

Improved Selective Refinement Network for Face Detection

Shifeng Zhang^{1*}, Rui Zhu^{2*}, Xiaobo Wang^{2*}, Hailin Shi^{2†}, Tianyu Fu², Shuo Wang², Tao Mei², Stan Z. Li¹
¹ CBSR & NLPR, Institute of Automation, Chinese Academy of Sciences, Beijing, China.
² JD AI Research, Beijing, China.

{shifeng.zhang, szli}@nlpr.ia.ac.cn

{zhurui10, wangxiaobo8, shihailin, futianyu, wangshuo30, tmei}@jd.com

Abstract

As a long-standing problem in computer vision, face detection has attracted much attention in recent decades for its practical applications. With the availability of face detection benchmark WIDER FACE dataset, much of the progresses have been made by various algorithms in recent years. Among them, the Selective Refinement Network (SRN) face detector introduces the two-step classification and regression operations selectively into an anchor-based face detector to reduce false positives and improve location accuracy simultaneously. Moreover, it designs a receptive field enhancement block to provide more diverse receptive field. In this report, to further improve the performance of SRN, we exploit some existing techniques via extensive experiments, including new data augmentation strategy, improved backbone network, MS COCO pretraining, decoupled classification module, segmentation branch and Squeeze-and-Excitation block. Some of these techniques bring performance improvements, while few of them do not well adapt to our baseline. As a consequence, we present an improved SRN face detector by combining these useful techniques together and obtain the best performance on widely used face detection benchmark WIDER FACE dataset.

1. Introduction

Face detection is the primary procedure for other face-related tasks including face alignment, face recognition, face animation, face attribute analysis and human computer interaction, to name a few. The accuracy of face detection systems has a direct impact on these tasks, hence the success of face detection is of crucial importance. Given an arbitrary image, the goal of face detection is to determine whether there are any faces in the image, and

if present, return the image location and extent of each face. In recent years, great progress has been made on face detection [15, 26, 17, 44, 8, 31, 29, 1, 51, 49] due to the development of deep convolutional neural network (CNNs) [30, 10, 32, 13] and the collection of WIDER FACE benchmark dataset [41]. This challenging dataset has a high degree of variability in scale, pose and occlusion as well as plenty of tiny faces in various complex scenes, motivating a number of robust CNN-based algorithms.

We first give a brief introduction to these algorithms on the WIDER FACE dataset as follows. ACF [39] borrows the concept of channel features to the face detection domain. Faceness [40] formulates face detection as scoring facial parts responses to detect faces under severe occlusion. MTCNN [47] proposes a joint face detection and alignment method using unified cascaded CNNs for multi-task learning. CMS-RCNN [55] integrates contextual reasoning into the Faster R-CNN algorithm to help reduce the overall detection errors. LDCF+ [24] utilizes the boosted decision tree classifier to detect faces. The face detection model for finding tiny faces [12] trains separate detectors for different scales. Face R-CNN [35] and Face R-FCN [37] apply Faster R-CNN [27] and R-FCN [6] in face detection and achieve promising results. ScaleFace [42] detects different scales of faces with a specialized set of deep convolutional networks with different structures. SSH [22] adds large filters on each prediction head to merge the context information. SFD [50] compensates anchors for small faces with a few strategies in SSD [21] framework. MSCNN [2] performs detection at multiple output layers so as to let receptive fields match objects of different scales. Based on RetinaNet [19], FAN [36] proposes an attention mechanism at anchor level to detect the occluded faces. Zhu *et al.* [54] propose an Expected Max Overlapping score to evaluate the quality of anchor matching. PyramidBox [33] takes advantage of the information around human faces to improve detection performance. FNet [45] employs several training and testing techniques to Faster R-CNN to perform face detection. Inspired by RefineDet [48], SRN [5] appends another binary

*These authors contributed equally to this work.

†Corresponding author

This work was mainly done at JD AI Research and supported by JD-Grapevine Plan.

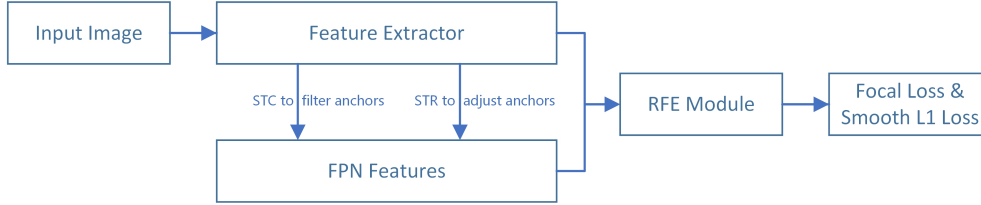


Figure 1. The brief overview of Selective Refinement Network. It consists of Selective Two-step Classification (STC), Selective Two-step Regression (STR) and Receptive Field Enhancement (RFE).

classification and regression stage in RetinaNet, in order to filter out most of simple negative anchors in the large feature maps and coarsely adjust the locations of anchors in the high level feature maps. FANet [46] aggregates higher-level features to augment lower-level features at marginal extra computation cost. DSFD [16] strengthens the representation ability by a feature enhance module. DFS [34] introduces a more effective feature fusion pyramid and a more efficient segmentation branch to handle hard faces. VIMFD [52] combines many previous techniques on SRN and achieves the state-of-the-art performance.

