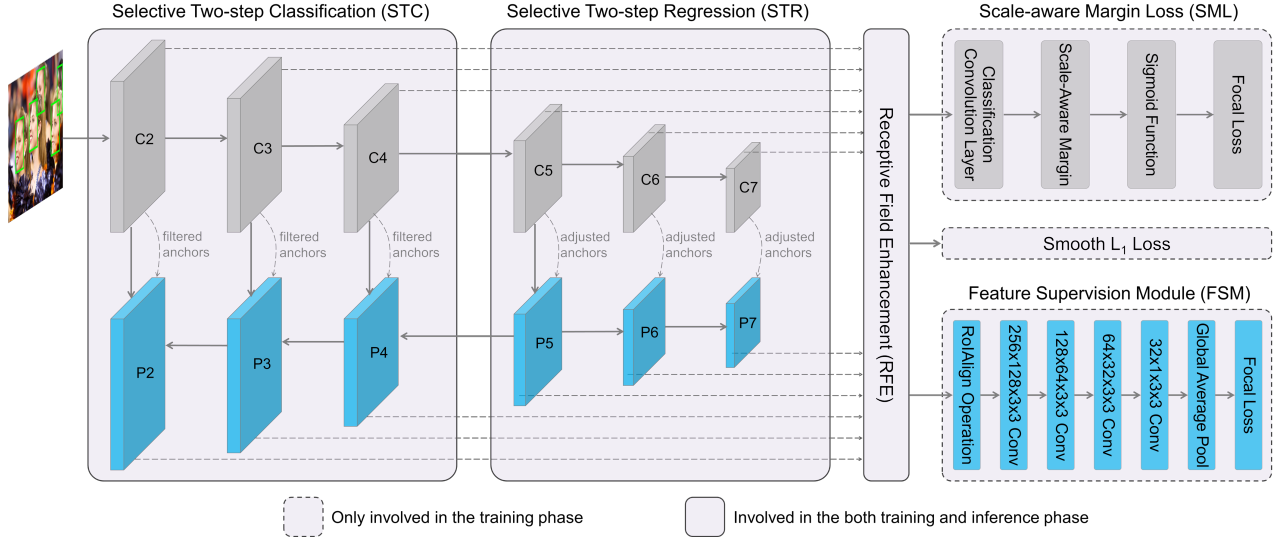


Fig. 3. Structure of RefineFace. It is based on RetinaNet with five proposed modules. Among them, SML and FAM are only involved in training without any overhead in inference, while STC, STR and RFE introduce a small amount of overhead.

2 RELATED WORK

Face detection has attracted much attention these years for its wide practical applications. The pioneering work of Viola and Jones [1] uses AdaBoost with Haar features to train a cascaded face detector and inspires several different approaches afterwards [16], [17], [18]. Besides, the Deformable Part Model (DPM) [19] is another popular framework in traditional face detection [9], [20], [21]. However, the aforementioned methods are unreliable in complex scenarios because of non-robust hand-crafted features and classifiers.

In recent years, face detection has been dominated by CNN-based methods. Li *et al.* [22] achieve promising accuracy and efficiency by separately training a series of CNN models and following work [23] realizes end-to-end optimization. Yang *et al.* [24] detect faces under severe occlusion and unconstrained pose variations via scoring facial parts responses according to their spatial structure and arrangement. Ohn-Bar *et al.* [25] utilize the boosted decision tree classifier to detect faces. Yu *et al.* [26] propose an IoU loss to directly regress the bounding box that is robust to objects of varied shapes and scales. Zhang *et al.* [27] use multi-task cascaded CNNs to jointly detect faces and landmarks. Yang *et al.* [28] apply a specialized set of CNNs with different structures to detect different scales of faces. Hu *et al.* [29] train some separate detectors for different scales to find tiny faces. Zhu *et al.* [30] introduce an EMO score to evaluate the quality of anchor setting. Hao *et al.* [31] develop a scale proposal stage to guide the zoom-in and zoom-out of image to detect normalized face. Song *et al.* [32] propose a scale estimation and spatial attention proposal module to pay attention to some specific scales and valid locations in the image pyramid. Shi *et al.* [33] detect rotated faces in a coarse-to-fine manner under a cascade-style structure. Bai *et al.* [34] utilize GAN [35] to detect blurry small faces via generating clear super-resolution ones.

In addition, generic object detection algorithms have inspired many face detection methods. CMS-RCNN [36] integrates contextual reasoning into Faster R-CNN [37] to

help reduce the overall detection errors. Face R-CNN [38], Face R-FCN [39] and FDNet [40] apply Faster R-CNN [37] and R-FCN [41] with some specific strategies to perform face detection. FaceBoxes [42] designs a CPU real-time face detector based on SSD [43]. S³FD [6] introduces some specific strategies onto SSD [43] to alleviate the matching problem of small faces. SSH [4] adds large filters on each prediction head to merge the context information. Pyramid-Box [5] takes advantage of the information around human faces to improve detection performance. FAN [44] utilizes the anchor-level attention mechanism onto RetinaNet [8] to detect the occluded faces. FANet [45] aggregates higher-level features like FPN [46] to augment lower-level features at marginal extra computation cost. DFS [47] introduces a more effective feature fusion pyramid and a more efficient segmentation branch to handle hard faces. DSFD [3] inherits the architecture of SSD [43] and introduces a feature enhance module to extend the single shot detector to dual shot detector. SRN [2] combines the multi-step detection in RefineDet [15] and the focal loss in RetinaNet [8] to perform efficient and accurate face detection. VIM-FD [48] and ISRN [49] improve SRN [2] with data augmentation, attention mechanism and training from scratch.

3 REFINEFACE

The overall architecture of RefineFace is shown in Figure 3. We adopt ResNet [50] with 6-level feature pyramid structure as backbone for RefineFace. The feature maps extracted from those four residual blocks are denoted as C2, C3, C4, and C5, respectively. C6 and C7 are extracted by two simple down-sample 3×3 convolution layers after C5. The lateral structure between the bottom-up and the top-down pathways is the same as [46]. P2, P3, P4, and P5 are the feature maps extracted from lateral connections, corresponding to C2, C3, C4, and C5 that are respectively of the same spatial sizes, while P6 and P7 are down-sampled by two 3×3 convolution layers after P5. The proposed RefineFace is based on our baseline face detector RetinaNet with five newly proposed modules:

- STR: It selects C5, C6, C7, P5, P6, and P7 to conduct two-step regression.
- STC: It selects C2, C3, C4, P2, P3, and P4 to perform two-step classification.
- SML: It adds the scale-aware margin to the classification loss to better distinguish faces from background across different scales.
- FSM: It contains one RoIAlign layer, four 3×3 convolution layers, one global average pooling layer and the Focal loss to let backbone learn more discriminative features for the classification task.
- RFE: It enriches the receptive field of features used to predict the classification and location of objects.

Without the above proposed five modules, it is our baseline face detector, consisting of C2-C5 and P2-P7 with the Focal loss and the smooth L_1 loss. As shown in Table 1, we aim to boost its regression and classification ability to obtain a new state-of-the-art method.

3.1 Selective Two-step Regression

Single-shot detectors conduct only one regression operation to get final detections from anchors. Comparing to two-stage detectors with multi-step regression, single-shot detectors with one-step regression lack of strong regression ability. This inadequacy causes lots of inaccurate detection results, which will be considered as false positives, especially under MS COCO-style evaluation standard. To this end, using cascade structure [14], [15] to conduct multi-step regression is an effective method to improve the regression ability for accurate detection bounding boxes.

