

Adversarial Occlusion-aware Face Detection

Yujia Chen^{1,2}, Lingxiao Song^{1,2,3}, Ran He^{1,2,3*}

¹National Laboratory of Pattern Recognition, CASIA

²Center for Research on Intelligent Perception and Computing, CASIA

³University of Chinese Academy of Sciences, Beijing 100190, China

Abstract

Occluded face detection is a challenging detection task due to the large appearance variations incurred by various real-world occlusions. This paper introduces an Adversarial Occlusion-aware Face Detector (AOFD) by simultaneously detecting occluded faces and segmenting occluded areas. Specifically, we employ an adversarial training strategy to generate occlusion-like face features that are difficult for a face detector to recognize. Occlusion is predicted simultaneously while detecting occluded faces and the occluded area is utilized as an auxiliary instead of being regarded as a hindrance. Moreover, the supervisory signals from the segmentation branch will reversely affect the features, helping extract more informative features. Consequently, AOFD is able to find the faces with few exposed facial landmarks with very high confidences and keeps high detection accuracy even for masked faces. Extensive experiments demonstrate that AOFD not only significantly outperforms state-of-the-art methods on the MAFA occluded face detection dataset, but also achieves competitive detection accuracy on benchmark dataset for general face detection such as FDDB.

1. Introduction

Face detection has been well studied in recent years. From the pioneering work of Viola-Jones face detector [27] to recent state-of-the-art CNN-based methods, the performance of face detectors has been improved remarkably. For example, the average precision has been boosted to over 98% [9, 21, 35] in the unconstrained FDDB dataset.

Although face detection algorithms have obtained quite good results under general scenarios, detecting faces in specific scenarios is still worth studying. For instance, one of the remaining challenges is partially occluded face detection. Facial occlusions occur frequently, e.g. facial accessories including sunglasses, masks and scarfs. Occluded



Figure 1. The proposed AOFD is able to detect various heavily-occluded faces, while the occlusion-aware segmentation branch predicts masks of the occluded area.

faces are only partially visible, and occluded regions have arbitrary appearances that may diverse from normal face regions. Hence occluded faces have significant intra-class variation, leading to difficulties in learning discriminative features for detection. A standard paradigm to address this problem is to enlarge the training dataset of occluded faces, but it can't solve this problem in essence. Moreover, the lack of large-scale occluded face datasets makes it harder to handle this obstacle.

In this paper, we propose a framework for occluded face detection, aiming at formulating a new strategy to tackle the problem of limited occluded face training data, and exploiting the power of CNN representations for the faces with occlusions as far as possible. Firstly, motivated by the remarkable success achieved by adversarial learning in recent years, a deep adversarial network is proposed in our approach to generate face samples with occlusions from the mask generator. A compact constraint is adopted to reinforce the realness of generated masks. Secondly, we introduce an occlusion-aware model by predicting the occlusion segments at the same time with detecting faces.

In all, the generator aims to make the model focus more on the exposed areas, while the segmentation branch is to

*corresponding author

extract more informative features of the occluded area. Intuitively, jointly solving these two tasks can be reciprocal.

To sum up, we make contributions in the following aspects:

- A novel adversarial framework is proposed to alleviate the lack of occluded training face images by generating occluded or masked face features. We employ a compact constraint to get more realistic occlusions.
- Mask prediction is conducted simultaneously while detecting occluded faces. The occluded area will NOT be regarded as a hindrance but an auxiliary of face detection.
- Experimental evaluations on the MAFA dataset demonstrate that the proposed AOFD can significantly improve the face detection accuracy under heavily occlusions. Besides, AOFD can also achieve competitive performance on the unconstrained face detection benchmark.

2. Related Work

We first briefly survey face detection algorithms, followed by a review of the state-of-the-art occluded face detection researches.

Face detection can be considered as a special task of object detection. Successful general face detection algorithms often show great performance on face recognition. The Viola-Jones [27] detector can be recognized as a milestone in the field of face detection. They innovatively adopted AdaBoost to train cascade classifier with Haar-like features, which first makes it possible to apply face detection in real-time applications. Following their work, lots of boosting-based models were proposed [15, 2, 20, 36], focusing on designing more sophisticated hand-crafted features or improving the boosting strategy. More Recently, CNN features [30] were utilized in this boosting framework. Another famous category of face detectors is DPM-based. Deformable part models [4] were proposed for object detection at first, which acquired impressive accuracy in complex environment. Inspired by this model, many extensions of DPM were developed to face detection [6] by modeling potential deformations among facial parts. However, DPM models suffered from the high computational complexity, making it difficult to be applied in real-world applications such as digital cameras, phones or other mobile devices.