In this report, we exploit some existing techniques from classification and detection tasks to further improve the performance of SRN, including data augmentation strategy, improved backbone network, MS COCO pretraining, decoupled classification module, segmentation branch and SE block. By conducting extensive experiments, we share some useful techniques that make SRN regain the state-of-the-art performance on WIDER FACE. Meanwhile, we list some techniques that do not work well in our model, probably because (1) we have a strong baseline that causes them to not work well, (2) combination of ideas is not trivial, (3) they are not robust enough for universality, and (4) our implementation is wrong. This does not mean that they are not applicable to other models or other datasets.

2. Review of Baseline

In this section, we present a simple review of our baseline Selective Refinement Network (SRN). As illustrated in Figure 1, it consists of the Selective Two-step Classification (STC), Selective Two-step Regression (STR) and Receptive Field Enhancement (RFE). These three module are elaborated as follows.

2.1. Selective Two-step Classification

For one-stage detectors, numerous anchors with extreme positive/negative sample ratio (*e.g.*, there are about $300k$ anchors and the positive/negative ratio is approximately 0.006% in SRN) leads to quite a few false positives. Hence it needs another stage like RPN to filter out some negative examples. Selective Two-step Classification, inherited from RefineDet, effectively rejects lots of negative anchors and alleviates the class imbalance problem.

Specifically, most of anchors (*i.e.*, 88.9%) are tiled on the

first three low level feature maps, which do not contain adequate context information. So it is necessary to apply STC on these three low level features. Other three high level feature maps only produce 11.1% anchors with abundant semantic information, which is not suitable for STC. To sum up, the application of STC on three low level features brings advanced results, while on three high level ones will bring ineffective results and more computational cost. STC module suppresses the amount of negative anchors by a large margin, leading the positive/negative sample ratio about 38 times increased (*i.e.*, from around 1:15441 to 1:404). The shared classification convolution module and the same binary Focal Loss are used in the two-step classification, since both of the targets are distinguishing the faces from the background.

2.2. Selective Two-step Regression

Multi-step regression like Cascade RCNN [3] can improve the accuracy of bounding box locations, especially in some challenging scenes, *e.g.*, MS COCO-style evaluation metrics. However, applying multi-step regression to the face detection task without careful consideration may hurt the detection results.

For SRN, the numerous small anchors from three low level feature maps will cause the loss to bias towards regression problem and hinder the essential classification problem. Meanwhile, the feature representations of three lower pyramid levels for small faces are coarse, leading to the obstacle to perform two-step regression. These concerns will not happen while performing two-step regression on the three high level features, whose detailed features of large faces with large anchor scales help regress to more accurate locations. In summary, Selective Two-step Classification and Regression is a specific and efficient variant of RefineDet on face detection task, especially for small faces and some false positives.

2.3. Receptive Field Enhancement

Current networks usually possess square receptive fields, which affect the detection of objects with different aspect ratios. To address this issue, SRN designs a Receptive Field Enhancement (RFE) to diversify the receptive field of features before predicting classes and locations, which helps to capture faces well in some extreme poses.

3. Description of Improvement

Here we share some existing techniques that make SRN regain the state-of-the-art performance on the WIDER FACE dataset, including data augmentation, feature extractor and training strategy.

3.1. Data Augmentation

We use the original data augmentation strategies of SRN including photometric distortions, randomly expanding by zero-padding operation, randomly cropping patches from images and resizing patches to 1024×1024 . Additionally, with probability of $1/2$, we utilize the data-anchor-sampling in PyramidBox [33], which randomly selects a face in an image and crops sub-image based anchor. These data augmentation methods are crucial to prevent over-fitting and construct a robust model.

3.2. Feature Extractor

The greatest challenge in WIDER FACE is to accurately detect plenty of tiny faces. We believe that the ResNet-50-FPN [18] backbone of SRN still remains considerable room to improve the accuracy, especially for the tiny faces. Root-ResNet from ScratchDet [56] aims to improve the detection performance of small object, but its training speed is much slower than ResNet. To balance training efficiency and detection accuracy, we improve the ResNet-50 by taking the advantages of Root-ResNet and DRN [43].

