

ACFD: Asymmetric Cartoon Face Detector

Bin Zhang^{12*}, Jian Li^{1*}, Yabiao Wang¹, Zhipeng Cui¹,
Yili Xia^{2†}, Chengjie Wang¹, Jilin Li¹, Feiyue Huang¹

¹Youtu Lab, Tencent

²School of Information Science and Engineering, Southeast University

{z-bingo, yili_xia}@seu.edu.cn

{swordli, caseywang, zhipengcui, jasoncjwang, jeronlinli, garyhuang}.tencent.com

Abstract

Cartoon face detection is a more challenging task than human face detection due to many difficult scenarios is involved. Aiming at the characteristics of cartoon faces, such as huge differences within the intra-faces, in this paper, we propose an asymmetric cartoon face detector, named ACFD. Specifically, it consists of the following modules: a novel backbone VoVNetV3 comprised of several asymmetric one-shot aggregation modules (AOSA), asymmetric bi-directional feature pyramid network (ABi-FPN), dynamic anchor match strategy (DAM) and the corresponding margin binary classification loss (MBC). In particular, to generate features with diverse receptive fields, multi-scale pyramid features are extracted by VoVNetV3, and then fused and enhanced simultaneously by ABi-FPN for handling the faces in some extreme poses and have disparate aspect ratios. Besides, DAM is used to match enough high-quality anchors for each face, and MBC is for the strong power of discrimination. With the effectiveness of these modules, our ACFD achieves the 1st place on the detection track of 2020 iCartoon Face Challenge under the constraints of model size 200MB, inference time 50ms per image, and without any pretrained models.

1 Introduction

Face detection is a long-standing and essential task for many downstream applications, such as face alignment, face recognition and face tracking. In order to detect face effectively and efficiently, many detection pipelines [Zhang *et al.*, 2017; Tang *et al.*, 2018; Chi *et al.*, 2019; Li *et al.*, 2019a; Zhang *et al.*, 2020b] have been proposed and achieved better and better performance on the challenging face detection benchmarks, such as WIDER Face [Yang *et al.*, 2016]. In particular, on the most challenging dataset WIDER Face, the average precision (AP) has been improved from the original 40% to more

*Equal contribution. This work was down when Bin Zhang was an intern at Tencent Youtu Lab.

†Contact Author

than 90%, and there are still many researchers are working on solving the remaining difficulties.

With the development of multimedia technology, more and more cartoon faces appear in the animations, comics, and even appear on the social media platforms as the personal profile. Therefore, the excellent detector for cartoon faces is also particularly important. Although the recent face detectors could handle most of the human faces, they cannot detect the cartoon faces accurately with many missing and false-positive results. That is because many scenarios in cartoons are more difficult than the real world, for instance, different faces may perform completely different characteristics, many faces are very similar to the negative samples (the bodies, background, etc.) and a considerable part of faces have the disparate aspect ratio of more than 3 or less than 1/3, which is almost not involved in human faces.

In this way, a large-scale and challenging cartoon person dataset is beneficial for designing an effective and robust deep learning based-approach. Several works [Gong *et al.*, 2020; Jha *et al.*, 2018; Wu *et al.*, 2019] transfer the human faces to cartoons directly by utilizing generative adversarial networks due to the lack of cartoon face dataset, but there is still a big gap. Recently, a manually annotated cartoon face dataset is proposed by iQIYI [Li *et al.*, 2019b] that contains 50000 images for detection tasks and more than 380000 images for face recognition, promoting the development of the community.

To solve the aforementioned difficulties in cartoon faces, we propose a novel asymmetric cartoon face detector (ACFD) with four improvements to enhance the diversity of features, regression, and classification ability of networks according to the characteristics of cartoon faces. Specifically, to provide features with more diverse receptive fields, we propose a more effective backbone network VoVNetV3 based on VoVNet [Lee *et al.*, 2019] and VoVNetV2 [Lee and Park, 2020] with better performance than ResNet [He *et al.*, 2016] and DenseNet [Huang *et al.*, 2017], through which the multi-scale features can be extracted and then fused and enhanced by the following ABi-FPN. For enhancing the regression and classification capacity, we utilize DAM to compensate high-quality anchors with strong regression ability providing better initialization for the regressor, furthermore, MBC is used to better distinguish faces from the dense predictions. The main contributions of the paper can be summarized as follow.