Generally speaking, boosting-based methods and DPM-based methods design features and optimize classifiers separately. The pipeline of these methods is divided into two stages, which is not an end-to-end architecture. Recently, benefitting from the prosperity of social network and big data, numerous deep learning based object detection algorithms have been proposed [7, 23, 12, 17, 1]. CNN-based

detectors therefore have become the mainstream in face detection gradually [3, 13]. CNN-based face detectors directly learn robust face representations from data and optimize classifiers in an end-to-end style. For example, [34] developed a deep cascaded multi-task framework that predict face and landmark location in a coarse-to-fine manner, and [33] further improved the performance of cascade models by optimizing feature selection algorithms.

Although many efforts have been made in face detection, the performance of occluded face detection is still far from satisfactory, and there are few works on occluded face detection as far as we know. [31] explicitly inferred face-ness score through local part responses via an attribute-aware model. But additional face-specific attribute annotations needed in this method were very difficult to collect. [22] introduced a specific grid loss layer into CNNs that minimized the error rates on each sub-block of the feature map independently, thus every sub-part is discriminative on its own. [19] introduced a partial face detection approach based on detection of facial segments. They mainly focused on detecting incomplete faces that captured by the front camera of smart phones. Recently, [5] combined pre-trained CNN features with local linear embedding (LLE-CNN) to get similarity-based descriptors for partially visible faces. They built a dataset for masked face detection specifically, named the MAFA that contains 35K occluded faces. [28] applied anchor-level attention on Feature Pyramid Networks [16].

As mentioned above, our work is also related to adversarial learning. Generative Adversarial Network (GAN) [8] has shown great performance in numerous computer vision applications including image style transfer [37, 11], image generation [25, 10] and so on. Adversarial learning provides a simple yet efficient way to train powerful models via the min-max two-player game between the generator and the discriminator. Most of the previous work focused on promoting generators. Recently, researchers began to pay attention to increase the capacity of discriminator by adversarial learning. [29] used adversarial learning in generating hard examples for object detection. [14] employed Perceptual GAN to enhance the representations for small objects. Inspired by these applications, we develop an adversarial occlusion-aware model, which can synthesize occlusion-like face features for boosting occluded face detectors.

3. Methods

In this section, we propose an AOFD method to tackle one of the most common and vicious problems in face detection-occlusion problem. We first analyze the occluded face detection problem (Sec. 3.1) and summarizes the overall architecture of AOFD (Sec. 3.2), and then introduce the mask generation and segmentation method in AOFD in Sec. 3.3 and Sec. 3.4, respectively.

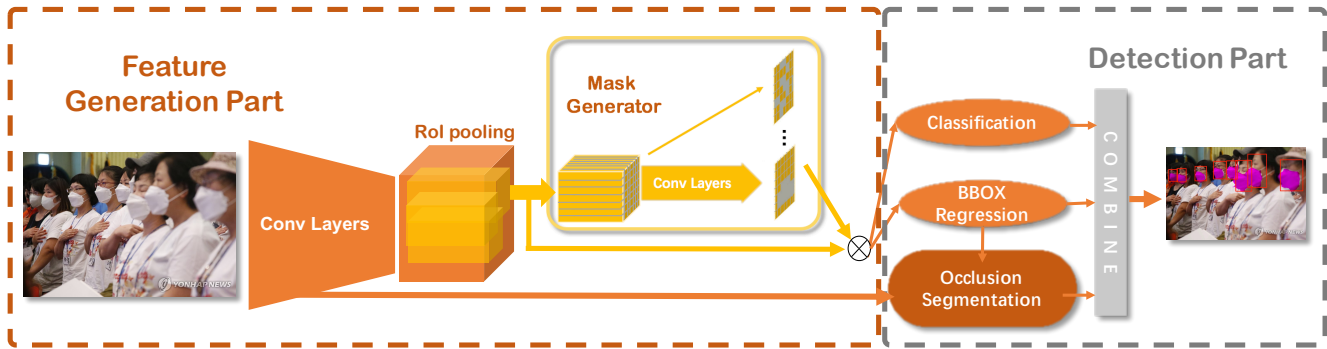


Figure 2. The overall architecture of the proposed method. Mask generator are operated directly on RoIs where one type of masks are selected for each proposal.

3.1. Problem Analysis

In real-world situations, we can generally classify face occlusion problems into three categories: facial landmark occlusion, occluded by faces and occluded by objects. Facial landmark occlusion includes conditions like wearing glasses and gauze masks. Occluded by faces is a complicated situation because a detector easily mis-recognize several faces into one or only detect a part of the faces. The segmentation method is proposed in order to mitigate this problem. When occluded by an object, usually more than half of a face will be directly masked. An original masking strategy is used to mimic these in-the-wild situations.

We also visualized features of occluded faces in Figure 3 (a), finding that occluded areas rarely respond. For some heavily occluded faces, useful information in feature maps is too scarce for a detector to identify. To tackle this problem, we may need to enhance representation ability of exposed area. Meanwhile, recognition of occluded area can also bespeak that “there is a face” on the condition that sufficient context information is provided. For the most complex problem where a face is occluded by another face, the context area should cover at least the nearby faces and a larger receptive field is required so that the integrity of the background information can be ensured. This idea is enlightened by human vision, that is, human need a large context to define a small or incomplete object. Besides, as the quality of features directly determines the results, segmentation is better conducted on image features than on RoIs in order to extract more informative feature maps.