Specifically, the downsampling operation (stride=2) to the image in the first 7×7 convolution layer of ResNet will cause the loss of important information, especially for small faces. After considering the motivation of Root-ResNet and DRN, we change the first conv layer’s stride from 2 to 1 and channel number from 64 to 16, as well as add two residual blocks (see Figure 2). One residual block is for enriching representational information while the other is for down-sampling, whose channel number are reduced to 16 and 32 to balance the parameters. This configuration can keep essential information of small faces without additional overhead.

3.3. Training Strategy

Because our ResNet-50-FPN backbone have been modified, we can not use the ImageNet pretrained model. One solution is like DRN that trains the modified backbone on ImageNet dataset [28] and then finetunes on WIDER FACE. However, He *et al.* [9] and ScratchDet have proved that the ImageNet pretraining is not necessary. Thus, we double the training epoch to 260 epochs and train the model with modified backbone from scratch. One of the key factor to train from scratch is the normalization. Due to the large input size (*i.e.*, 1024×1024), one 24G GPU only can be input up to 5 images, causing Batch Normalization [14] to not work

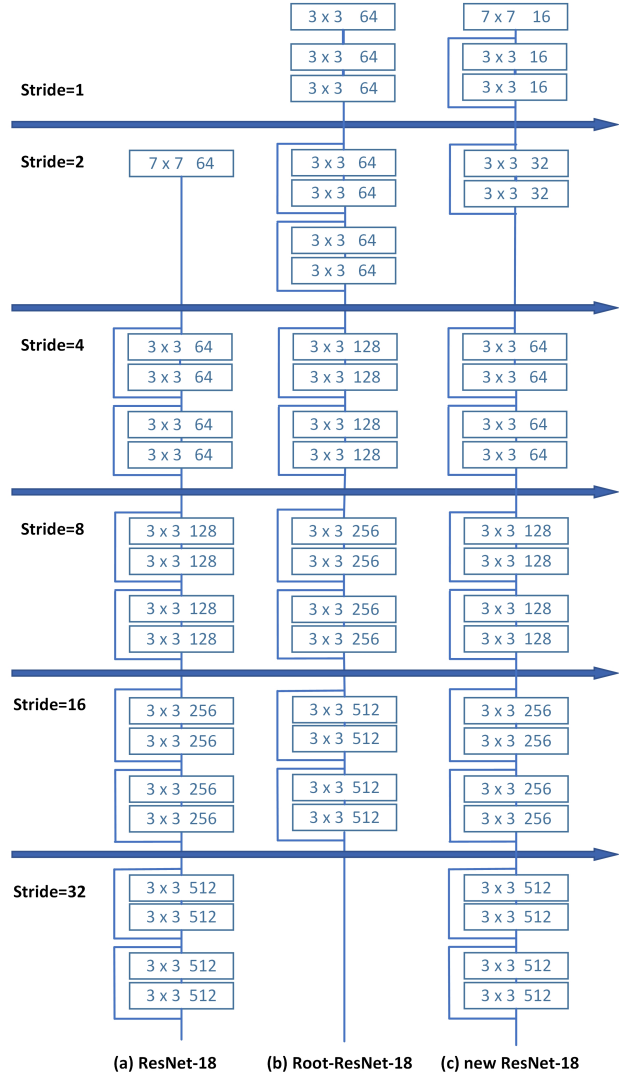


Figure 2. Network structure illustration. (a) ResNet-18: original structure. (b) Root-ResNet-18: replacing the 7×7 conv layer with three stacked 3×3 conv layers and changing the stride 2 to 1. (c) New-ResNet-18: combining DRN with Root-ResNet-18 to have a training speed/accuracy trade-off backbone for SRN.

well during training from scratch. To this end, we utilize Group Normalization [38] with group=16 to train this modified ResNet-50 backbone from scratch.

Besides, recent work FA-RPN [23] demonstrates that pretraining the model on the MS COCO dataset [20] is helpful to improve the performance of face detector on the WIDER FACE dataset. We attribute this promotion to a number of examples from people category and the objects with similar small scale (*i.e.*, ground truth area < 32) in the MS COCO dataset. So we also apply this pretraining strategy.