- Proposing a novel backbone network VoVNetV3 to ob-

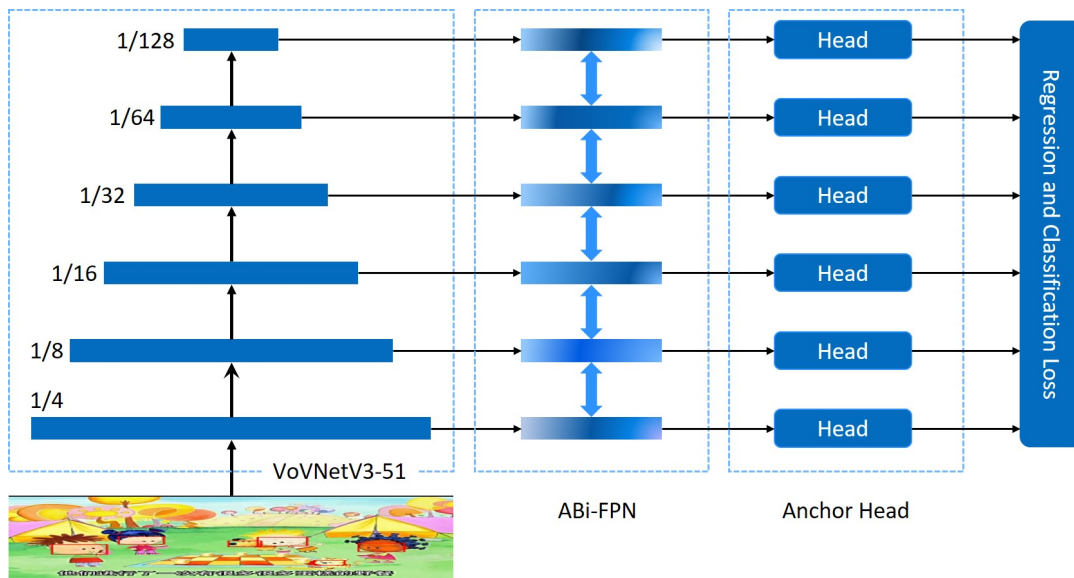


Figure 1: Pipeline of our ACFD, which is an anchor-based one-stage detector, where feature maps of backbone network with strides from 4 to 128 are feed into ABi-FPN and used for predicting.

tain features with more diverse receptive fields.

- Designing an ABi-FPN to fuse multi-scale features and enhance the semantic information simultaneously.
- Presenting a DAM strategy to match high-quality anchors for each faces dynamically.
- Introducing a MBC module to distinguish faces within the dense predictions.
- Achieving the 1st place on the detection track of 2020 iCartoon Face Challenge.

2 Related Work

Face Detection. Face detectors usually inherit the basic settings of generic object detection. Review the recent state-of-the-art face detectors [Zhang *et al.*, 2017; Tang *et al.*, 2018; Chi *et al.*, 2019; Li *et al.*, 2019a; Zhang *et al.*, 2020b], all of them are anchor-based one-stage detectors and consist of backbone, neck and detection head. In detail, the backbone network designed for image classification [He *et al.*, 2016; Simonyan and Zisserman, 2014] is used to extract multi-scale features, and the neck network is utilized to fuse and enhance the multi-scale features sequentially. We propose backbone network VoVNetV3 and ABi-FPN to extract multi-scale features with strong representation, aggregate, and enhance the context respectively.

Anchor Match. Traditionally, an anchor would be assigned as positive if its Intersection over Union (IoU) with a ground-truth bounding box is larger than a predefined threshold [Ren *et al.*, 2015]. To suit face detection, [Zhang *et al.*, 2017] proposes an improved strategy helping the outer faces to match enough anchors. [Chi *et al.*, 2019] and [Zhang *et al.*, 2020b] introduce selective two-step regression and classification to refine and filter the predefined anchors. [Kong *et al.*, 2019] and [Liu *et al.*, 2019] observe that taken regressed bounding boxes into match phase would bring considerable perfor-

mance improvements. Inspired by these works, we propose a dynamic anchor match strategy to compensate enough high-quality anchors for the outer faces.

Loss Design. Smooth-L1 and cross-entropy losses [Ren *et al.*, 2015] are widely used in object detection due to the simplicity and effectiveness. Nowadays, IoU-based losses [Rezatofighi *et al.*, 2019; Zheng *et al.*, 2020] are proposed to strengthen the connections between regression and classification tasks, focal loss [Lin *et al.*, 2017b] is proposed to alleviate the extreme foreground-background class imbalance. Unlike these methods, we transfer the margin loss [Wang *et al.*, 2018; Deng *et al.*, 2019] in face recognition to face detection for further improving the power of discrimination due to many cartoon faces are very similar to the background.

3 Methodology

In this section, we briefly introduce the pipeline of the proposed ACFD at first in Sec. 3.1, then detailly describe the proposed VoVNetV3 in Sec. 3.2, asymmetric bi-directional feature pyramid network in Sec. 3.3, dynamic anchor match in Sec. 3.4, and binary margin classification loss in Sec. 3.5, respectively.