3.2. Overall Architecture

In order to detect faces with heavy occlusion, Adversarial Occlusion-aware Face Detector (AOFD) is designed with the view of (1) effectively utilizing the exposed facial areas, and (2) transferring the interference of the occlusions into beneficial information. For the first problem, we find that undetected faces are typically those with their characteristic part of face occluded, such as eyes and mouth. One feasible way is to mask the distinctive part of face in train-

ing set, forcing the detector to learn what possibly a face looks like even if there is less exposed area. To this end, a mask generator is designed in an adversarial way to generate a mask for each positive sample. It will generate different masks with faces of different poses. A masking strategy is applied for a better utilization of the mask generator as well. More details are illustrated in Sec. 3.3.

For the latter problem, We believe that finding common occlusions is helpful to detect incomplete faces behind them. Thus, an occlusion segmentation branch is introduced to segment occluded areas including hair, glasses, scarves, hands and other objects. This is not an easy task due to few training samples. Therefore, we labeled 374 training samples downloaded from internet for occlusion segmentation and came up with an original training strategy. This dataset is denoted as SFS (small dataset for segmentation). More details are listed in Sec. 4.

As is demonstrated in Figure 2, a mask generator is added after a region of interest (RoI) pooling layer, followed by a classification branch and a bounding box regression branch. Finally, a segmentation branch is in responsible to segment the occluded area inside each bounding box. The final result combining classification, bounding box regression and occlusion segmentation will be output in the end. The overall loss of our architecture takes the following multi-task form:

$$L = \alpha L_c + \beta L_b + \mu L_s \quad (1)$$

where L_c denotes a binary softmax loss for classification, L_b denotes a smooth L1 loss for bounding box regression. We apply a binary softmax loss for segmentation branch, which is L_s . During training, the coefficients α , β and μ are set 1, 1, and $1e^{-5}$ respectively.

3.3. Mask Generator

Mask generator: Since human face is very structural, facial features tend to appear in similar locations. However, with different poses, expressions and occlusions, distinguished facial area varies significantly. Our aim is to find

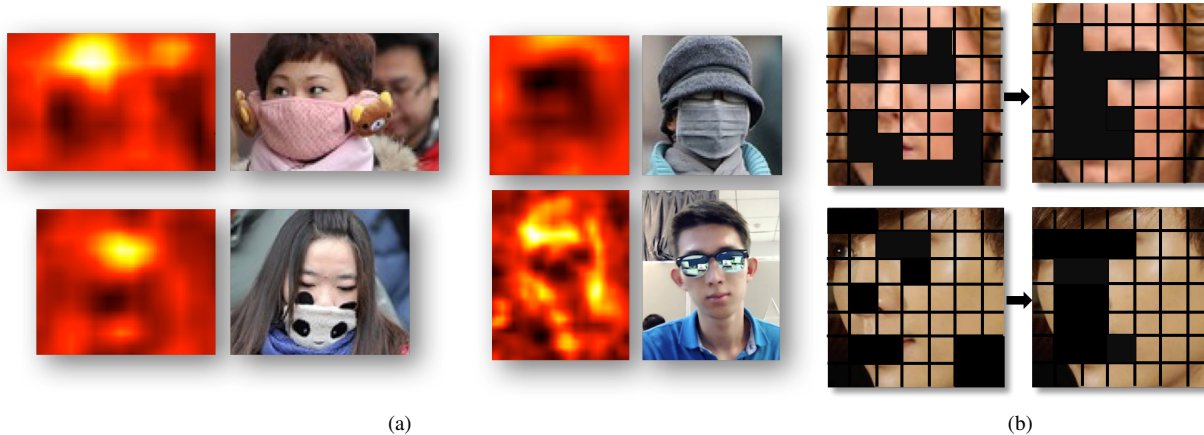


Figure 3. (a) Heat map of the features from conv5_3 extracted from Faster RCNN. Occluded areas have less information in the features (the black areas). (b) The compact constraint has made the generated masks more accurate and efficient. A quarter of the minimum values is selected as the mask and is marked black. The masks are operated directly on the 7 by 7 RoIs. We map these masks to their corresponding receptive fields in the original image space for a better visualization.

this distinguished area and to generate a customized mask. We visualize some mask examples in Figure 3 (b). As we have observed, occluded area in features rarely respond in real images. To simulate this characteristic, masks are directly operated on RoIs. Therefore, the generator, which contains four convolutional layers with a straight mapping, is designed simply as it can be regarded as a binary prediction problem. Besides, the peculiarity of our mask generator, to distinguish from [29], is the original generating procedure and the mask forms. Since face structures are inherently different from those of objects, they need to be learned in a more subtle and flexible way, or no plausible mask can be obtained.