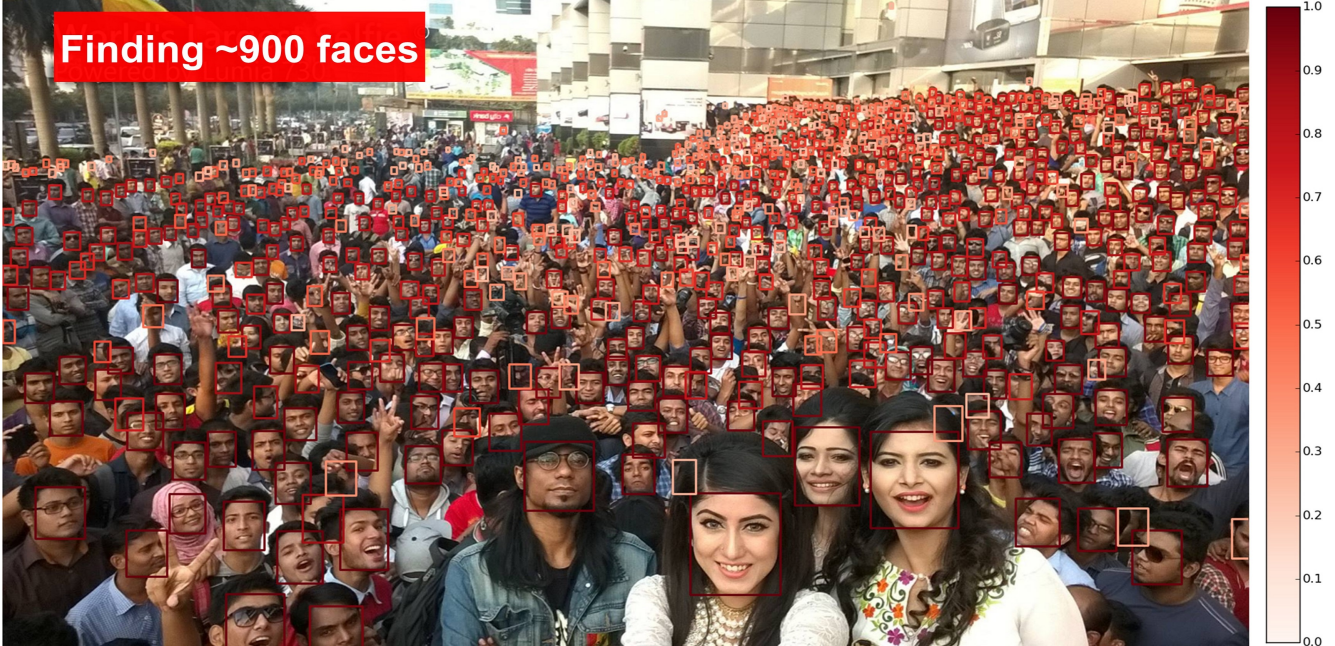


Figure 3. A qualitative result. Our detector successfully finds about 900 faces out of the reported 1000 faces in the above image. The confidences of the detections are presented in the color bar on the right hand. Best view in color.

3.4. Implementation Detail

Anchor Setting and Matching. Two anchor scales (*i.e.*, $2S$ and $2\sqrt{2}S$, where S represents the total stride size at each pyramid level) and one aspect ratios (*i.e.*, 1.25) cover the input images (*i.e.*, 1024×1024), with the anchor scale ranging from 8 to 362 pixels across pyramid levels. We assign anchors with $\text{IOU} > \theta_p$ as positive, anchors with IOU in $[0, \theta_n)$ as negative and others as ignored examples. Empirically, we set $\theta_n = 0.3$ and $\theta_p = 0.7$ for the first step, and $\theta_n = 0.4$ and $\theta_p = 0.5$ for the second step.

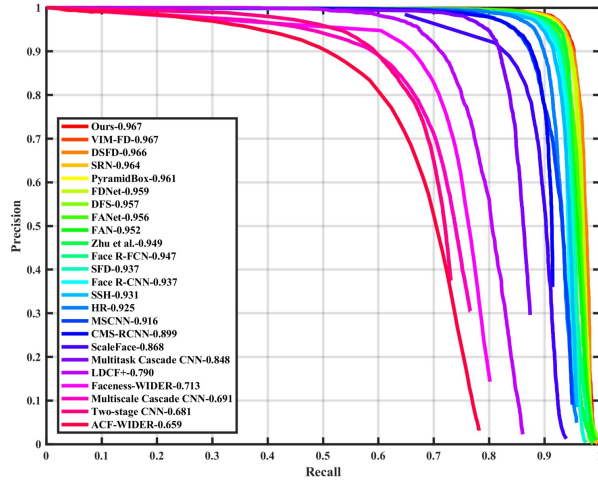
Optimization. At the training process, we simply sum the STC loss and the STR loss. We pretrain the new-designed backbone network with GroupNorm on MS COCO and finetune on WIDER FACE training set using SGD with 0.9 momentum, 0.0001 weight decay, and batch size 32. After 5 epochs warming up, the learning rate is set to 10^{-2} for the first 230 epochs, and decayed to 10^{-3} and 10^{-4} for another 20 and 10 epochs, respectively. Our method is implemented with the PyTorch library [25].

Inference. During the inference phase, the STC first filters the anchors on the first three feature maps with the positive confidence scores smaller than the threshold $\theta = 0.01$, and then the STR adjusts the anchors on the last three feature maps. The second step keeps top 2000 high detections among these refined anchors. Finally, we apply the non-maximum suppression (NMS) with jaccard overlap of 0.4 to generate 750 high confident results per image. The multi-scale testing strategy is used during the inference phase.