However, blindly adding multi-step regression to the face detection task is often counterproductive. Specifically, experimental results in Table 2 indicate that applying two-step regression in the three lower pyramid levels impairs the performance. The reasons behind this phenomenon are twofold: 1) the three lower pyramid levels are associated with plenty of small anchors to detect small faces. These small faces are characterized by very coarse feature representations, so it is difficult for these small anchors to perform two-step regression; 2) in the training phase, if we let the network pay too much attention to the difficult regression task on the low pyramid levels, it will cause larger regression loss and hinder the more important classification task. In contrast, the three higher pyramid levels are associated with a small number of large anchors to detect large faces with detailed features, conducting two-step regression in these three levels is feasible and will improve the performance as shown in Table 2.

Based on the above analyses, we selectively perform two-step regression on the three higher pyramid levels. As shown in Figure 4(a), the STR coarsely adjusts the locations and sizes of anchors from high levels of detection layers to provide better initialization for the subsequent regressor, which can enhance the regression ability to regress more accurate locations of bounding boxes. The loss function of STR consists of two parts, which is shown below:

$$\mathcal{L}_{STR} = \frac{1}{N_{s1}} \sum_{i \in \Psi} [l_i^* = 1] \mathcal{L}_r(x_i, g_i^*) + \frac{1}{N_{s2}} \sum_{i \in \Phi} [l_i^* = 1] \mathcal{L}_r(t_i, g_i^*), \quad (1)$$

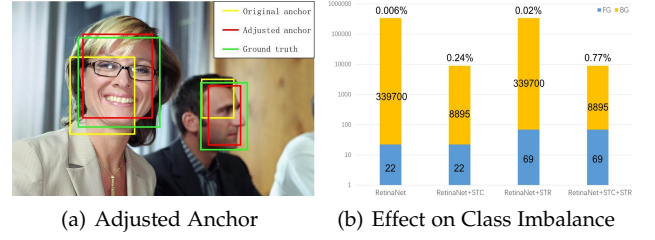


Fig. 4. (a) STR provides better initialization for the subsequent regressor. (b) STC increases the positives/negatives ratio by about 38 times.

where i is the index of anchor in a mini-batch, l_i^* and g_i^* are the ground truth class label and size of anchor i , x_i is the refined coordinates of anchor i in the first step, t_i is the coordinates of the bounding box in the second step, N_{s1} and N_{s2} are the numbers of positive anchors in the first and second steps, Ψ represents a collection of samples selected for two-step regression, and Φ represents a sample set in the second step. Similar to Faster R-CNN [37], we use the smooth L_1 loss as the regression loss \mathcal{L}_r . The Iverson bracket indicator function $[l_i^* = 1]$ outputs 1 when the condition is true, *i.e.*, $l_i^* = 1$ (the anchor is not the negative), and 0 otherwise. Hence $[l_i^* = 1] \mathcal{L}_r$ indicates that the regression loss is ignored for negative anchors.

3.2 Selective Two-step Classification

It is necessary for single-shot anchor-based face detectors to tile plenty of anchors over the image to detect faces of various scales, which causes the extreme class imbalance between the positive and negative samples. For example, in our RefineFace structure with the 1024×1024 input resolution, if we tile 2 anchors at each anchor point, the total number of samples will reach 300k. Among them, the number of positive samples is only a few dozen or less. To solve this issue, the two-step classification is introduced in RefineDet [15]. It is a kind of cascade classification implemented through a two-step network architecture, in which the first step filters out most negative anchors using a preset threshold $\theta = 0.99$ to reduce the search space for the subsequent step. Thus, the two-step classification can enhance the classification ability to reduce false positives.

However, it is unnecessary to perform two-step classification in all pyramid levels. Since the anchors tiled on the three higher levels (*i.e.*, P5, P6, and P7) only account for 11.1% and the associated features are much more adequate. Therefore, the classification task is relatively easy in these three higher pyramid levels. It is thus dispensable to apply two-step classification on the three higher pyramid levels, and if applied, it will lead to an increase in computation cost. In contrast, the three lower pyramid levels (*i.e.*, P2, P3, and P4) have the vast majority of samples (88.9%) and lack of adequate features. It is urgently needed for these low pyramid levels to do two-step classification in order to alleviate the class imbalance problem and reduce the search space for the subsequent classifier.

Therefore, our STC module selects C2, C3, C4, P2, P3, and P4 to perform two-step classification. As the statistical result shown in Figure 4(b), the STC increases the positive/negative sample ratio by approximately 38 times,

from around 1:15441 to 1:404. In addition, we use the Focal loss in both two steps to make full use of samples. Unlike RefineDet [15], the RefineFace shares the same classification module in the two steps, since they have the same task to distinguish the face from the background. The experimental results of applying the two-step classification on each pyramid level are shown in Table 4. Consistent with our analysis, the two-step classification on the three lower pyramid levels helps to improve performance, while on the three higher pyramid levels is ineffective.

The loss function for STC consists of two parts, *i.e.*, the loss in the first step and the second step. For the first step, we calculate the focal loss for those samples selected to perform two-step classification. And for the second step, we just focus on those samples that remain after the first step filtering. With these definitions, we define the loss function:

$$\mathcal{L}_{\text{STC}} = \frac{1}{N_{s_3}} \sum_{i \in \Omega} \mathcal{L}_{\text{FL}}(p_i, l_i^*) + \frac{1}{N_{s_4}} \sum_{i \in \Phi} \mathcal{L}_{\text{FL}}(q_i, l_i^*), \quad (2)$$

where p_i and q_i are the predicted confidence of the anchor i being a face in the first and second steps, N_{s_3} and N_{s_4} are the numbers of positive anchors in the first and second steps, Ω represents a collection of samples selected for two-step classification, l_i^* and Φ are the same as defined in STR. The binary classification loss \mathcal{L}_{FL} is the sigmoid focal loss over two classes (face *vs.* background).

3.3 Scale-aware Margin Loss

To further improve the classification ability of our baseline, we borrow the idea of the margin-based loss function [51], [52], [53], [54] from face recognition to face detection. The margin-based idea is a promising strategy in the face recognition task to improve the classification ability, which adds an extra margin to the classification loss to enhance the discrimination ability. Take the binary classification with sigmoid function for example. Supposed x is the prediction value before the sigmoid output for a sample. Then the margin-based prediction is:

$$y = \text{sigmoid}(x - m), \quad (3)$$

where m is the margin added to x and y is the prediction probability. After that, y is used in the classification loss, which can make the decision boundary more discriminative. Inspired by the margin-based loss function, we would like to add an extra margin to the sigmoid loss to improve the classification ability of face detectors.