Masking strategy: The generated mask is a one-channel heat map where 0 represents masked area and 1 otherwise. During training, each pixel value will be squeezed to zero or one. We select a quarter of the minimum values as the mask when training the generator and one-third of the minimum values when training the overall model.

Heavily occluded samples after masking will become an extremely hard source for training, making the model difficult to converge. To this end, three types of masks are proposed and jointly training with the original features. The first type is to use mask generator, which corresponds to facial landmark occlusion. The second type is to mask half of the features, whether left, right, top or bottom, and the third type is randomly dropping half of the pixels. This masking strategy embodies in-the-wild occlusion types analyzed in Sec. 3.1.

Loss function: When training the mask generator, we employ an adversarial training method. We aim to increase classification loss as much as possible. Since a masked area is limited and a distinguished facial area is comparatively salient in feature maps, the model can easily con-

verge. However, we find it not enough because the occluded area is sometimes strip-like or sporadic, while it is supposed to be more compact in real situations. Recall that the areas with longer or irregular edges will have a larger value for each pixel using a kernel of an edge detector. A kernel to make the occluded area sleeker and more circular is designed as a compact constraint for generated masks. The loss function is:

$$L_g = \gamma L_{com} - \eta L_c \quad (2)$$

where L_g denotes the loss for generator, L_{com} denotes a compact loss, and γ and η are coefficients. γ is set $1e^{-6}$ and η is set 1 in order to balance the derivatives. The compact loss is computed with a convolutional layer in a way as follows:

$$L_{com} = \sum((1 - mask) * kernel) \quad (3)$$

where $*$ denotes a convolutional operation, $mask$ is the first type of mask generated by the mask generator and the last item is the designed kernel, which is

$$\begin{bmatrix} -\frac{1}{8} & -\frac{1}{8} & -\frac{1}{8} \\ -\frac{1}{8} & 1 & -\frac{1}{8} \\ -\frac{1}{8} & -\frac{1}{8} & -\frac{1}{8} \end{bmatrix}$$

In this way, strip-like or sporadic areas will get very high penalty and more reasonable masks can be obtained.

3.4. Segmentation

Design: Previous works on segmentation have proved that CNNs is capable of comprehending the semantic information of a picture and elaborately conduct a pixel-wise

classification. When combining detection with segmentation, it is usually designed in RoI level to achieve higher accuracy.

Considering segmenting each RoI, one problem in occluded face situation is that the overlap of two bounding boxes will have different meanings. For example, if one face is occluded by another face, part of the front face should be regarded as an occlusion for the back face, while there shouldn't be any occluded area for the front face. Since our destination is to utilize the effective information contained in the occluded area to confirm if there is a back face and then make the exposed area more distinguished, ample context information is required. Moreover, it is the features that really matter in the bounding box classification and regression branch. With reasons above, segmentation is conducted in image level to directly affect the image feature maps. Therefore, the detector is able to find faces with more informative features embodying image-level signals like the appearance of an occlusion or a person. We call this method as an occlusion-aware method.

The segmentation branch is designed in a fully convolutional way (Figure 4). In order to obviate noise, it follows a bounding box regression branch and only areas inside bounding boxes are maintained. Bounding boxes are enlarged in scale with a factor of 1.3 before dropping the noise. Because it is only an auxiliary task for face detection, we didn't compare it with other segmentation methods. The results (Figure 1, Sec.4) should have shown the effectiveness. Although the final results have proved the feasibility of this method, the edges of segmentation seem to be a bit rough. This is caused by the limited size of the SFS training set. Nevertheless, we have verified the possibility to train the model with very limited training samples.

Loss function: We choose softmax loss instead of L1 or L2 loss used in some segmentation and image generation tasks [10, 14] because it helps stabilize the training process.

4. Experiments

In this section, we qualitatively evaluate the proposed method with state-of-the-art methods. We first introduce detailed information during training (Sec. 4.1), and then test AOFD on several comparative benchmarks (Sec. 4.2). A series of ablative studies are conducted to verify the effectiveness of our method (Sec. 4.3).

4.1. Training Details

There are two stages in the training procedure. Firstly, based on Faster RCNN [23], we train the mask generator with the loss in equation (2). Secondly, the detector and the segmentation branch are trained jointly with the parameters in the mask generator fixed.

In the second stage, due to the limitation of training data for segmentation, an unordinary training strategy is needed.

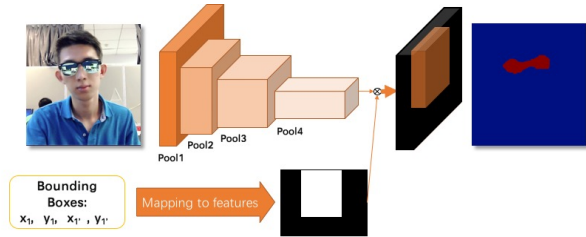


Figure 4. The image segmentation branch in AOFD. Only areas inside bounding boxes are maintained so as to reduce noise.