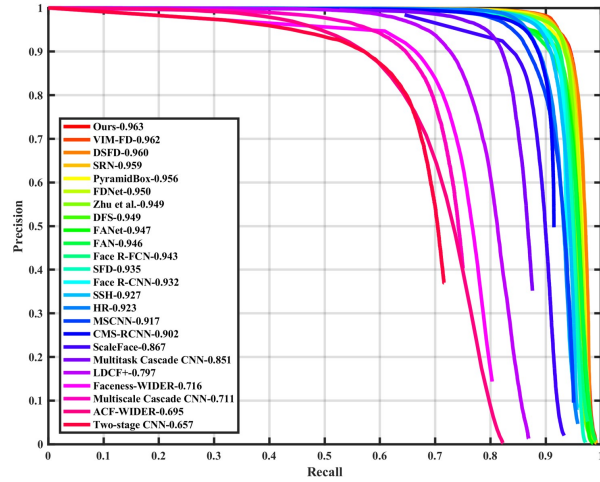
4. Result on WIDER FACE

The WIDER FACE dataset contains 32,203 images and 393,703 annotated faces bounding boxes including high degree of variability in scale, pose, facial expression, occlusion and lighting condition. It is split into the training (40%), validation (10%) and testing (50%) subsets by randomly sampling from each scene category (totally 61 classes), and defines three levels of difficulty: Easy, Medium, Hard, based on the detection rate of Edge-Box [57]. Following the evaluation protocol in WIDER FACE, we only train the model on the training set and test on both the validation and testing sets. To obtain the evaluation results on the testing set, we submit the detection results to the authors for evaluation.

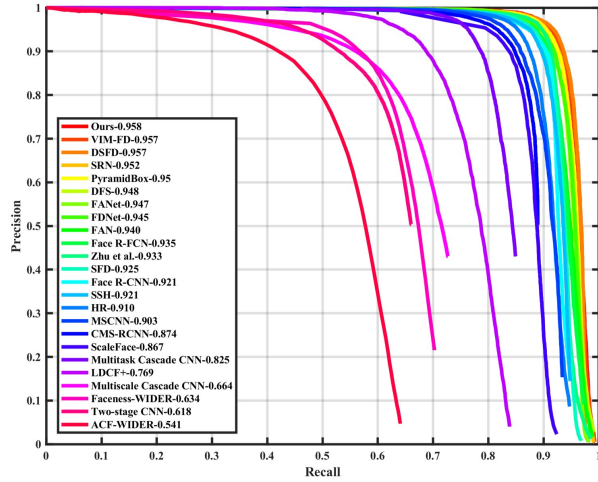
As shown in Figure 4, we compare our method (namely ISRN) with 23 state-of-the-art face detection methods [41, 39, 40, 47, 55, 24, 12, 35, 37, 42, 22, 50, 2, 36, 54, 33, 45, 5, 46, 16, 34, 52]. We find that our model achieves the state-of-the-art performance based on the average precision (AP) across the three evaluation metrics, especially on the Hard subset which contains a large amount of tiny faces. Specifically, it produces the best AP scores in all subsets of both validation and testing sets, *i.e.*, 96.7% (Easy), 95.8% (Medium) and 90.9% (Hard) for validation set, and 96.3% (Easy), 95.4% (Medium) and 90.3% (Hard) for testing set, surpassing all approaches, which demonstrates the superiority of our face detector. We show one qualitative result of the World Largest Selfie in Figure 3. Our detector successfully finds about 900 faces out of the reported 1,000 faces.