However, the constant margin used in face recognition is not suitable for face detection, because the binary classification task in face detection needs to classify samples of different scales. Specifically, for our anchor-based baseline, different scales of anchors are classified to detect various scales of faces. Small anchors have a large amount and are classified using coarser features, while large anchors have a small amount and are classified using detailed features. Thus, different scales of anchors have different degrees of classification difficulty and the classification of smaller anchors is more difficult. To elaborate on it, we visualize the decision boundary of our baseline detector at different scales on the WIDER FACE validation Hard subset as: feeding each image to the trained model and recording all the

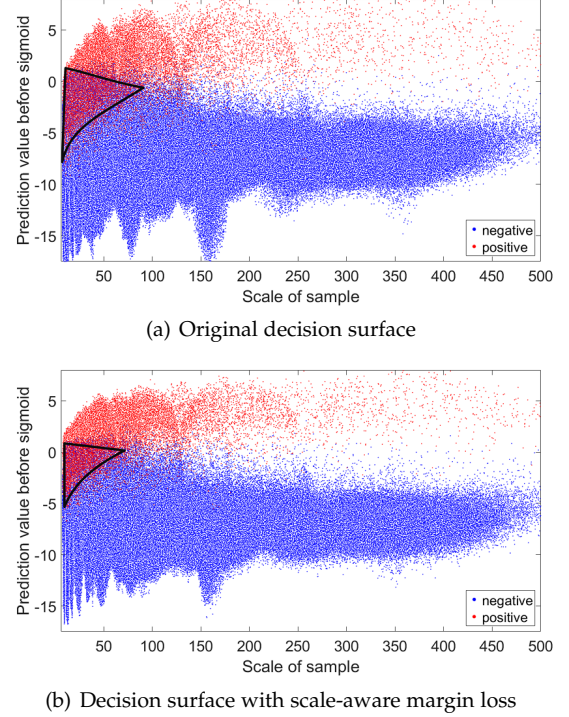


Fig. 5. Visualization of classification ability. Red and blue points are positive and negative samples. The black curve enclosed area indicates the mixed region where face and background are indistinguishable.

detections along with their scales (*i.e.*, the square root of width \times height) and prediction scores before sigmoid, then showing all positive and negative samples in Figure 5(a). We can see that the decision boundary of the baseline detector becomes blurred as the scales become smaller. It means the classification ability becomes weaker as the scale of faces becomes smaller.

To handle this problem, we propose a scale-aware margin loss (SML) function, which adjusts the margin for each sample by its scale. The margin for each sample is given as:

$$m = \alpha / \sqrt{wh}, \quad (4)$$

where α is a hyper-parameter to scale the margin, and w and h are the width and height of the sample. Larger faces have a more discriminative decision boundary so that they do not need too large margin; on the contrary, smaller faces have a blurred decision boundary so they need a larger margin to enhance the classification ability. After applying the scale-aware margin, the mixed region enclosed by the black curve becomes smaller as shown in Figure 5(b), suggesting the classification ability is enhanced for these small faces.

3.4 Feature Supervision Module

Another reason for the limited classification ability of single-shot face detector is that the learned features in the backbone network are not sufficiently discriminative because of misalignment. Single-shot detectors perform face detection based on regular and dense anchors using fully convolutional network. As shown in Figure 6, the features used to classify the anchors in single-shot face detectors are misaligned, *i.e.*, extracted from the corresponding receptive field and not tailored to the exact boundary of features in the

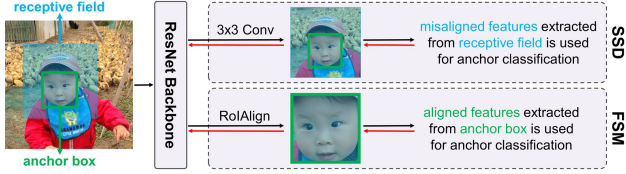


Fig. 6. Single-Shot Detectors (SSD) classify anchor using misaligned features extracted from receptive field, while FSM classifies anchor using aligned features extracted from anchor box. Black and red arrows mean forward and backward.

anchor box. This shortcoming will make the learned features in the backbone not discriminative enough.

To solve this issue, we design a feature supervision module (FSM) to enable the single-shot backbone network to learn more discriminative features. This module is appended after the backbone network and classifies the anchors using aligned features extracted from anchor box. The backbone will be updated by the classification loss of FSM to learn more discriminative features. Specifically, this module first uses RoIAlign [55] to extract the features at each detection, and then performs an extra binary classification based on the resultant features. This module has three characteristics: 1) we want to use it to enhance the classification ability, hence it only performs the binary classification; 2) it has a lightweight fully convolutional subnetwork with a relatively small loss, since it is an auxiliary module and should not over-dominate the training of the face detector; 3) it is not involved in the inference phase and will not introduce any additional overhead. With this additional supervision module, the backbone network is forced to learn more discriminative features for classification without any extra overhead during the inference phase.

To train the feature supervision module, we apply NMS with a threshold of 0.7, add ground truth boxes, and distribute 512 prediction proposals to the pyramid level from which they come to sample their RoI features. As shown in Figure 3, the RoIAlign operation is performed at the assigned feature layers, yielding 5×5 resolution features, which are fed into three subsequent convolutional layers, a prediction convolutional layer and a global average pooling layer to classify between face and background.

3.5 Receptive Field Enhancement

At present, most detection networks utilize ResNet and VGGNet as the basic feature extraction module, while both of them possess square receptive fields. The singleness of the receptive field affects the detection of objects with different aspect ratios. This issue seems unimportant in face detection task, because the aspect ratio of face annotations is about 1:1 in many datasets. Nevertheless, statistics show that the WIDER FACE training set has a considerable part of faces that have an aspect ratio of more than 2 or less than 0.5. Consequently, there is a mismatch between the receptive field of network and the aspect ratio of faces.

To address this issue, we propose a module named Receptive Field Enhancement (RFE) to diversify the receptive field of features before predicting classes and locations. In particular, RFE module replaces the middle two convolution layers in the class subnet and the box subnet of RetinaNet.

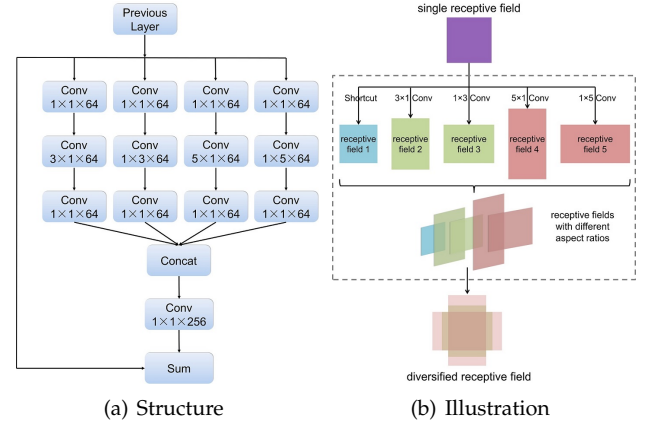


Fig. 7. Structure and illustration of Receptive Field Enhancement (RFE).

The structure of RFE is shown in Figure 7(a). Our RFE module adopts a four-branch structure, which is inspired by the Inception block [56]. To be specific, first, we use a 1×1 convolution layer to decrease the channel number to one quarter of the previous layer. Second, we use $1 \times k$ and $k \times 1$ ($k = 3$ and 5) convolution layer to provide rectangular receptive field. Through another 1×1 convolution layer, the feature maps from four branches are concatenated together. Additionally, we apply a shortcut path to retain the original receptive field from previous layer. As shown in Figure 7(b), the RFE provides more diverse receptive fields that is helpful for detecting extreme-pose faces.

3.6 Training and Inference

Data Augmentation. To prevent over-fitting and construct a robust model, several data augmentation strategies are used to adapt to face variations, described as follows: 1) applying some photometric distortions to the training images; 2) expanding the images with a random factor in the interval $[1, 2]$ by mean-padding; 3) cropping two square patches and randomly selecting one for training. One patch is with the size of the image's shorter side and the other one is with the size determined by multiplying a random number in the interval $[0.5, 1.0]$ by the image's shorter side; 4) flipping the selected patch randomly and resizing it to 1024×1024 to get the final training sample.

Anchor Design. At each location of a detection layer, we associate two scales of anchors (corresponding to $2S$ and $2\sqrt{2}S$ in the original image, where S is the downsampling factor of the detection layer) with one 1.25 aspect ratio. So, there are $A = 2$ anchors at each location of a detection layer, and it covers the scale of $8 - 362$ pixels in an input image.