3.1 Pipeline of ACFD

The pipeline of the proposed ACFD is indicated in Fig. 1. We adopt the proposed VoVNetV3-51 as the backbone network of ACFD, which consists of 6 stages to generate feature maps with strides from 4 to 128. Then, the multi-scale pyramid features extracted from the backbone network are feed into the proposed ABi-FPN for further aggregating and refining the context information. Finally, the dense predictions are obtained by the corresponding anchor head networks.

3.2 VoVNetV3 Backbone Network

VoVNet [Lee *et al.*, 2019] is a computation and energy efficient backbone network that can efficiently present diversi-

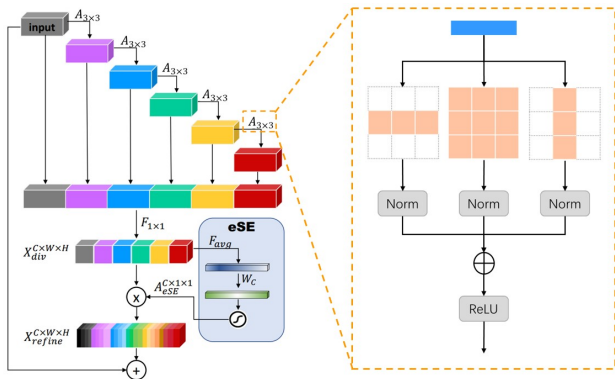


Figure 2: The architecture of AOSA module. $F_{1 \times 1}$ is convolution layer with kernel size 1, $A_{3 \times 3}$ denotes ACB with kernel size 3, F_{avg} is global average pooling, W_C is fully-connected layer, \oplus and \otimes indicate element-wise addition and multiplication.

fied feature representations owing to the utilization of one-shot aggregation (OSA) modules. For further boosting the performance, VoVNetV2 [Lee and Park, 2020] is proposed by adding the residual connection to address the limitation of VoVNet during optimization [He *et al.*, 2016], and employing an efficient squeeze-excitation (eSE) attention module for modeling the interdependency between the channel of feature maps to enhance its representation. In order to further enrich the diversity of feature maps, a more effective backbone network VoVNetV3 is proposed on the basis of VoVNet and VoVNetV2, which is comprised of several asymmetric one-shot aggregation (AOSA) modules, as presented in Fig. 2.

Asymmetric One-Shot Aggregation Module. Unlike DenseNet [Huang *et al.*, 2017], the OSA modules of VoVNet and VoVNetV2 generate feature maps in a relatively sparse-connected manner, where each feature map is connected to the subsequent convolution layer so as to produce the feature with a large receptive field, and concatenated with the final output feature map only once, receptively. As a result, OSA modules could generate features with rich receptive fields, however, both of them only process square receptive fields and would affect the detection of faces with different aspect ratios. Inspired by [Szegedy *et al.*, 2017; Ding *et al.*, 2019], the proposed AOSA overcomes the above limitation by replacing all 3×3 convolutions in OSA with asymmetric convolution blocks (ACB) as shown in Fig. 2. In this way, two additional convolutions with kernel 1×3 and 3×1 are added in parallel on 3×3 convolution layer to extract features with rectangle receptive fields.

We utilize VoVNetV3-51 with 6 stages of outputs in this paper, its configuration is presented at Table. 1.

3.3 ABi-FPN

At present, most of the face detectors utilize ResNet [He *et al.*, 2016] and VGG [Simonyan and Zisserman, 2014] to extract the multi-scale features, while both of them can only possess square receptive fields and are harmful to the faces with extreme aspect ratios. It is especially important in cartoon face detection due to about 10% of faces with ratios of larger than 2.0 or smaller than 0.5. As for solving the limitation of networks, recent state-of-the-arts face detectors [Tang

Layer	Output	Stride	Repeat	Channel
image	640×640	—	—	—
conv1	320×320	2	1	64
conv2	320×320	1	1	64
conv3	160×160	2	1	128
stage1	160×160	1	1	128/256
max-pool	80×80	2	1	256
stage2	80×80	1	1	160/512
max-pool	40×40	2	1	512
stage3	40×40	1	2	192/768
max-pool	20×20	2	1	768
stage4	20×20	1	2	224/1024
max-pool	10×10	2	1	1024
stage5	10×10	1	1	128/128
max-pool	5×5	2	1	128
stage6	5×5	1	1	128/128

Table 1: The configuration of proposed VoVNetV3-51, in which the channel in format $-/-$ means the splitted and concatenated channel in AOSA.



Figure 3: Illustration of our DAM, where red boxes mean ground-truth faces, blue boxes are anchors, and green boxes are the corresponding regressed bounding boxes. (a) Assigning anchors with IoU greater than the predefined threshold as positive samples, (b) assigning unmatched anchors in the first step with strong regression ability as positive.

et al., 2018; Li *et al.*, 2019a; Zhang *et al.*, 2020b] add an additional module followed by feature fusion modules to refine the receptive fields, which is effective but inefficient. Instead of adding the extra modules for processing, we propose an efficient and effective module named ABi-FPN to fuse multi-scale features, enrich semantic information, and refine the receptive fields of features simultaneously. Specifically, ACB is also employed to replace the convolution layers of BiFPN [Tan *et al.*, 2020], in which the receptive fields of features would be more diverse as the aggregation of multi-scale features.