We first train on SFS for 10k iterations, causing overfitting in segmentation. Then the model is trained on the combination of WIDER FACE training set and SFS for 50k iterations with loss weights for segmentation set $1e^{-7}$ and finally the model is tuned only on SFS for 3 epochs. Derivatives from WIDER FACE training set will be zero for the segmentation branch during training. Because there are far more training images from the WIDER FACE training set than SFS, the segmentation branch can only be trained every several iterations while the features are changing all the time when training on the combined dataset. In this way, the overall loss can get rid of local minima and the overfitting problem can be solved. The basic learning rate is 0.001. AOFD runs 5 FPS on a TITAN X GPU, which is similar to the original Faster RCNN.

Experiment settings: AOFD is based on a Faster RCNN with a VGG16 backbone [26]. For anchors of PRN, we use three aspect ratios (1.7, 1 and 1.3) and four scales (64^2 , 128^2 , 256^2 and 512^2). Batch size is set 1. An RoI is treated as foreground if its intersection over union (IoU) with any ground truth bounding box is higher than 0.5. To balance the number of foreground and background training samples, the ratio of foreground RoIs to background RoIs is set 1:3. During training, the short side of an input image is resized to either 512 or 1024 on condition that long side is no longer than 1024.

Small dataset for segmentation (SFS): SFS contains 376 images with 1138 labeled faces downloading from the Internet. There is at least one occluded face in each image and over 80% of the faces are occluded.

4.2. Evaluation on benchmarks

Our model is trained on the WIDER FACE [32] training set and evaluated on the Fddb and MAFA [5] databases. Although the MAFA database dose not release its training set, we still obtain state-of-the-art results on the MAFA testing set without fine tuning the model to adjust the variance between different annotation protocols.

Fddb (Face Detection Data Set and Benchmark) is an unconstrained dataset for face detection. It has 2, 845 images with 5, 171 faces. The detection results of different methods are shown in Figure 6. [18] and several other methods obtain higher continuous score because they

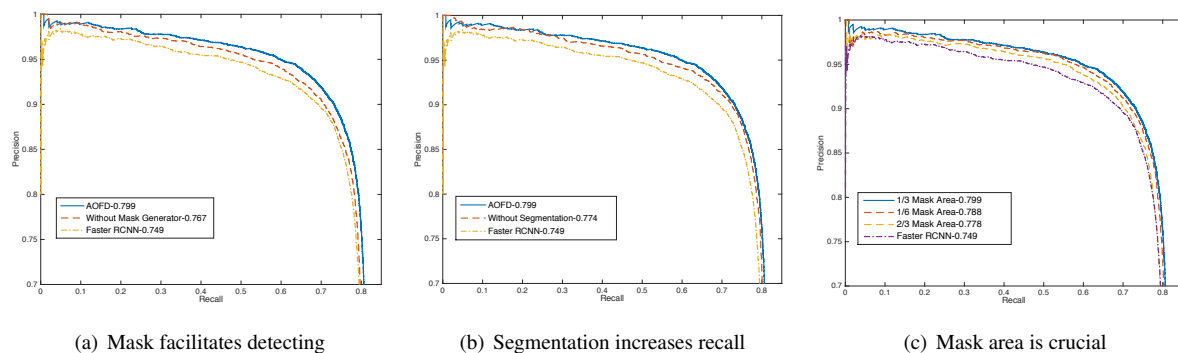


Figure 5. PR curves of ablative studies on the whole MAFA testing set without OHEM.

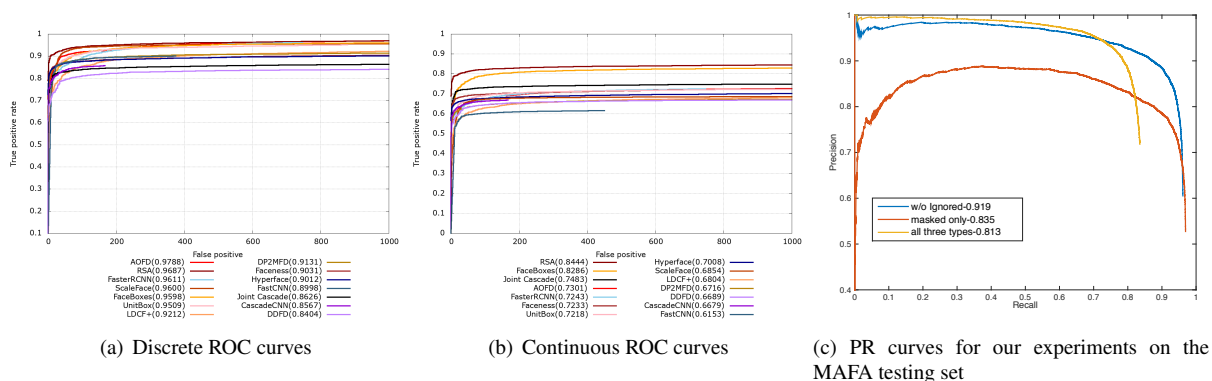


Figure 6. Results on Fddb ((a)(b)) and MAFA testing set ((c)).

have transformed the rectangle bounding boxes into ellipse ones. The fact that we didn’t carry out extra training in the Fddb training set may lead to the increase of localization errors because of the difference of annotation criterion. However, in the comparison with state-of-the-art methods, we observe that AOFD outperforms all the other methods in terms of discrete score, demonstrating its strong ability to detect nearly all large faces even if faces with short side less than around 15 pixels are mostly neglected due to the anchor setting.