Sample Matching. During the training phase, the preset anchors and proposals need to be assigned as positive and negative samples for training. We assign anchors to ground-truth boxes using an intersection-over-union (IoU) threshold of θ_p , and to background if their IoU is in $[0, \theta_n)$. If an anchor is unassigned, which may happen with an IoU in $[\theta_n, \theta_p)$, it is discarded during training. Empirically, we set $\theta_n = 0.3$ and $\theta_p = 0.7$ for the first step in STR and STC, $\theta_n = 0.4$ and $\theta_p = 0.5$ for the second step in STR and STC, and $\theta_n = 0.4$ and $\theta_p = 0.7$ for the feature supervision module.

Loss Function. We use a hybrid loss function to train the proposed model in an end-to-end fashion as $\mathcal{L} = \mathcal{L}_{STR} + \mathcal{L}_{STC} + \mathcal{L}_{FSM}$, where \mathcal{L}_{STR} is the loss function for STR, \mathcal{L}_{STC} is the loss function of STC, and \mathcal{L}_{FSM} is the Focal loss for the binary classification in FAM.

Optimization. We use the ImageNet pretrained model to initialize the backbone network. All the parameters in the newly added convolution layers are initialized by the “xavier” method. The stochastic gradient descent (SGD) algorithm is applied to fine-tune the RefineFace model with 0.9 momentum, 0.0001 weight decay and 32 batch size. We use the warmup strategy to gradually ramp up the learning rate from 3.125×10^{-4} to 1×10^{-2} at the first 5 epochs. After that, it switches to the regular learning rate schedule, *i.e.*, dividing by 10 at 100 and 120 epochs and ending at 130 epochs. We use the PyTorch [57] library to implement the proposed RefineFace.

Inference. During the inference phase, the STC first filters the regularly tiled anchors on the selected pyramid levels with the negative confidence scores larger than the threshold $\theta = 0.99$, and then STR adjusts the locations and sizes of selected anchors. After that, the second step takes over these refined anchors and outputs top 5,000 detection results whose confidence scores are all higher than the threshold of 0.05. Finally, we apply the non-maximum suppression (NMS) algorithm with Jaccard overlap of 0.4 to generate top 750 high confident detections per image as final results.

4 EXPERIMENTS

In this section, we conduct extensive experiments on WIDER FACE [7], AFW [9], PASCAL Face [10], FDDB [11] and MAFA [58] to verify the effectiveness of our RefineFace. Notably, our final model trained on WIDER FACE is directly evaluated on other datasets without finetuning.

4.1 WIDER FACE Dataset

It has 393,703 annotated faces with variations in pose, scale, facial expression, occlusion and lighting condition in 32,203 images. These images are split into three subsets: training (40%), validation (10%) and testing (50%) sets. Each subset has three levels of difficulty: Easy, Medium and Hard based on the detection rate of EdgeBox [59]. All the models are trained on the training subset and tested on both the validation and testing subsets. Since the annotations of testing subsets are held-out, we submit the detection results to the collectors to report final evaluation results.

4.1.1 Model Analysis

To demonstrate the effectiveness of our proposed modules in RefineFace, each of them is added to our baseline to examine how it affects the final performance on the WIDER FACE dataset. All the experiments are based on ResNet-50 and use the same parameter settings for a fair comparison, except for specified changes. As listed in the first column of Table 1, our baseline detector based on RetinaNet achieves 95.1% (Easy), 93.9% (Medium) and 88.0% (Hard) on the validation subset, which is better than most of face detectors on the WIDER FACE dataset. It can be considered as a strong

TABLE 1
Effectiveness of our proposed modules. Based on ResNet-50, all models are trained on WIDER FACE training subset and evaluated with AP (%) on validation subset.

Module	RefineFace									
STR	✓						✓	✓	✓	✓
STC		✓					✓	✓	✓	✓
RFE			✓					✓	✓	✓
SML				✓					✓	✓
FSM					✓					✓
<i>Easy</i>	95.1	95.9	95.3	95.5	95.5	95.8	96.1	96.4	96.6	96.9
<i>Medium</i>	93.9	94.8	94.4	94.3	94.6	94.5	95.0	95.3	95.6	95.9
<i>Hard</i>	88.0	88.8	89.4	88.3	89.1	88.7	90.1	90.2	90.7	91.1

single-shot face detector baseline and in the following, we will enhance its regression and classification ability to set a new state-of-the-art performance.

Selective Two-step Regression. We add the STR module to our baseline detector to verify its effectiveness. As shown in Table 1, it produces much better results than the baseline, with 0.8%, 0.9% and 0.8% AP improvements on the Easy, Medium, and Hard subsets. Experimental results of applying two-step regression to each pyramid level (see Table 2) confirm our previous analysis. Inspired by the detection evaluation metric of MS COCO, we use 4 IoU thresholds $\{0.5, 0.6, 0.7, 0.8\}$ to compute the AP, so as to prove that the STR module can produce more accurate localization. As shown in Table 3, the STR module produces consistently accurate detection results than the baseline method. The gap between the AP across all three subsets increases as the IoU threshold increases, which indicates that the STR module enhances the regression ability of our baseline and produces more accurate detections.

TABLE 2
AP (%) of the two-step regression applied to each pyramid level.

STR	B	P2	P3	P4	P5	P6	P7
<i>Easy</i>	95.1	94.8	94.3	94.8	95.4	95.7	95.6
<i>Medium</i>	93.9	93.4	93.7	93.9	94.2	94.4	94.6
<i>Hard</i>	88.0	87.5	87.7	87.0	88.2	88.2	88.4

TABLE 3
AP (%) at different IoU thresholds on the WIDER FACE Hard subset.

IoU	0.5	0.6	0.7	0.8
RetinaNet	88.0	76.4	57.8	28.5
RetinaNet+STR	88.8	83.4	66.5	38.2

Selective Two-step Classification. Experimental results of applying two-step classification to each pyramid level are shown in Table 4, indicating that applying two-step classification to the low pyramid levels improves the performance, especially on tiny faces. Therefore, the STC module selectively applies the two-step classification on the low pyramid levels (*i.e.*, P2, P3, and P4), since these levels are associated with lots of small anchors, which are the main source of false positives. As shown in Table 1, we find that after using the STC module, the AP scores of the detector are

improved from 95.1%, 93.9% and 88.0% to 95.3%, 94.4% and 89.4% on the Easy, Medium and Hard subsets, respectively. In order to verify whether the improvements benefit from reducing the false positives, we count the number of false positives under different recall rates. As listed in Table 5, our STC effectively reduces the false positives across different recall rates, demonstrating the effectiveness of the STC module. In addition, coupled with the STR module as listed in the seventh column of Table 1, the performance is further improved to 96.1%, 95.0% and 90.1% on the Easy, Medium and Hard subsets, respectively.

TABLE 4

AP (%) of the two-step classification applied to each pyramid level.

STC	B	P2	P3	P4	P5	P6	P7
<i>Easy</i>	95.1	95.2	95.2	95.2	95.0	95.1	95.0
<i>Medium</i>	93.9	94.2	94.3	94.1	93.9	93.7	93.9
<i>Hard</i>	88.0	88.9	88.7	88.5	87.8	88.0	87.7

TABLE 5

Number of false positives at different recall rates.