3.4 Dynamic Anchor Match (DAM)

Recently, some works [Kong *et al.*, 2019; Liu *et al.*, 2019] observe an interesting phenomenon that some unmatched anchors have strong regression ability, as presented in Fig. 3 (b), the anchor with strong regression capacity could obtain a bounding box with the large IoU score, despite its own IoU is very small. Inspired by the observation, we propose a DAM strategy to make full use of these anchors with strong regres-

sion ability, as a result, it matches enough high-quality anchors for each ground-truth face. Firstly, the same as traditional match strategy [Liu *et al.*, 2016] shown as Fig. 3 (a), the anchors with IoU greater than a threshold are assigned as positive. Secondly, the anchors would be compensated as positive if IoUs of their related regressed bounding boxes are higher than another threshold. The details can be seen in Algorithm. 1.

Algorithm 1 Dynamic anchor match strategy.

Input: A, B, T_1, T_2, G, R, L

A a set of predefined anchors, B the related bounding boxes. T_1, T_2 IoU thresholds for first and second match steps. G ground-truth boxes.

R, L matched bounding boxes and labels of each anchor.

Output: R, L

```

1: for  $a_i, b_i$  in  $A, B$  do
2:   # first step
3:    $AnchorIoU_i \leftarrow IoU(a_i, G)$ 
4:    $GtIdx, AnchorMax \leftarrow argmax(AnchorIoU_i)$ 
5:   if  $AnchorMax \geq T_1$  then
6:      $R[i] \leftarrow G[GtIdx], L[i] \leftarrow 1$ 
7:   else
8:     # second step
9:      $BboxIoU_i \leftarrow IoU(b_i, G)$ 
10:     $GtIdx, BboxMax \leftarrow argmax(BboxIoU_i)$ 
11:    if  $BboxMax \geq T_2$  then
12:      # 2 means compensated anchor
13:       $R[i] \leftarrow G[GtIdx], L[i] \leftarrow 2$ 
14:    else
15:       $L[i] \leftarrow 0$ 
16:    end if
17:  end if
18: end for

```

3.5 Margin Binary Classification Loss (MBC)

As described in the above section, our DAM strategy could match enough high-quality anchors with strong regression ability for each face. However, these high-quality anchors are usually far from the ground-truth face and may dominate the loss during propagation. Therefore, we separately compute and weight the losses of matched anchor in the first step and compensated high-quality anchors in the second step, the regression and classification losses are reformatted as follow,

$$\begin{aligned} \ell_{reg} &= \frac{1}{N_1} \sum_{i \in \psi_1} \mathcal{L}_{smoothL1}(x_i, x_i^*) + \frac{\lambda_{reg}}{N_2} \sum_{i \in \psi_2} \mathcal{L}_{smoothL1}(x_i, x_i^*) \\ \ell_{cls} &= \frac{1}{N_1} \sum_{i \notin \psi_2} \mathcal{L}_{Focal}(p_i^m, p_i^*) + \frac{\lambda_{cls}}{N_2} \sum_{i \in \psi_2} \mathcal{L}_{Focal}(p_i^m, p_i^*) \end{aligned} \quad (2)$$

where N_1 and N_2 are numbers of matched anchors ψ_1 and compensated anchors ψ_2 , λ_{reg} and λ_{cls} are the corresponding weighted coefficients.

Furthermore, to improve the classification ability of our network so that it can distinguish faces who are similar to the background, we transfer the widely used margin-based loss function [Wang *et al.*, 2018; Deng *et al.*, 2019] in face

recognition to face detection. The margin-based losses share the same idea that maximizing the inter-class variance, minimizing the intra-class variance, and enhance the capacity of discrimination by adding an extra hard margin. In our margin binary classification application, suppose p_i is the output of the network, then the margin-based prediction is formulated as follow,

$$p_i^m = [p_i^o = 1] \cdot (p_i - m) + [p_i^o = 0] \cdot p_i \quad (3)$$

in which, p_i^m is the corresponding one-hot label, m is a hard margin, and p_i^o is used for the computation of classification loss.

4 Experiments

In this section, we firstly present the details of our implementation, then the effectiveness of our proposed ACFD is verified through a set of experiments, finally, we introduce the engineering tricks that are used for optimizing the time-consuming of our ACFD.

4.1 Implementation Details

Dataset. During ablation studies, we split 50000 images of iCartoon Face [Li *et al.*, 2019b] into 45000 images for training and 5000 for reporting the effectiveness. All images are used for training for the final submitted model.