Furthermore, our AOFD also obtains a higher recall rate at 1000 FPs on Fddb than other Faster RCNN methods with similar settings by a large margin (Figure 5 and Figure 6). The superior performance also reveals that applying masking strategy and training with a segmentation task are valuable attempts to enhance the model’s capacity.

MAFA is designed for the evaluation of masked face detection, which contains 35806 face annotations with a minimum size of 32×32 . Since the MAFA testing set uses squares to label faces, the rectangle bounding boxes in our results are transformed into squares to match the annotation.

There are three types of annotations in the MAFA dataset: masked, unmasked and ignored. Blurry or deformed faces or those with side length less than 32 pixels

Methods	All	‘masked’ only	w/o ‘Ignored’
AOFD	81.3%	83.5%	91.9%
FAN	-	76.5%	88.3%
LLE-CNNs	-	-	76.4%
MTCNN	-	-	60.8%

Table 1. Average precision on the MAFA testing set.

are labeled as ‘Ignored’. But we find that many ‘ignored’ faces are also acceptable. Since the other methods ([28] [5] [34]) didn’t count those annotations labeled as ‘Ignored’, we report our results on both MAFA subsets and the whole testing set for comparison (Table 1).

As shown in Table 1, the average precision achieves the highest 91.9% (threshold 0.5) if we only evaluate on the faces with ‘masked’ and ‘unmasked’ labels. The result outperforms LLE-CNNs [5] by a large margin and is also better than the state-of-the-art Face Attention Network (FAN) [28]. Since AOFD is proposed to address occlusion problem, we also evaluate our model on faces labeled as ‘masked’ only. AOFD achieves 83.5% and has 7% im-

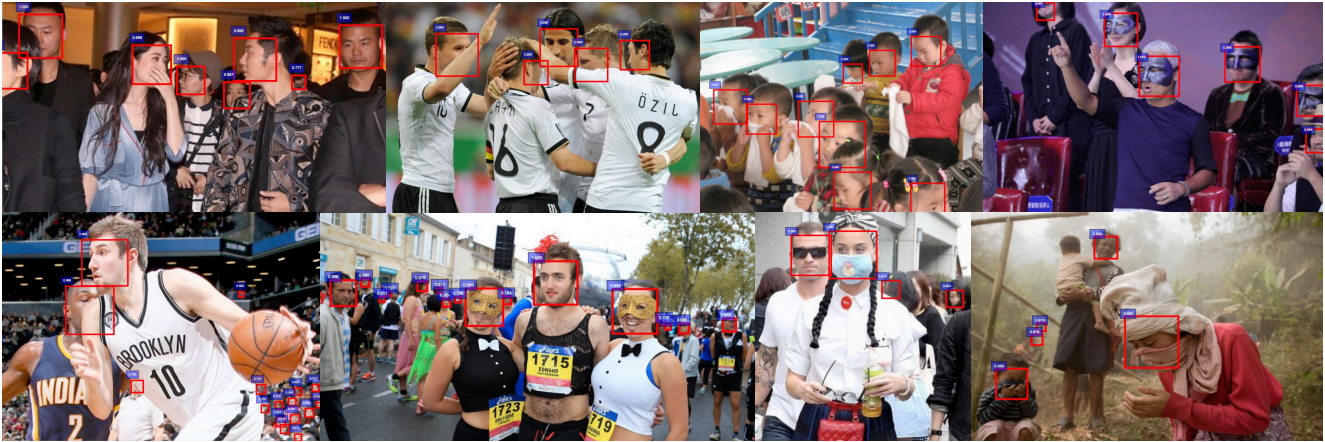


Figure 7. Qualitative results of AOFD on the test set of the MAFA dataset.

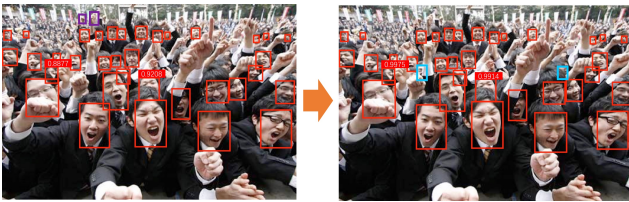


Figure 8. Not only is AOFD able to detect occluded faces with higher confidences, but also capable of increasing the average precision. The **purple** bounding boxes are false detections by Faster RCNN and the **blue** ones are new true positives detected by AOFD.

provement over the state-of-the-art result obtained by FAN.