Recall (%)	10	30	50	80	90	95
# FP of RetinaNet	3	24	126	2,801	27,644	466,534
# FP of RetinaNet+STC	1	20	101	2,124	13,163	103,586

Receptive Field Enhancement. The RFE is used to diversify the receptive fields of detection layers in order to capture faces with extreme poses. Comparing the detection results between first and fourth columns in Table 1, adding RFE to our baseline improves the AP performances by 0.4%, 0.4% and 0.3% for the Easy, Medium, and Hard subsets, respectively. Even though using RFE after STR and STC, it still consistently increases the AP scores in different subsets, *i.e.*, from 96.1% to 96.4% for Easy, from 95.0% to 95.3% for Medium and from 90.1% to 90.2% for Hard. These improvements can be mainly attributed to the diverse receptive fields, which is useful to capture various pose faces for better detection accuracy.

Scale-aware Margin Loss. We first only apply the scale-aware margin to the classification loss function of our baseline. Comparing the first and fifth columns in Table 1, we can observe that it improves the AP scores by 0.4%, 0.7% and 1.1% for the Easy, Medium, and Hard subsets respectively, benefiting from better discrimination between face and background, especially for small faces as shown in Figure 5(b). There is a hyper-parameter in Equation 4 to scale the margin, we conduct several experiments to study its effect. We train the model with different α in [3, 7, 11, 15, 19, 23] on the WIDER FACE training set, then test on the validation set. As shown in Table 6, we observe that the proposed detector is relatively insensitive to the variations of α . Too small value (*e.g.*, $\alpha = 3$) will make the margin not work, while too large (*e.g.*, $\alpha = 23$) would cause it difficult for training to be optimized. Thus, we choose $\alpha = 15$ based on the validation performance in our experiments. Finally, we apply SML after STR, STC and RFE with a high starting point and the AP performances are still improved from 96.4%, 95.3% and 90.2% to 96.6%,

95.6% and 90.7% on the Easy, Medium and Hard subsets respectively. These results demonstrate its effectiveness.

TABLE 6
AP (%) of different α in the scale-aware margin loss.

α	0	3	7	11	15	19	23
<i>Easy</i>	95.1	95.1	95.3	95.5	95.5	95.3	95.0
<i>Medium</i>	93.9	93.9	94.3	94.5	94.6	94.4	94.0
<i>Hard</i>	88.0	88.2	88.6	89.0	89.1	89.1	88.9

Feature Supervision Module. To verify the effectiveness of the feature supervision module, we append it after our baseline and train the whole network end-to-end. As listed in the sixth column of Table 1, it boosts the AP results of our baseline by 0.7%, 0.6% and 0.7% for Easy, Medium and Hard subsets. These improvements come from the more discriminative features learned by the backbone network with the help of FSM. The output size of RoIAlign is a hyper-parameter and we conduct several experiments to select it. As shown in Table 7, the moderate size 5×5 has the best results and too large or too small size will impair the performance. Besides, except the convolution (Conv) type, FAM can also be designed in the fully connected (FC) type. As shown in Table 8, the Conv type achieves better performances with less parameters than the FC type, benefiting from that the Conv type shares parameters and retains spatial information. Notably, this module is only involved during training without any additional overhead during inference. Finally, FSM is added after other four modules, which still boost the AP performances from 96.6%, 95.6%, 90.7% to 96.9%, 95.9%, 91.1% on the Easy, Medium and Hard subsets respectively.

TABLE 7
AP (%) of different output sizes of RoIAlign in FAM.

Size	3×3	5×5	7×7
<i>Easy</i>	95.5	95.8	95.7
<i>Medium</i>	94.3	94.5	94.4
<i>Hard</i>	88.6	88.7	88.5

TABLE 8
AP (%) of different design types of FAM.

Type	# parameter	<i>Easy</i>	<i>Medium</i>	<i>Hard</i>
None	0	95.1	93.9	88.0
FC	829,472	95.7	94.3	88.5
Conv	387,360	95.8	94.5	88.7

4.1.2 Tradeoff Analysis

All the experiments in last subsection are based on ResNet-50 and in this section we will analysis the trade-off between speed and accuracy of different backbone networks. The inference time of RefineFace are measured on a single NVIDIA GTX 1080-Ti with CUDA 9.0 and cuDNN v7.0. Frame-per-second (FPS) and millisecond (ms) are used to compare the speed. As listed in Table 9, our best model with ResNet-152 backbone can run at 17.7 FPS (56.6 ms) for the VGA-resolution image and ResNet-101 achieves promising

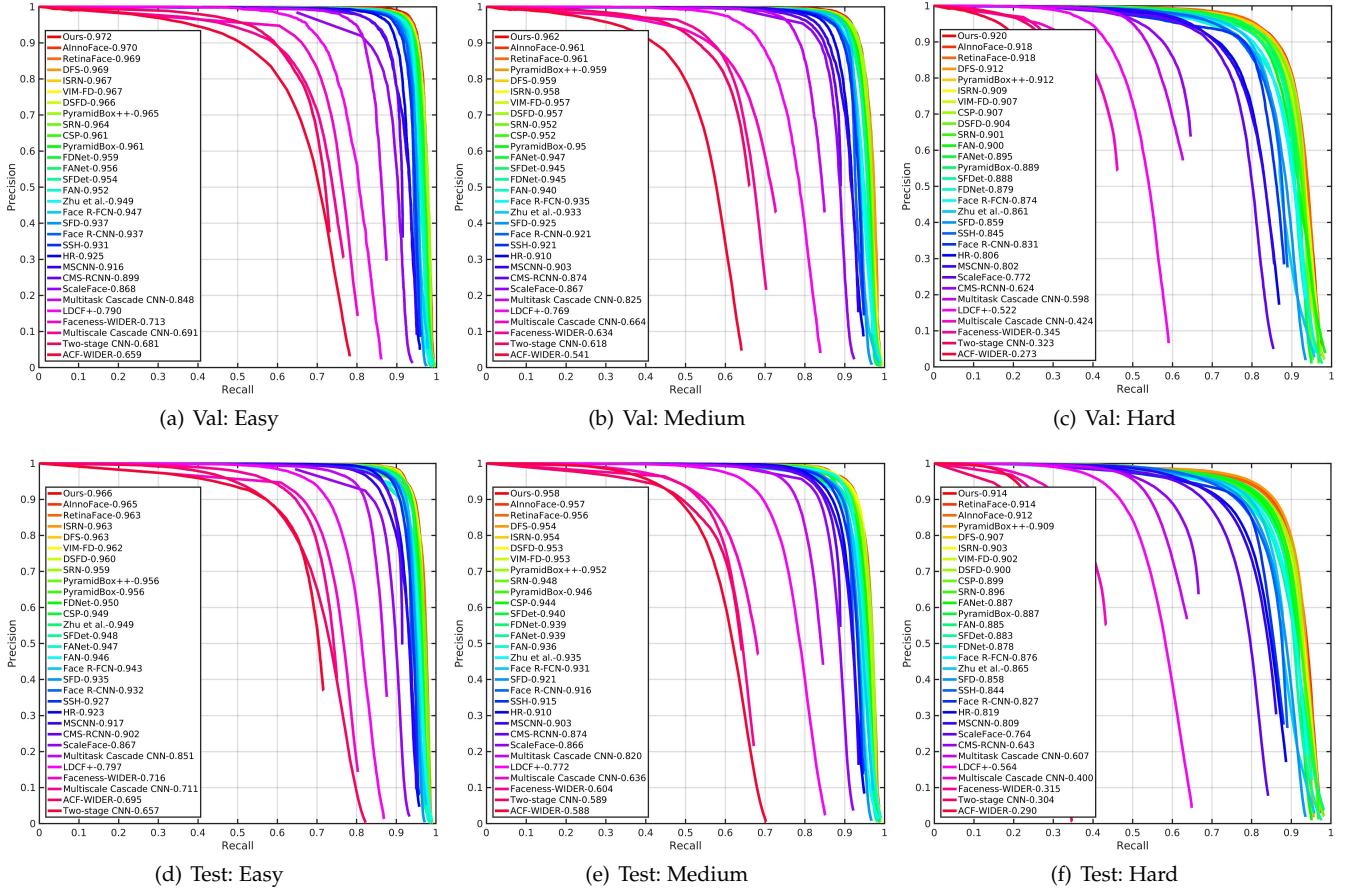


Fig. 8. Precision-recall curves on WIDER FACE validation and testing subsets.

performances with near real-time speed. We also utilize more lightweight backbones to balance between speed and accuracy. Based on more efficient backbones, the proposed model can achieve real-time speed (≥ 25 FPS) for the VGA-resolution image, *i.e.*, 28.5 FPS (35.1 ms) for ResNet-50 and 37.3 FPS (26.8 ms) for ResNet-18. Among these backbones, ResNet-152 has the best performance with the slowest speed, ResNet-18 reaches the fastest speed with the lowest AP, while ResNet-50 has the best balance between speed and accuracy, which is why most of models [2], [44], [47], [48], [49] use ResNet-50 as the backbone network to build the face detectors.

