# Improved Face Detection and Alignment using Cascade Deep Convolutional Network

Weilin Cong[†], Sanyuan Zhao[†], Hui Tian[‡], and Jianbing Shen[†]

[†]Beijing Key Laboratory of Intelligent Information Technology, School of Computer Science,Beijing Institute of Technology, Beijing 100081, P.R.China
[‡]China Mobile Research Institute, Xuanwu Men West Street, Beijing

## Abstract

Real-world face detection and alignment demand an advanced discriminative model to address challenges by pose, lighting and expression. Recent studies have utilized the relation between face detection and alignment to make models computationally efficiency, but they ignore the connection between each cascade CNNs. In this paper, we combine detection, calibration and alignment in each Cascade CNN and propose an HEM method for *End-to-End* cascade network training, which give computers more space to automatic adjust weight parameter and accelerate convergence. Experiments on FDDB and AFLW demonstrate considerable improvement on accuracy and speed.

**Keywords:** Face detection, facial key points detection, cascade CNNs, End-to-End

# 1   Introduction

Face detection and alignment have play an important role in many face-based applications, challenges mainly come from various poses, uncontrollable illuminations and exaggerated expressions. Inspired by the remarkable performance of deep learning, some of the convolution based feature extraction methods have been proposed in recent years. In face detection area, Li *et al.*[2] use cascade CNNs for face detection. An extra calibration stage is added after each detection stage which cost extra computing expense on bounding box calibration. Yang *et al.*[5] extracte the features of hair, eyes, nose, mouth and neck regions with five separate CNNs, combine the result of five CNNs and use the position information to improve face detection result. Due to it's complex structure, this approach is time costly in practice. Zhang *et al.*[9] proposed a framework that joint detection and alignmnet in single CNN, however it is not an end-to-end structure and the model is very complicated to train. Face alignment also attracts extensive interests. Sun *et al.*[3] use three-level CNNs extract global and local features to estimate the positions of facial key points. Zhou *et al.*[4] use coarse-to-fine four-level cascade architecture extend 5 points detection to 68 key points detection in real-time detection.

Since detection always followed by alignment, we adopt a structure to calculate detection, calibration with alignmnet simultaneously and propose an End-to-End strategy to unify the Cascade structure. Hard Example Proposal Network is used to generate high-quality proposals.

This model is tiny but fast. The summary size of models is 2MB, much smaller than well known ImageNet Challenge models like VGG16(is about 510MB) and AlexNet(is about

250MB). For $800 \times 600$ images, our detector can detect faces and key points 90 fps on Titan X, faster than most state-of-the-art models[1]. The contributions of this paper are summarized as follows: (1)We propose an Hard Sample Proposal Structure to automatic generate higher quality training samples simultaneously for each cascade CNN. (2)We propose an End-to-End architecture and combine detection, calibration with alignment to make system more robust and accurate. (3)Extensive experiments on benchmarks show impressive improvement than the state-of-the-art face detection and alignmnent tasks.



Figure 1: Selected examples of face detection and key points detection on FDDB test set. A score threshold of [0.6,0.6,0.7] is used for 12net, 24net and 48net to display these images. The running time for obtaining these results is 15ms per image in average.

## 2 Approach

### 2.1 Overall framework

The overall pipeline of our framework is shown in Figure 2. Given a test image, 12net use fully convolution method look through the image only one time to remove most of the negative candidates and calibrate positive candidates. Local-NMS is applied to eliminate candidate windows with high overlap ratio. Then remaining candidate windows are cropped and scaled to $24 \times 24$ for 24net to remove more negative candidates and calibrate positive ones, Local-NMS is also applied to further reduce the number of candidates. Finally, 24net's outputs are cropped and scaled to $48 \times 48$ for 48net to estimate the location of the last candidate windows and facial key points. Global-NMS will eliminate the candidate windows with higher proportion IoU and output the final results.