Figure 6(c) further shows the PR (Precision-Recall) curves of the three experimental settings. If we only count the faces annotated as ‘masked’ (the orange curve in Figure 6(c)), precision witnesses a sharp drop at the beginning. This is caused by unmasked and unlabeled face detections which are regarded as FPs when evaluating masked faces. More results on MAFA are presented in Figure 7.

Furthermore, we have studied the main obstacle of our model to achieve a higher AP. The minimum IoU threshold for a true positive proposal is modified from 0.5 to 0.45, from which we observe that a slight decrease of IoU threshold can boost AP from 91.9% to 93.8%. This explains that the precision of bounding boxes can still be further improved

4.3. Model Analysis

To better understand the function of each part of our model, we ablate each component to observe AOFD’s performance. In this way, the mask generator and segmentation branch are removed one after the another. We delve into the optimal area of the mask as well and find that the mask area is crucial for the functioning of mask generator. Besides, the efficiency of the compact constraint and the comparison with online hard example mining (OHEM) [24] are also dis-

Settings	Recall rate at 1000 FPs
AOFD	97.88%
w/o segmentation	97.13%
w/o generator	96.85%

Table 2. Results of the ablative studies on FDDDB.

cussed in this section.

Mask facilitates detecting: State-of-the-art detectors are able to detect some of occluded faces, but with lower confidence. As shown in Figure 8, AOFD can increase the confidence of occluded faces by a large margin. Without the mask generator, AOFD pays less attention to exposed area or face structure, and the recall rate at 1000 false positives on FDDDB drops by 1.3%(Table 2). The sharp decline (3.2%) of average precision on the MAFA testing set in Figure 5(a) reveals the value of the mask generator as well. It is also observed that AOFD’s results would drop by around 1% with only random and square-like occlusions. Since faces have unique structure characteristics such as facial symmetry, generating adaptive occlusions is essential in order to fool the detector.

Segmentation increases recall: With the segmentation branch, the result in Table 2 witnesses an increase of 0.75%. This improvement is relatively slight because there are not many heavily occluded faces in the FDDDB testing set. The drop of average precision from 79.9% to 77.4% (Figure 5(b)) will be more convincing to confirm the effectiveness of the segmentation branch.

Mask area is crucial: We find that the mask would vitiate the detector if a mask area is too large. Nevertheless, it would be of no use if it is too small. Figure 5 gives a brief overview of our experiments, from which we find occluding one-third of features is an ideal area for a mask.

Compact constraint matters: We propose a compact

Methods	AP on MAFA	Recall on FDDB
single OHEM	75.9%	96.54%
single AOFD	79.9%	97.12%
AOFD with OHEM	81.3%	97.88%

Table 3. Comparing AOFD with OHEM on the whole MAFA and FDDB testing set.

constraint L_c to help generate more practical masks. As is mentioned in Sec. 3.3, the generated masks are discrete or sporadic and are not plausible, e.g., two pixels occlusion on the mouth, three pixels occlusion on the eyes and others on the corners (Figure 3 (b)). In our initial experiments, the average precision is 0.785 when masking 1/3 of RoIs without the compact constraint, which is similar to having 1/6 masking area in Figure 5(c). However, masks become harder and more reasonable with L_c , which can account for the increase of performance.

Comparing with OHEM: We compare online hard example mining [24] with our methods in Table 3. We can see that the performance of training a Faster RCNN with OHEM is generally worse than a single AOFD without OHEM. But the combination of these two methods leads to a better performance. Although a harder training procedure means a more robust detector under this condition, the measurement of hard level needs to be carefully handled. For example, the decrease incurred by too large masking area demonstrated in Figure 5(c).

5. Conclusion

This paper has proposed a face detection model named AOFD to address the long-standing issue of face occlusions. A novel masking strategy has been integrated into AOFD to increase training complexity, and can plastically mimic different situations of face occlusions. The multitask training method with a segmentation branch provides a feasible solution and verifies the possibility to train an auxiliary task with very limited training data. The superior performance on both general face detection and masked face detection benchmarks demonstrates the effectiveness of AOFD.