Data Augmentation. To prevent over-fitting and augment the difficult faces during training such as small and blur faces, the faces too large to location accurately, and etc. A series of data augmentations are employed, summarized as follow: (1) color distort for training images, (b) expand the images with a random range $[1, 4]$ by mean-padding to augment the small faces, (c) crop the images with a random size at a random position to augment the big faces, (d) random tile the faces to anchor scales, finally, resize the images to 640×640 for feeding into the network.

Anchor Design. We associate only one anchor for each location of the detection layers with scale 4 and ratio 1:1. So, there are total of 34125 anchors per image covered faces of size 16–512 pixels within the training phase.

Optimization. Without ImageNet-pretrained model, we use "kaiming" method [He *et al.*, 2015] to initialize all the parameters. The SGD algorithm is applied to train the model with momentum 0.9, weight decay 5×10^{-4} , batch size 16×4 V100 GPUs, and warmup strategy to increase the learning rate from 1.0×10^{-6} to 0.04 at the first 1000 iterations. Then, for the final submit, the learning rate is divided by 10 at 200, 250, 280 epochs, and ended at 300 epoch; for the ablation studies, it is divided by 10 at 100, 150, 180 epochs and ended at 200 epoch.

Other Hyper-parameters. Empirically, we define two thresholds of dynamic anchor match as $T_1 = 0.35$ and $T_2 = 0.7$. Correspondingly, the weighted coefficients in regression and classification losses are $\lambda_{reg} = \lambda_{cls} = 0.7$, and the margin used in classification loss is 0.2.

Module	ResNet50	SE-ResNet50	Res2Net50	ResNeSt50	EfficientNet-B3	VoVNetV2-39	VoVNetV3-51
mAP (%)	90.18	90.23	89.59	89.97	88.63	90.37	90.74

Table 2: mAP (%) of ACFD with different backbone networks.

Inference. Multi-scale test is employed for the final submitting, we resize the images to three predefined scales that 480×645 , 640×860 and 800×1075 for predicting. During inference, the network outputs the top-1000 predictions with confidence scores higher than 0.08 for each scale, finally, we apply the non-maximum suppression (NMS) algorithm with IoU threshold 0.55 to generate top-100 high confidence detections as the final results.

4.2 Model Analysis

In this subsection, extensive experiments are conducted to demonstrate the effectiveness of our proposed modules in ACFD. For the fair comparison, we use the same parameter settings apart from the specific changes to the components. In order to better understand our ACFD, the experiments are carried out on the different baselines, by ablating different modules to perform how it affects the final performance.

VoVNetV3-51. Firstly, we compare our VoVNetV3-51 with the state-of-the-art and widely used backbone networks, e.g., ResNet50 [He *et al.*, 2016], SE-ResNet50 [Hu *et al.*, 2018], Res2Net50 [Gao *et al.*, 2019], ResNeSt50 [Zhang *et al.*, 2020a], EfficientNet-B3 [Tan and Le, 2019], VoVNetV2-39 [Lee and Park, 2020]. To better suit the framework of ACFD, two extra convolution layers with stride 2 followed by batch normalization and ReLU activation are added to these backbones for generating 6-level features with stride from 4 to 128. As presented in Table 2, ResNet50 and SE-ResNet50 perform the comparable performance on cartoon face detection with mAP 90.18 and 90.23, and far surpass the recent state-of-the-art Res2Net50, ResNeSt50, EfficientNet-B3 of 89.59, 89.97, 88.63 respectively. Owing to the superiority of OSA module, VoVNetV2-39 achieves the better performance. By introducing asymmetric convolution blocks, the proposed VoVNetV3 achieves the highest mAP score 90.74 since it can generate features with more diverse receptive fields, this is beneficial for the cartoon faces with large inter-class variances.

A-BiFPN. Next, take ACFD with ResNet50 backbone as the baseline, a series of simulations are carried to verify the effectiveness of proposed ABi-FPN by comparing with the plain FPN [Lin *et al.*, 2017a], BiFPN [Tan *et al.*, 2020] and SEPC [Wang *et al.*, 2020]. As presented in Table 3, benefit from the ability of BiFPN to aggregate multi-scale feature maps by top-down and bottom-up paths, it outperforms the plain FPN and SEPC by 1.8% and 1.3% points. The proposed ABi-FPN further enhances this ability by employing asymmetric convolution blocks to capture the features with more diverse receptive fields, and achieves the better mAP score of 90.36%.