### 2.2 Loss Function

We consider three parts for loss functions: face vs non-face classification loss, candidate windows regression loss and facial key points detection loss. We minimize a muiti task loss defined as:

$$L(cls_i, reg_i, pts_i) = \alpha \sum L_{cls}(cls_i, cls_i^*) + \beta \sum L_{reg}(reg_i, reg_i^*) + \gamma \sum L_{pts}(pts_i, pts_i^*)) \quad (1)$$

Where $\alpha$, $\beta$, $\gamma$ are used to balance different task loss. In 12net and 24net we keep $\alpha = 1, \beta = 0.5, \gamma = 0.5$ while in 48net $\alpha = 1, \beta = 0.5, \gamma = 1$.

The face vs non-face classification task can be regarded as binary classification problem. For each sample $x_i$, we use logistic regression loss after softmax layer:

$$L_{cls}(y_i, y_i^*) = -(y_i^* log(y_i) + (1 - y_i^*)log(1 - y_i)) \quad (2)$$

Where $y_i$ is the probability of being a face, $y_i^* \in \{0, 1\}$ denote ground truth labels.

---

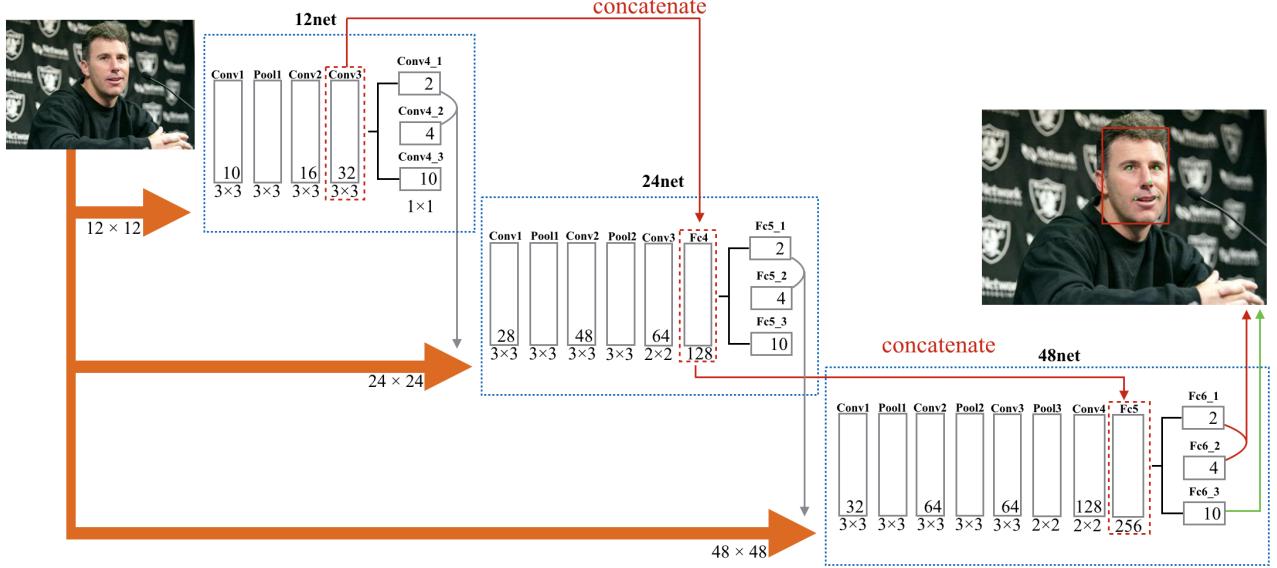[1]FRCN can detect object 7 fps, SSD detect object 59 fps

2

Figure 2: Pipeline of our detector: 12net produce multi candidate windows, 24net refine these candidate windows and 48net produces final bounding box and facial landmarks position. Each single CNN produce three different outputs: face vs non-face classification result, candidate windows calibration result and key points detection result. An End-to-End strategy is adopted between 12net and 24net, 24net and 48net.