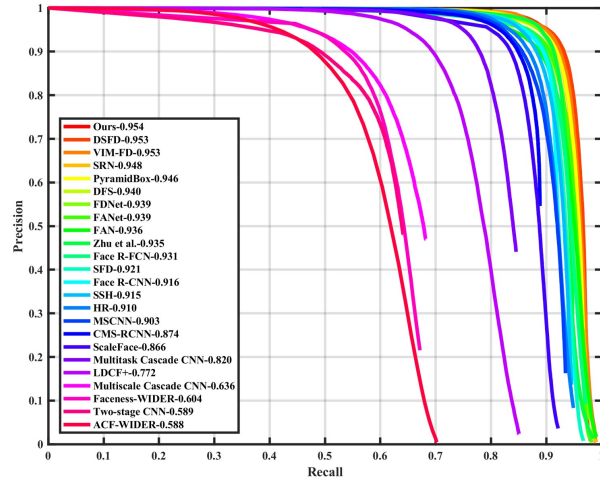
(a) Val: Easy



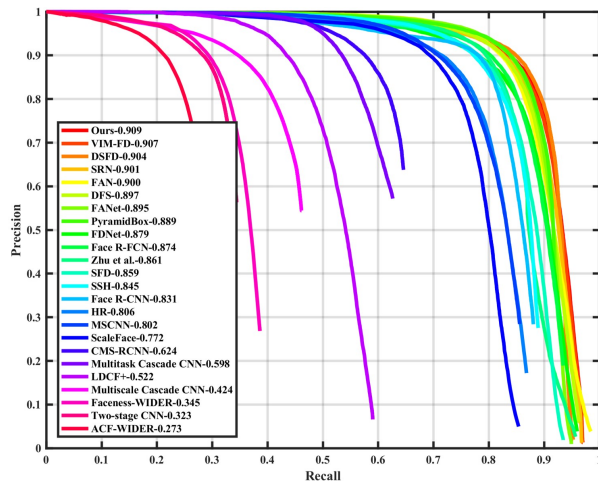
(b) Test: Easy



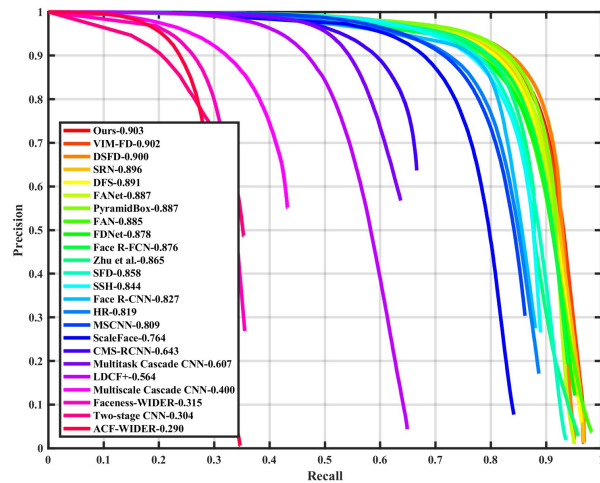
(c) Val: Medium



(d) Test: Medium



(e) Val: Hard



(f) Test: Hard

Figure 4. Precision-recall curves on WIDER FACE validation and testing subsets.

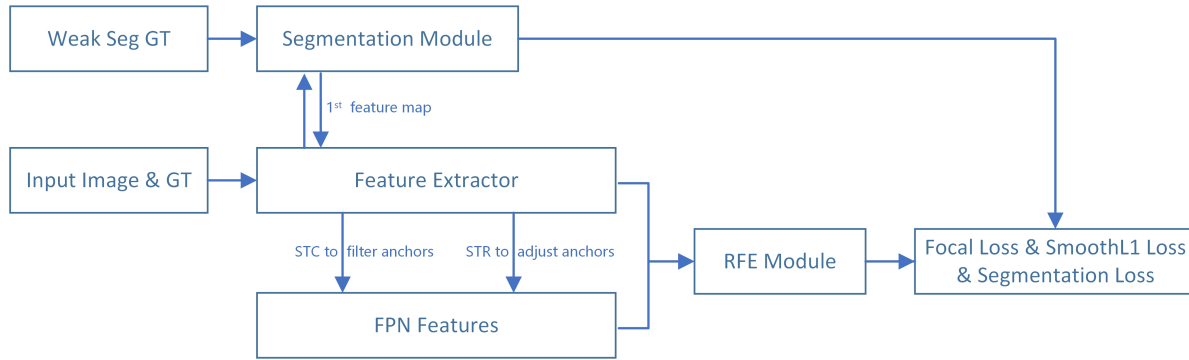


Figure 5. The brief overview of Selective Refinement Network with segmentation branch.

5. Things We Tried That Did Not Work Well

This section lists some techniques that do not work well in our model, probably because (1) we have a strong baseline, (2) combination of ideas is not trivial, (3) they are not robust enough for universality, and (4) our implementation is wrong. This does not mean that they are not applicable to other models or other datasets.

Decoupled Classification Refinement (DCR) [4]. It is an extra-stage classifier for the classification refinement. During the training process, DCR samples hard false positives with high confidence scores from the base Faster R-CNN detector, and then trains a stronger classifier. At the inference time, it simply multiplies the score from the base detector and another score from DCR to rerank detection results. Faster R-CNN with DCR gets great improvement on MS COCO [20] and PASCAL VOC [7] datasets. Therefore, we try to use DCR to suppress the false positives at the beginning and conduct some inquiring experiments based on our SRN baseline. However, with the help of RPN proposals and ROIs, the sampling strategies of DCR for two-stage detectors are much easier to design than the one for one-stage methods. SRN face detector produces too much boxes so we try a lot of sampling heuristics. Besides, we attempt some different crop size of the training examples due to the large scale variance and numerous small faces on WIDER FACE. Considering the scale of training set (positive and negative examples cropped from WIDER FACE training set) and network overfitting, we also try DCR backbones with different order of magnitude. With the setting of crop size=20, DRN-22 backbone and sampling strategy (positive examples: $0.5 < \text{IOU} < 0.8$ and negative examples: $\text{IOU} < 0.3$), our best result is also slight lower than our baseline detector. Further experiments about DCR need to be conducted for face detection task on WIDER FACE.