Dynamic Anchor Match. Then, several experiments are conducted to evaluate the superiority of the proposed DAM, as shown in Table 4. We select the model with ResNet50

Module	FPN	BiFPN	SEPC	ABi-FPN
mAP (%)	88.30	90.18	88.80	90.36

Table 3: mAP (%) of different feature pyramid networks on ACFD.

backbone, the plain FPN and classic match strategy as the baseline, it gets mAP of 87.65%. In particular, the worse performance is obtained while compensating anchors by using regressed bounding boxes with a lower IoU threshold $T_2 = 0.35$, this is due to many low-quality anchors are matched for regression causing too many false-positive results. On the other hand, when compensating anchors with a higher IoU threshold but summing the losses by weights $\lambda_{reg} = \lambda_{cls} = 1.0$, the performance is not the best due to the loss of compensated anchors may dominate. By choosing $T_1 = 0.35, T_2 = 0.7, \lambda_{reg} = \lambda_{cls} = 0.7$, the model equipped with DAM achieves mAP score of 88.90%, it is about 1.3% points higher than the baseline.

T_1	T_2	λ_{res}	λ_{cls}	mAP (%)
-	-	-	-	87.65
0.35	0.35	1.0	1.0	87.56
0.35	0.7	1.0	1.0	88.60
0.35	0.35	0.5	0.5	87.54
0.35	0.7	0.5	0.5	88.75
0.35	0.7	0.7	0.7	88.90

Table 4: mAP (%) of ACFD matching anchors by proposed DAM with different parameters.

Margin Binary Classification. At the next stage, the effectiveness of different margins in MBC is verified by simulations shown in Table 5, in which $m = 0$ means margin loss is disabled. Too small margin value makes the model not work and too large value makes the models hard to optimized, $m = 0.2$ is adopted in our ACFD.

margin (m)	0	0.1	0.2	0.3	0.5
mAP (%)	90.48	90.51	90.73	90.67	90.41

Table 5: mAP (%) of ACFD with different margin parameters.

Ablation Study. Finally, we take model with ResNet50 backbone and the plain FPN as the baseline and demonstrate the effectiveness of proposed modules by adding each of them to the baseline. As listed in the first and second rows of Table 6, the baseline achieves mAP of 87.85% and 87.13% with and without the multi-scale test. By adding all modules to the baseline, mAP scores 90.94% and 90.33% are obtained with and without the multi-scale test, which surpasses the baseline more than 3.0 points mAP. This final model achieves 92.91% on the leaderboard of competition and wins the first

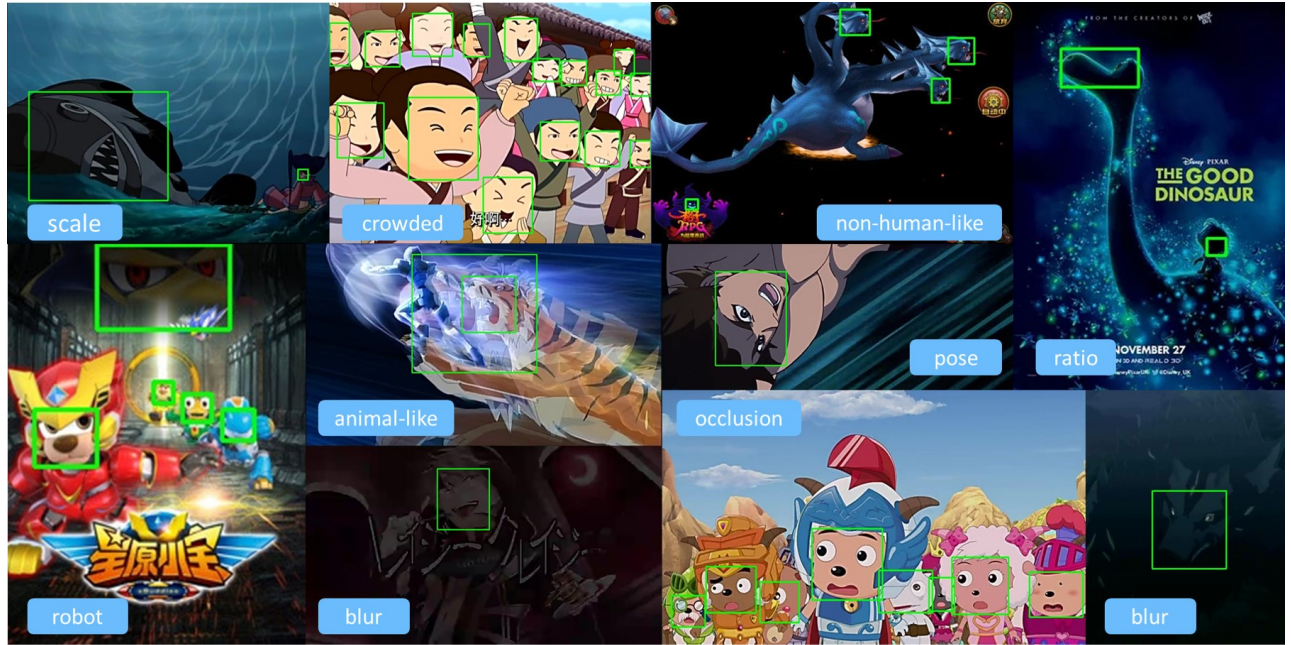


Figure 4: Illustration of our ACFD to various large variations on iCartoonFace dataset. Green bounding boxes indicate the detector confidence is above 0.7.

place with a large margin. Fig. 4 shows more examples to demonstrate the effectiveness of ACFD.