TABLE 9

Tradeoff analysis between speed and accuracy. Speed is measured with a VGA-resolution (640×480) input image. All batch normalization (BN) layers are merged into the convolution layers during inference.

Backbone	Speed		Accuracy (%)		
	FPS	ms	Easy	Medium	Hard
ResNet-18	37.3	26.8	96.3	95.1	90.2
ResNet-50	28.5	35.1	96.9	95.9	91.1
ResNet-101	22.8	43.9	97.1	96.1	91.6
ResNet-152	17.7	56.6	97.2	96.2	92.0

4.1.3 Performance Analysis

As shown in Figure 8, we compare RefineFace with 29 state-of-the-art face detection methods [2], [3], [4], [5], [6],

[7], [24], [25], [27], [28], [29], [30], [36], [38], [39], [40], [44], [45], [47], [48], [49], [60], [61], [62], [63], [64], [65], [66] on both the validation and testing subsets. The proposed RefineFace achieves the best AP performance in all subsets of both validation and testing sets, *i.e.*, 97.2% (Easy), 96.2% (Medium) and 92.0% (Hard) for validation set, and 96.6% (Easy), 95.8% (Medium) and 91.4% (Hard) for testing set. It outperforms all compared state-of-the-art methods based on the average precision (AP) across the three subsets, demonstrating the superiority of the proposed face detector. Notably, among all the published methods [2], [3], [4], [5], [6], [7], [24], [25], [27], [29], [30], [60], [61], [63], [64], our method outperforms the previous best method DSFD [3] by a large margin. Although there are several unpublished technical reports have very promising performances, some of them [49], [65] use additional data, some of them [47], [48], [62], [66] apply some existing time-consuming tricks including segmentation, attention and context. In contrast, the proposed method presents five new modules with a small amount of additional overhead and the five modules are complementary to existing methods.

4.2 AFW dataset

It contains 473 labeled faces in 205 images, which are collected from Flickr and have cluttered backgrounds with large variations in both face viewpoints and appearances (*e.g.*, ages, sunglasses, make-ups, skin colors, expressions,

etc.). As shown by the precision-recall curves in Figure 9, we compare RefineFace against three commercial face detectors (*i.e.*, Face.com, Face++ and Picasa) and nine state-of-the-art methods [9], [10], [17], [20], [24], [67], [68]. The proposed method improves the AP score of state-of-the-art results by 1.55% compared with the second best method STN [67].

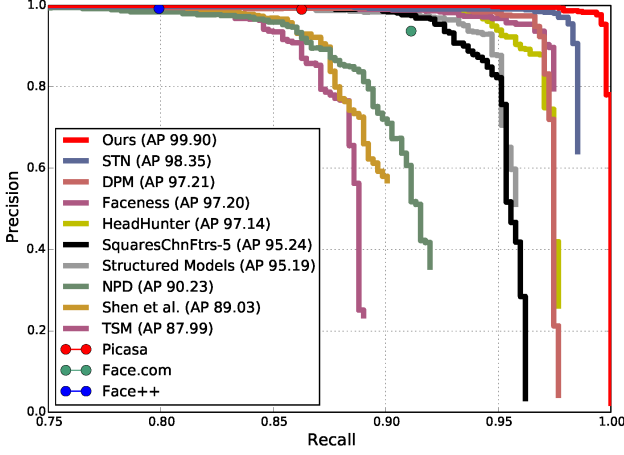


Fig. 9. Precision-recall curves on the AFW dataset.

4.3 PASCAL Face dataset

It is collected from PASCAL person layout test subset and consists of 1,335 labeled faces in 851 images with large face appearance and pose variations. Figure 10 shows the precision-recall curves of the proposed RefineFace compared with 9 state-of-the-art methods [9], [10], [20], [24], [67], [69] and 3 commercial face detectors (*i.e.*, SkyBiometry, Face++ and Picasa). The RefineFace model outperforms the state-of-the-art methods with the top AP score (99.45%).

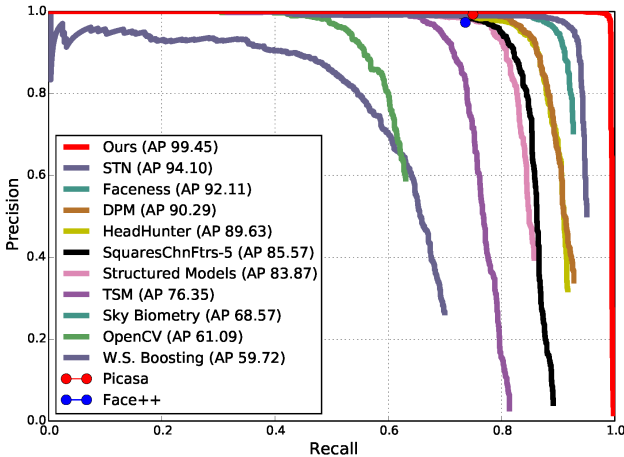


Fig. 10. Precision-recall curves on the PASCAL Face dataset.

4.4 FDDB dataset

It has 5,171 faces annotated in 2,845 images with a wide range of challenges including low image resolutions, severe occlusions and difficult poses. Since there are lots of unlabelled faces on FDDB, which results in many false positive

faces with high scores. Hence, we use the new annotations [64] to evaluate the proposed detector and compare it with several state-of-the-art methods [3], [5], [6], [17], [24], [25], [26], [27], [29], [39], [42], [45], [70], [71], [72], [73], [74], [75], [76], [77], [78], [79], [80], [81], [82], [83], [84] in Figure 11. The proposed face detector achieves 99.11% true positive rate when the number of false positives is equal to 1,000, setting a new state-of-the-art result. It indicates the superior performance of RefineFace in presence of various scales, large appearance variations, heavy occlusions and severe blur degradations in unconstrained scenarios.

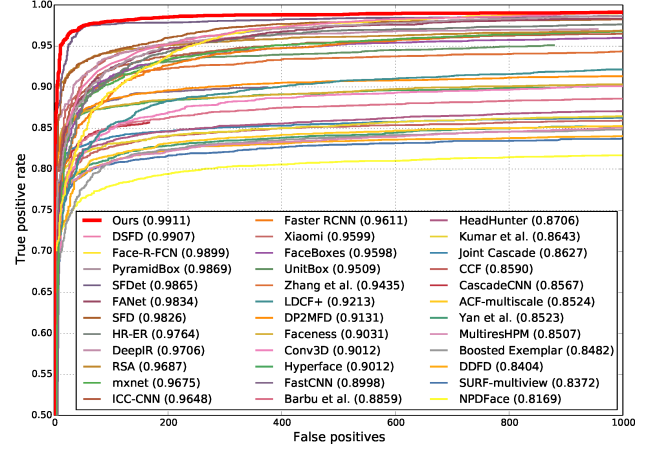


Fig. 11. Receiver operating characteristic (ROC) curves with discrete scores on the FDDB dataset. The number in the legend is the true positive rate (TPR) at the false positives (FP) equals to 1,000.