Segmentation Branch [53]. A segmentation branch is added on SSD in DES, which is applied to the low level feature map and supervised by weak bounding-box level segmentation ground truth. The low level feature map is reweighed by the output $H \times W \times 1$ map of segmentation

branch. This enhancement can be regarded as a element-wise attention mechanism. As shown in Figure5, we apply segmentation branch to the first low level feature map of SRN but the final results drop a bit on three metrics.

Squeeze-and-Excitation (SE) Block [11]. It adaptively reweighs channel-wise features by using global information to selectively emphasise informative features and suppress useless ones. It can be regarded as a channel-wise attention mechanism with a squeeze-and-excitation $1 \times 1 \times C$ feature map, and the original feature will be reweighed to generate more representational one. As shown in Figure6, we apply SE block to the final detection feature map of SRN, but the final results drop 0.2%, 0.2% and 0.4% respectively on Easy, Medium and Hard metric.

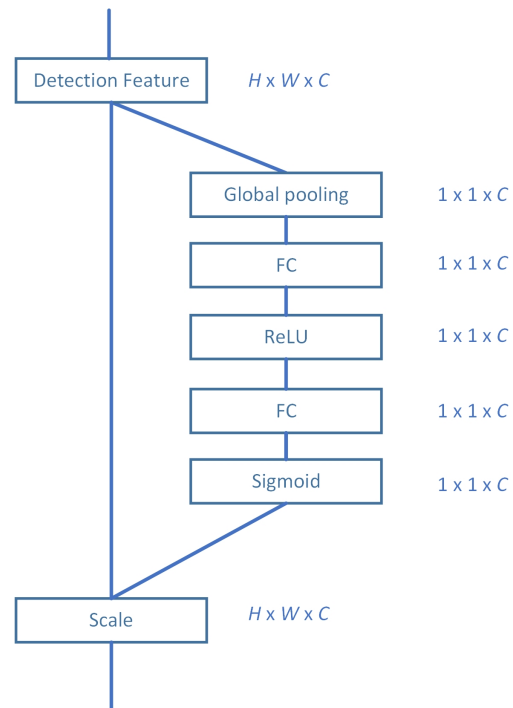


Figure 6. Applying SE block to the detection feature.

6. Conclusion

To further boost the performance of SRN, we exploit some existing techniques including new data augmentation strategy, improved backbone network, MS COCO pretraining, decoupled classification module, segmentation branch and SE block. By conducting extensive experiments on the WIDER FACE dataset, we find that some of these techniques bring performance improvements, while few of them do not well adapt to our baseline. By combining these useful techniques together, we present an improved SRN detector and obtain the state-of-the-art performance on the widely used face detection benchmark WIDER FACE dataset.

References

- [1] Y. Bai, Y. Zhang, M. Ding, and B. Ghanem. Finding tiny faces in the wild with generative adversarial network. In *CVPR*, 2018.
- [2] 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.
- [3] Z. Cai and N. Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *CVPR*, 2018.
- [4] B. Cheng, Y. Wei, H. Shi, R. Feris, J. Xiong, and T. Huang. Revisiting rnn: On awakening the classification power of faster rnn. In *ECCV*, 2018.
- [5] 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.
- [6] J. Dai, Y. Li, K. He, and J. Sun. R-fcn: Object detection via region-based fully convolutional networks. In *NIPS*, 2016.
- [7] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 2010.
- [8] Z. Hao, Y. Liu, H. Qin, J. Yan, X. Li, and X. Hu. Scale-aware face detection. In *CVPR*, 2017.
- [9] K. He, R. Girshick, and P. Dollár. Rethinking imagenet pre-training. *arXiv preprint arXiv:1811.08883*, 2018.
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [11] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. *arXiv preprint arXiv:1709.01507*, 7, 2017.
- [12] P. Hu and D. Ramanan. Finding tiny faces. In *CVPR*, 2017.
- [13] G. Huang, Z. Liu, K. Q. Weinberger, and L. van der Maaten. Densely connected convolutional networks. In *CVPR*, 2017.
- [14] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- [15] H. Li, Z. Lin, X. Shen, J. Brandt, and G. Hua. A convolutional neural network cascade for face detection. In *CVPR*, 2015.
- [16] J. Li, Y. Wang, C. Wang, Y. Tai, J. Qian, J. Yang, C. Wang, J. Li, and F. Huang. Dsf: Dual shot face detector. *arXiv preprint arXiv:1810.10220*, 2018.
- [17] 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.
- [18] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017.
- [19] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar. Focal loss for dense object detection. In *ICCV*, 2017.
- [20] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.
- [21] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *ECCV*, 2016.
- [22] M. Najibi, P. Samangouei, R. Chellappa, and L. S. Davis. Ssh: Single stage headless face detector. In *ICCV*, 2017.
- [23] M. Najibi, B. Singh, and L. S. Davis. Fa-rpn: Floating region proposals for face detection. *arXiv preprint arXiv:1812.05586*, 2018.
- [24] 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.
- [25] A. Paszke, S. Gross, S. Chintala, and G. Chanan. Pytorch, 2017.
- [26] H. Qin, J. Yan, X. Li, and X. Hu. Joint training of cascaded CNN for face detection. In *CVPR*, 2016.
- [27] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015.
- [28] 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. *IJCV*, 2015.
- [29] X. Shi, S. Shan, M. Kan, S. Wu, and X. Chen. Real-time rotation-invariant face detection with progressive calibration networks. In *CVPR*, 2018.
- [30] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [31] 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.
- [32] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015.
- [33] X. Tang, D. K. Du, Z. He, and J. Liu. Pyramidbox: A context-assisted single shot face detector. *arXiv preprint arXiv:1803.07737*, 2018.
- [34] 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. *arXiv preprint arXiv:1811.08557*, 2018.
- [35] H. Wang, Z. Li, X. Ji, and Y. Wang. Face r-cnn. *arXiv preprint arXiv:1706.01061*, 7, 2017.
- [36] J. Wang, Y. Yuan, and G. Yu. Face attention network: An effective face detector for the occluded faces. *arXiv preprint arXiv:1711.07246*, 2017.