References

- [1] Y. Chen and C. Li. Gm-net: Learning features with more efficiency. *arXiv preprint arXiv:1706.06792*, 2017.
- [2] P. Dollár, R. Appel, S. Belongie, and P. Perona. Fast feature pyramids for object detection. *IEEE TPAMI*, 36(8):1532–1545, 2014.
- [3] S. S. Farfadi, M. J. Saberian, and L.-J. Li. Multi-view face detection using deep convolutional neural networks. In *ACM ICMR*, pages 643–650, 2015.
- [4] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *IEEE TPAMI*, 32(9):1627–1645, 2010.
- [5] S. Ge, J. Li, Q. Ye, and Z. Luo. Detecting masked faces in the wild with l1e-cnns. In *IEEE CVPR*, 2017.
- [6] G. Ghiasi and C. C. Fowlkes. Occlusion coherence: Localizing occluded faces with a hierarchical deformable part model. In *IEEE CVPR*, pages 2385–2392, 2014.
- [7] R. Girshick. Fast r-cnn. In *IEEE ICCV*, pages 1440–1448, 2015.
- [8] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014.
- [9] P. Hu and D. Ramanan. Finding tiny faces. In *IEEE CVPR*, 2017.
- [10] R. Huang, S. Zhang, T. Li, and R. He. Beyond face rotation: Global and local perception gan for photorealistic and identity preserving frontal view synthesis. In *IEEE ICCV*, 2017.
- [11] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *IEEE CVPR*, 2017.
- [12] Z. Lei, R. Chu, R. He, S. Liao, and S. Z. Li. Face recognition by discriminant analysis with gabor tensor representation. In S.-W. Lee and S. Z. Li, editors, *Advances in Biometrics*, pages 87–95, 2007.
- [13] H. Li, Z. Lin, X. Shen, J. Brandt, and G. Hua. A convolutional neural network cascade for face detection. In *IEEE CVPR*, pages 5325–5334, 2015.
- [14] J. Li, X. Liang, Y. Wei, T. Xu, J. Feng, and S. Yan. Perceptual generative adversarial networks for small object detection. In *IEEE CVPR*, 2017.
- [15] J. Li and Y. Zhang. Learning surf cascade for fast and accurate object detection. In *IEEE CVPR*, pages 3468–3475, 2013.
- [16] T. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie. Feature pyramid networks for object detection. *CoRR*, 2016.
- [17] R. Liu, S. Z. Li, X. Yuan, and R. He. Online Determination of Track Loss Using Template Inverse Matching. In *The Eighth International Workshop on Visual Surveillance - VS2008*, 2008.
- [18] Y. Liu, H. Li, J. Yan, F. Wei, X. Wang, and X. Tang. Recurrent scale approximation for object detection in cnn. *arXiv preprint arXiv:1707.09531*, 2017.
- [19] U. Mahbub, V. M. Patel, D. Chandra, B. Barbelo, and R. Chellappa. Partial face detection for continuous authentication. In *IEEE ICIP*, pages 2991–2995, 2016.
- [20] M. Mathias, R. Benenson, M. Pedersoli, and L. Van Gool. Face detection without bells and whistles. In *ECCV*, pages 720–735, 2014.
- [21] M. Najibi, P. Samangouei, R. Chellappa, and L. Davis. Ssh: Single stage headless face detector. In *IEEE ICCV*, 2017.
- [22] M. Opitz, G. Waltner, G. Poier, H. Possegger, and H. Bischof. Grid loss: Detecting occluded faces. In *ECCV*, pages 386–402, 2016.

- [23] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, pages 91–99, 2015.
- [24] A. Shrivastava, A. Gupta, and R. Girshick. Training region-based object detectors with online hard example mining. In *CVPR*, pages 761–769, 2016.
- [25] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb. Learning from simulated and unsupervised images through adversarial training. In *IEEE CVPR*, 2017.
- [26] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [27] P. Viola and M. Jones. Rapid object detection using a boosted cascade of simple features. In *IEEE CVPR*, volume 1, pages I–I, 2001.
- [28] J. Wang, Y. Yuan, and G. Yu. Face attention network: An effective face detector for the occluded faces. *arXiv preprint arXiv:1711.07246v1*, 2017.
- [29] X. Wang, A. Shrivastava, and A. Gupta. A-fast-rcnn: Hard positive generation via adversary for object detection. In *IEEE CVPR*, 2017.
- [30] B. Yang, J. Yan, Z. Lei, and S. Z. Li. Convolutional channel features. In *IEEE ICCV*, pages 82–90, 2015.
- [31] S. Yang, P. Luo, C.-C. Loy, and X. Tang. From facial parts responses to face detection: A deep learning approach. In *IEEE ICCV*, pages 3676–3684, 2015.
- [32] S. Yang, P. Luo, C.-C. Loy, and X. Tang. Wider face: A face detection benchmark. In *IEEE CVPR*, pages 5525–5533, 2016.
- [33] B. Yu, M. Fang, D. Tao, and J. Yin. Submodular asymmetric feature selection in cascade object detection. In *AAAI*, pages 1387–1393, 2016.
- [34] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Signal Processing Letters*, 23(10):1499–1503, 2016.
- [35] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Z. Li. S3fd: Single shot scale-invariant face detector. In *IEEE ICCV*, 2017.
- [36] C. Zhu and Y. Peng. Group cost-sensitive boosting for multi-resolution pedestrian detection. In *AAAI*, pages 3676–3682, 2016.
- [37] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *IEEE ICCV*, 2017.