4.5 MAFA Dataset

It is a benchmark dataset for detecting occluded faces and contains 35,806 masked faces in 30,811 images collected from Internet. This dataset covers 60 cases of occluded faces in daily scenarios with three degrees of occlusions, four types of masks and five face orientations. In this dataset, if the face is severely blurred, deformed, or smaller than 32 pixels, it is ignored. The left are divided into masked faces and unmasked faces. Therefore, there are three subsets on the MAFA dataset: (1) the whole subset (containing 6,354 masked faces, 996 unmasked faces and 2,683 ignored faces), (2) the masked subset (including masked faces), and (3) the unignored subset (consisting of masked faces and unmasked faces). We directly evaluate the model trained on WIDER FACE with all three subsets and report the average precision (AP) scores against eight state-of-the-art methods in Table 10. Our method sets a new state-of-the-art result on all three subsets, *i.e.*, 83.9% on the whole subset, 96.2% on the masked subset and 95.7% on the unignored subset. These results demonstrate that the proposed method is robust to various occlusions in our daily scenarios.

5 CONCLUSION

This paper proposes a state-of-the-art single-shot face detector by enhancing the regression and classification ability. On the one hand, boosting the regression ability can improve the location accuracy and reduce the LOC error, for this purpose we design the STR to coarsely adjust the locations

TABLE 10
AP (%) on the MAFA testing set.

Methods	Whole set	Masked set	Unignored set
TSM [9]	-	-	41.6
HeadHunter [20]	-	-	50.9
HPM [72]	-	-	60.0
MTCNN [27]	-	-	60.8
LLE-CNNs [58]	-	-	76.4
FAN [44]	-	76.5	88.3
AOFD [85]	81.3	83.5	91.9
SFDet [64]	81.5	94.0	93.7
Ours	83.9	96.2	95.7

and sizes of anchors from high level detection layers to provide better initialization for the subsequent regressor. On the other hand, enhancing the classification ability can improve the recall efficiency and reduce the CLS error, to this end we first use the STC to filter out most simple negatives from low level detection layers to reduce the search space for the subsequent classifier, then apply the SML to better distinguish faces from background at various scales and the FSM to let the backbone network learn more discriminative features for classification. In addition, we introduce the RFE to provide more diverse receptive field to better capture faces in some extreme poses. Experiments are conducted on most of challenging face detection datasets to demonstrate the effectiveness of RefineFace. In the future, we plan to design a lightweight architecture with the help of automatic machine learning (AutoML) methods to make RefineFace run in real-time not only on the GPU device but also on the CPU device.

REFERENCES

- [1] P. A. Viola and M. J. Jones, "Robust real-time face detection," *IJCV*, 2004.
- [2] C. Chi, S. Zhang, J. Xing, Z. Lei, S. Z. Li, and X. Zou, "Selective refinement network for high performance face detection," in *AAAI*, 2019.
- [3] J. Li, Y. Wang, C. Wang, Y. Tai, J. Qian, J. Yang, C. Wang, J. Li, and F. Huang, "Dsfd: dual shot face detector," in *CVPR*, 2019.
- [4] M. Najibi, P. Samangouei, R. Chellappa, and L. S. Davis, "SSH: single stage headless face detector," in *ICCV*, 2017.
- [5] X. Tang, D. K. Du, Z. He, and J. Liu, "Pyramidbox: A context-assisted single shot face detector," in *ECCV*, 2018.
- [6] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Z. Li, "S³FD: Single shot scale-invariant face detector," in *ICCV*, 2017.
- [7] S. Yang, P. Luo, C. C. Loy, and X. Tang, "WIDER FACE: A face detection benchmark," in *CVPR*, 2016.
- [8] T. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *ICCV*, 2017.
- [9] X. Zhu and D. Ramanan, "Face detection, pose estimation, and landmark localization in the wild," in *CVPR*, 2012.
- [10] J. Yan, X. Zhang, Z. Lei, and S. Z. Li, "Face detection by structural models," *IVC*, 2014.
- [11] V. Jain and E. Learned-Miller, "Fddb: A benchmark for face detection in unconstrained settings," University of Massachusetts, Amherst, Tech. Rep., 2010.
- [12] T. Lin, M. Maire, S. J. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: common objects in context," in *ECCV*, 2014.
- [13] S. Gidaris and N. Komodakis, "Object detection via a multi-region and semantic segmentation-aware CNN model," in *ICCV*, 2015.
- [14] Z. Cai and N. Vasconcelos, "Cascade R-CNN: delving into high quality object detection," in *CVPR*, 2018.
- [15] S. Zhang, L. Wen, X. Bian, Z. Lei, and S. Z. Li, "Single-shot refinement neural network for object detection," in *CVPR*, 2018.
- [16] S. C. Brubaker, J. Wu, J. Sun, M. D. Mullin, and J. M. Rehg, "On the design of cascades of boosted ensembles for face detection," *IJCV*, 2008.
- [17] S. Liao, A. K. Jain, and S. Z. Li, "A fast and accurate unconstrained face detector," *TPAMI*, 2016.
- [18] M. Pham and T. Cham, "Fast training and selection of haar features using statistics in boosting-based face detection," in *ICCV*, 2007.
- [19] P. F. Felzenszwalb, R. B. Girshick, D. A. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *TPAMI*, 2010.
- [20] M. Mathias, R. Benenson, M. Pedersoli, and L. J. V. Gool, "Face detection without bells and whistles," in *ECCV*, 2014.
- [21] J. Yan, Z. Lei, L. Wen, and S. Z. Li, "The fastest deformable part model for object detection," in *CVPR*, 2014.
- [22] H. Li, Z. Lin, X. Shen, J. Brandt, and G. Hua, "A convolutional neural network cascade for face detection," in *CVPR*, 2015.
- [23] H. Qin, J. Yan, X. Li, and X. Hu, "Joint training of cascaded CNN for face detection," in *CVPR*, 2016.
- [24] S. Yang, P. Luo, C. C. Loy, and X. Tang, "From facial parts responses to face detection: A deep learning approach," in *ICCV*, 2015.
- [25] E. Ohn-Bar and M. M. Trivedi, "To boost or not to boost? on the limits of boosted trees for object detection," in *ICPR*, 2016.
- [26] J. Yu, Y. Jiang, Z. Wang, Z. Cao, and T. S. Huang, "Unitbox: An advanced object detection network," in *ACMMM*, 2016.
- [27] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," *SPL*, 2016.
- [28] S. Yang, Y. Xiong, C. C. Loy, and X. Tang, "Face detection through scale-friendly deep convolutional networks," *CoRR*, 2017.
- [29] P. Hu and D. Ramanan, "Finding tiny faces," in *CVPR*, 2017.
- [30] C. Zhu, R. Tao, K. Luu, and M. Savvides, "Seeing small faces from robust anchors perspective," in *CVPR*, 2018.
- [31] Z. Hao, Y. Liu, H. Qin, J. Yan, X. Li, and X. Hu, "Scale-aware face detection," in *CVPR*, 2017.
- [32] G. Song, Y. Liu, M. Jiang, Y. Wang, J. Yan, and B. Leng, "Beyond trade-off: Accelerate fcn-based face detector with higher accuracy," in *CVPR*, 2018.
- [33] X. Shi, S. Shan, M. Kan, S. Wu, and X. Chen, "Real-time rotation-invariant face detection with progressive calibration networks," in *CVPR*, 2018.
- [34] Y. Bai, Y. Zhang, M. Ding, and B. Ghanem, "Finding tiny faces in the wild with generative adversarial network," in *CVPR*, 2018.
- [35] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *NIPS*, 2014.
- [36] C. Zhu, Y. Zheng, K. Luu, and M. Savvides, "CMS-RCNN: contextual multi-scale region-based CNN for unconstrained face detection," *CoRR*, 2016.
- [37] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal networks," *TPAMI*, 2017.
- [38] H. Wang, Z. Li, X. Ji, and Y. Wang, "Face r-cnn," *CoRR*, 2017.
- [39] Y. Wang, X. Ji, Z. Zhou, H. Wang, and Z. Li, "Detecting faces using region-based fully convolutional networks," *CoRR*, 2017.
- [40] C. Zhang, X. Xu, and D. Tu, "Face detection using improved faster RCNN," *CoRR*, 2018.
- [41] J. Dai, Y. Li, K. He, and J. Sun, "R-FCN: object detection via region-based fully convolutional networks," in *NIPS*, 2016.
- [42] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Z. Li, "Faceboxes: A CPU real-time face detector with high accuracy," in *IJCB*, 2017.
- [43] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. E. Reed, C. Fu, and A. C. Berg, "SSD: single shot multibox detector," in *ECCV*, 2016.
- [44] J. Wang, Y. Yuan, and G. Yu, "Face attention network: An effective face detector for the occluded faces," *CoRR*, 2017.
- [45] J. Zhang, X. Wu, J. Zhu, and S. C. H. Hoi, "Feature agglomeration networks for single stage face detection," *CoRR*, 2017.
- [46] T. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie, "Feature pyramid networks for object detection," in *CVPR*, 2017.
- [47] W. Tian, Z. Wang, H. Shen, W. Deng, B. Chen, and X. Zhang, "Learning better features for face detection with feature fusion and segmentation supervision," *CoRR*, 2018.
- [48] Y. Zhang, X. Xu, and X. Liu, "Robust and high performance face detector," *CoRR*, 2019.