- [37] Y. Wang, X. Ji, Z. Zhou, H. Wang, and Z. Li. Detecting faces using region-based fully convolutional networks. *arXiv preprint arXiv:1709.05256*, 2017.
- [38] Y. Wu and K. He. Group normalization. In *ECCV*, 2018.
- [39] B. Yang, J. Yan, Z. Lei, and S. Z. Li. Aggregate channel features for multi-view face detection. In *IJCB*, 2014.
- [40] S. Yang, P. Luo, C.-C. Loy, and X. Tang. From facial parts responses to face detection: A deep learning approach. In *ICCV*, 2015.
- [41] S. Yang, P. Luo, C.-C. Loy, and X. Tang. Wider face: A face detection benchmark. In *CVPR*, 2016.
- [42] S. Yang, Y. Xiong, C. C. Loy, and X. Tang. Face detection through scale-friendly deep convolutional networks. *arXiv preprint arXiv:1706.02863*, 2017.
- [43] F. Yu, V. Koltun, and T. Funkhouser. Dilated residual networks. In *CVPR*, 2017.
- [44] J. Yu, Y. Jiang, Z. Wang, Z. Cao, and T. S. Huang. Unitbox: An advanced object detection network. In *ACMMM*, 2016.
- [45] C. Zhang, X. Xu, and D. Tu. Face detection using improved faster rcnn. *arXiv preprint arXiv:1802.02142*, 2018.
- [46] J. Zhang, X. Wu, J. Zhu, and S. C. Hoi. Feature agglomeration networks for single stage face detection. *arXiv preprint arXiv:1712.00721*, 2017.
- [47] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *SPL*, 2016.
- [48] S. Zhang, L. Wen, X. Bian, Z. Lei, and S. Z. Li. Single-shot refinement neural network for object detection. In *CVPR*, 2018.
- [49] 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.
- [50] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Z. Li. S³FD: Single shot scale-invariant face detector. In *ICCV*, 2017.
- [51] S. Zhang, X. Zhu, Z. Lei, X. Wang, H. Shi, and S. Z. Li. Detecting face with densely connected face proposal network. *Neurocomputing*, 2018.
- [52] Y. Zhang, X. Xu, and X. Liu. Robust and high performance face detector. *arXiv preprint arXiv:1901.02350*, 2019.
- [53] Z. Zhang, S. Qiao, C. Xie, W. Shen, B. Wang, and A. L. Yuille. Single-shot object detection with enriched semantics. In *CVPR*, 2018.
- [54] C. Zhu, R. Tao, K. Luu, and M. Savvides. Seeing small faces from robust anchors perspective. In *CVPR*, 2018.
- [55] C. Zhu, Y. Zheng, K. Luu, and M. Savvides. Cms-rcnn: contextual multi-scale region-based cnn for unconstrained face detection. *arXiv preprint arXiv:1606.05413*, 2016.
- [56] R. Zhu, S. Zhang, X. Wang, L. Wen, H. Shi, L. Bo, and T. Mei. Scratchdet: Exploring to train single-shot object detectors from scratch. *arXiv preprint arXiv:1810.08425*, 2018.
- [57] C. L. Zitnick and P. Dollár. Edge boxes: Locating object proposals from edges. In *ECCV*, 2014.