VoVNetV3	ABi-FPN	DAM	MBC	ms-test	mAP (%)
					87.13
				✓	87.85
		✓			88.47
✓		✓			89.46
✓	✓	✓			89.39
✓	✓	✓			90.17
✓	✓	✓	✓		90.33
✓	✓	✓	✓	✓	90.94

Table 6: Effectiveness of our proposed modules. Taken a model with ResNet50 backbone as baseline, all models are trained on 45000 training images and evaluated with mAP (%) on 5000 validating images.

4.3 Time-Consuming Optimization

In order to meet the competition requirements that the inference time for a single image should not exceed 50 ms, we utilize some engineering tricks to optimize our ACFD. Firstly, due to the convolution and normalization operations in asymmetric convolution blocks are linear, so the horizontal and vertical branches can be merged to become a classic convolution block, as shown in Fig. 5. Furthermore, the convolution and normalization operations can be merged into a single convolution operation as follow,

$$a = \frac{\gamma}{\sqrt{\delta^2 + \epsilon}}, W' = W * a, B' = (B - \mu) * a + \beta, \quad (4)$$

where μ , δ , γ , and β are mean, variance and affine parameters of origin normalization operation, W , B and W' , B' are parameters of convolution layer before and after merging. Finally, we convert the PyTorch model to the TensorRT version

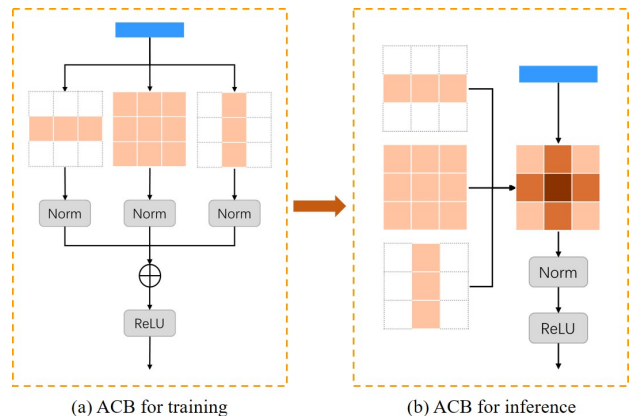


Figure 5: Illustration of asymmetric convolution blocks (a) original structure for training and (b) merged three convolutions into one for inference.

by torch2trt* tool and evaluate the model with a large batch size to further accelerate the inference speed.

5 Conclusion

Aiming at the difficulties of cartoon face detection, a novel asymmetric cartoon face detector (ACFD) is proposed in this paper. In particular, a strong backbone network VoVNetV3 is introduced to extract feature maps with more diverse receptive fields benefit from the asymmetric convolution blocks; then, the multi-scale features are better aggregated and enhanced by the proposed ABi-FPN. In addition, we propose a dynamic anchor match strategy to match high-quality anchors for each face by making full use of the regressed bounding boxes; furthermore, margin binary classification loss is used to enhance the discrimination of network to better distinguish faces from the dense predictions.