- [49] 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," *CoRR*, 2019.
- [50] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [51] J. Deng, J. Guo, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," *CoRR*, 2018.
- [52] W. Liu, Y. Wen, Z. Yu, and M. Yang, "Large-margin softmax loss for convolutional neural networks," in *ICML*, 2016.
- [53] F. Wang, J. Cheng, W. Liu, and H. Liu, "Additive margin softmax for face verification," *SPL*, 2018.
- [54] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu, "Cosface: Large margin cosine loss for deep face recognition," in *CVPR*, 2018.
- [55] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask R-CNN," in *ICCV*, 2017.
- [56] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. E. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *CVPR*, 2015.
- [57] A. Paszke, S. Gross, S. Chintala, and G. Chanan, "Pytorch," 2017.
- [58] S. Ge, J. Li, Q. Ye, and Z. Luo, "Detecting masked faces in the wild with lle-cnns," in *CVPR*, 2017.
- [59] C. L. Zitnick and P. Dollár, "Edge boxes: Locating object proposals from edges," in *ECCV*, 2014.
- [60] Z. Cai, Q. Fan, R. S. Feris, and N. Vasconcelos, "A unified multi-scale deep convolutional neural network for fast object detection," in *ECCV*, 2016.
- [61] B. Yang, J. Yan, Z. Lei, and S. Z. Li, "Aggregate channel features for multi-view face detection," in *IJCB*, 2014.
- [62] F. Zhang, X. Fan, G. Ai, J. Song, Y. Qin, and J. Wu, "Accurate face detection for high performance," *CoRR*, 2019.
- [63] W. Liu, S. Liao, W. Ren, W. Hu, and Y. Yu, "High-level semantic feature detection: A new perspective for pedestrian detection," in *CVPR*, 2019.
- [64] S. Zhang, L. Wen, H. Shi, Z. Lei, S. Lyu, and S. Z. Li, "Single-shot scale-aware network for real-time face detection," *IJCV*, 2019.
- [65] J. Deng, J. Guo, Y. Zhou, J. Yu, I. Kotsia, and S. Zafeiriou, "Retinaface: Single-stage dense face localisation in the wild," *CoRR*, 2019.
- [66] Z. Li, X. Tang, J. Han, J. Liu, and R. He, "Pyramidbox++: High performance detector for finding tiny face," *CoRR*, 2019.
- [67] D. Chen, G. Hua, F. Wen, and J. Sun, "Supervised transformer network for efficient face detection," in *ECCV*, 2016.
- [68] X. Shen, Z. Lin, J. Brandt, and Y. Wu, "Detecting and aligning faces by image retrieval," in *CVPR*, 2013.
- [69] Z. Kalal, J. Matas, and K. Mikolajczyk, "Weighted sampling for large-scale boosting," in *BMVC*, 2008.
- [70] A. Barbu, N. Lay, and G. Gramajo, "Face detection with a 3d model," *CoRR*, 2014.
- [71] S. S. Farfate, M. J. Saberian, and L. Li, "Multi-view face detection using deep convolutional neural networks," in *ICMR*, 2015.
- [72] G. Ghiasi and C. C. Fowlkes, "Occlusion coherence: Detecting and localizing occluded faces," *CoRR*, 2015.
- [73] H. Jiang and E. G. Learned-Miller, "Face detection with the faster R-CNN," in *FG*, 2017.
- [74] V. Kumar, A. M. Namboodiri, and C. V. Jawahar, "Visual phrases for exemplar face detection," in *ICCV*, 2015.
- [75] H. Li, Z. Lin, J. Brandt, X. Shen, and G. Hua, "Efficient boosted exemplar-based face detection," in *CVPR*, 2014.
- [76] J. Li and Y. Zhang, "Learning SURF cascade for fast and accurate object detection," in *CVPR*, 2013.
- [77] Y. Li, B. Sun, T. Wu, and Y. Wang, "Face detection with end-to-end integration of a convnet and a 3d model," in *ECCV*, 2016.
- [78] Y. Liu, H. Li, J. Yan, F. Wei, X. Wang, and X. Tang, "Recurrent scale approximation for object detection in CNN," in *ICCV*, 2017.
- [79] R. Ranjan, V. M. Patel, and R. Chellappa, "A deep pyramid deformable part model for face detection," in *BTAS*, 2015.
- [80] —, "Hyperface: A deep multi-task learning framework for face detection, landmark localization, pose estimation, and gender recognition," *TPAMI*, 2019.
- [81] X. Sun, P. Wu, and S. C. H. Hoi, "Face detection using deep learning: An improved faster RCNN approach," *Neurocomputing*, 2018.
- [82] D. Triantafyllidou and A. Tefas, "A fast deep convolutional neural network for face detection in big visual data," in *INNS Conference on Big Data*, 2016.
- [83] S. Wan, Z. Chen, T. Zhang, B. Zhang, and K. Wong, "Bootstrapping face detection with hard negative examples," *CoRR*, 2016.
- [84] K. Zhang, Z. Zhang, H. Wang, Z. Li, Y. Qiao, and W. Liu, "Detecting faces using inside cascaded contextual cnn," in *ICCV*, 2017.
- [85] Y. Chen, L. Song, and R. He, "Masquer hunter: Adversarial occlusion-aware face detection," *CoRR*, 2017.