*<https://github.com/z-bingo/torch2trt>

References

- [Chi *et al.*, 2019] Cheng Chi, Shifeng Zhang, and et. al. Selective refinement network for high performance face detection. In *AAAI*, volume 33, pages 8231–8238, 2019.
- [Deng *et al.*, 2019] Jiankang Deng, Jia Guo, and et. al. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, pages 4690–4699, 2019.
- [Ding *et al.*, 2019] Xiaohan Ding, Yuchen Guo, and et. al. ACNet: Strengthening the kernel skeletons for powerful cnn via asymmetric convolution blocks. In *ICCV*, pages 1911–1920, 2019.
- [Gao *et al.*, 2019] Shanghua Gao, Ming-Ming Cheng, and et. al. Res2net: A new multi-scale backbone architecture. *TPAMI*, 2019.
- [Gong *et al.*, 2020] Julia Gong, Yannick Hold-Geoffroy, and et. al. Autotoon: Automatic geometric warping for face cartoon generation. In *WACV*, pages 360–369, 2020.
- [He *et al.*, 2015] Kaiming He, Xiangyu Zhang, and et. al. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*, pages 1026–1034, 2015.
- [He *et al.*, 2016] Kaiming He, Xiangyu Zhang, and et. al. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [Hu *et al.*, 2018] Jie Hu, Li Shen, and et. al. Squeeze-and-excitation networks. In *CVPR*, pages 7132–7141, 2018.
- [Huang *et al.*, 2017] Gao Huang, Zhuang Liu, and et. al. Densely connected convolutional networks. In *CVPR*, pages 4700–4708, 2017.
- [Jha *et al.*, 2018] Saurav Jha, Nikhil Agarwal, and et. al. Bringing cartoons to life: Towards improved cartoon face detection and recognition systems. *arXiv*, 2018.
- [Kong *et al.*, 2019] Tao Kong, Fuchun Sun, and et. al. Consistent optimization for single-shot object detection. *arXiv*, 2019.
- [Lee and Park, 2020] Youngwan Lee and Jongyoul Park. Centermask: Real-time anchor-free instance segmentation. In *CVPR*, pages 13906–13915, 2020.
- [Lee *et al.*, 2019] Youngwan Lee, Joong-won Hwang, and et. al. An energy and gpu-computation efficient backbone network for real-time object detection. In *CVPR Workshops*, pages 0–0, 2019.
- [Li *et al.*, 2019a] Jian Li, Yabiao Wang, and et. al. Dsfed: dual shot face detector. In *CVPR*, pages 5060–5069, 2019.
- [Li *et al.*, 2019b] Shichao Li, Yi Zheng, and et. al. icartoon-face: A benchmark of cartoon person recognition. *arXiv*, 2019.
- [Lin *et al.*, 2017a] Tsung-Yi Lin, Piotr Dollár, and et. al. Feature pyramid networks for object detection. In *CVPR*, pages 2117–2125, 2017.
- [Lin *et al.*, 2017b] Tsung-Yi Lin, Priya Goyal, and et. al. Focal loss for dense object detection. In *ICCV*, pages 2980–2988, 2017.
- [Liu *et al.*, 2016] Wei Liu, Dragomir Anguelov, and et. al. Ssd: Single shot multibox detector. In *ECCV*, pages 21–37. Springer, 2016.
- [Liu *et al.*, 2019] Yang Liu, Xu Tang, and et. al. Hambox: Delving into online high-quality anchors mining for detecting outer faces. *arXiv*, 2019.
- [Ren *et al.*, 2015] Shaoqing Ren, Kaiming He, and et. al. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [Rezatofighi *et al.*, 2019] Hamid Rezatofighi, Nathan Tsoi, and et. al. Generalized intersection over union: A metric and a loss for bounding box regression. In *CVPR*, pages 658–666, 2019.
- [Simonyan and Zisserman, 2014] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv*, 2014.
- [Szegedy *et al.*, 2017] Christian Szegedy, Sergey Ioffe, and et. al. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*, 2017.
- [Tan and Le, 2019] Mingxing Tan and Quoc V Le. EfficientNet: Rethinking model scaling for convolutional neural networks. *arXiv*, 2019.
- [Tan *et al.*, 2020] Mingxing Tan, Ruoming Pang, and et. al. Efficientdet: Scalable and efficient object detection. In *CVPR*, pages 10781–10790, 2020.
- [Tang *et al.*, 2018] Xu Tang, Daniel K Du, and et. al. Pyramidbox: A context-assisted single shot face detector. In *ECCV*, pages 797–813, 2018.
- [Wang *et al.*, 2018] Hao Wang, Yitong Wang, and et. al. Cosface: Large margin cosine loss for deep face recognition. In *CVPR*, pages 5265–5274, 2018.
- [Wang *et al.*, 2020] Xinjiang Wang, Shilong Zhang, and et. al. Scale-equalizing pyramid convolution for object detection. In *CVPR*, 2020.
- [Wu *et al.*, 2019] Ruizheng Wu, Xiaodong Gu, and et. al. Landmark assisted cyclegan for cartoon face generation. *arXiv*, 2019.
- [Yang *et al.*, 2016] Shuo Yang, Ping Luo, and et. al. Wider face: A face detection benchmark. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5525–5533, 2016.
- [Zhang *et al.*, 2017] Shifeng Zhang, Xiangyu Zhu, and et. al. S3FD: Single shot scale-invariant face detector. In *ICCV*, pages 192–201, 2017.
- [Zhang *et al.*, 2020a] Hang Zhang, Chongruo Wu, and et. al. ResNeSt: Split-attention networks. *arXiv*, 2020.
- [Zhang *et al.*, 2020b] Shifeng Zhang, Cheng Chi, and et. al. Refineface: Refinement neural network for high performance face detection. *TPAMI*, 2020.
- [Zheng *et al.*, 2020] Zhaohui Zheng, Ping Wang, and et. al. Distance-iou loss: Faster and better learning for bounding box regression. In *AAAI*, pages 12993–13000, 2020.