For each candidate window, we predict the offset between it and the nearest ground truth. For each sample $x_i$, we use SmoothL1 loss:

$$L_{reg}(y_i, y_i^*) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{else} \end{cases} \tag{3}$$

Where $y_i$ is the predicted offset, $y_i^* = \{x, y, w, h\}$ is the ground truth offset.

Similar to the candidate windows regression task, facial key points detection is seen as a regression problem and we minimize the SmoothL1 loss:

$$L_{pts}(y_i, y_i^*) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{else} \end{cases} \tag{4}$$

Where $y_i$ is the facial key points predicted by the network and $y_i^* = \{x_1, y_1, x_2, y_2, ..., x_5, y_5\}$ is the ground truth five facial landmarks positive.

## 2.3 Multi-task Training

It has been proved by Li *et al.*[2] that adding a calibration function CNN after detection function CNN can improve candidate window localization effectiveness and significantly reduce the number of candidate windows. However, calibration CNNs requires extra computational expense and cause cascade structure too complex to handle. In order to reduce the amount of calculation and simplified the structure of the cascade CNNs, we joint detection with calibration similar to Faster-RCNN.

In addition, Chen *et al.*[11] use a face alignment model detect the face key points and combine face alignment with face detection, the result has proved that aligned face shapes provide better performance to face detection. Since seperating key points location and face detection will cause more computational expense, we combine key points detection to provide

Table 1: Compare Large-batch training and Mini-batch training

| Evaluation / Network | 1000 Times Calculate(large/mini batch) | Accuarcy(large/mini batch) |
|---|---|---|
| 12net | 0.117s / 0.103s | 94.4% / 94.4% |
| 24net | 1.521s / 1.510s | 95.4% / 95.1% |
| 48net | 4.701s / 4.490s | 95.2% / 95.2% |

the cascade CNNs with information from different aspect, and this approach lead the model a better performance. The result can be shown in Figure 4(a).

In order to train this multi-loss CNN, for each iteration batch we have to select different training data for different loss and only backpropagrate NN with specific data type and restrain others. There should be four different kinds of training data(including *Positives*, *Negatives*, *Part faces* and *Key Points Detection* data) for three tasks. Negatives and positives are used for face vs non-face classification task, positives and part faces are used for candidate windows regression and landmark faces are used for facial landmarks detection. In this paper, Positives, Negatives, Part Faces data are generated on *Wider Face*[20] and key points detection data is generated on *CelebA*[21].

We can either a) have a large batch size( = 256) and randomly select different kinds of data for this batch or b) have several iteration mini-batchs( = 64) and for each mini-batch have same kind of training data. According to Table 1, training with mini-batchs can have less calculations, faster training speed and similar accuracy to large batch option. So we choose mini batch in our training approach.

## 2.4   Hard Example Proposal Network

Deep learning process has achieved significant advances because of big datasets, but still includes many heuristics and hyperparameters that are difficult to select. On the other hand, big datasets always contain an large number of easy examples and a small number of hard examples, using too many easy examples will not improve much to the result of our CNN models. Therefore, reference to object detection hem methods[8, 7], we design a *Hard Example Proposal Network* for cascade, which generate the hard training samples simultaneously for three seperate CNNs.

As shown in Figure 3, Hard Example Proposal Network regular generate rectangular object proposals, each of them have an object class according to the IoU. The proposals which have IoU ratio less than 0.3 to all ground truth is classified as Negatives, the proposals which have IoU ratio above 0.7 to a ground truth is classified as Positives, the proposals with IoU between 0.3 and 0.7 is recognized as Part faces. Hard Example Proposal Network also crop the ground truth and label five landmarks' position as Landmark faces.

In Hard Example Proposal Network, three different loss layers respectively compute loss for their training data, sort them from large to small based on loss. The network will automatically select the top N(*batch size* = N) as hard examples and set the loss of all non-hard examples to 0. If hard example is not enough, then randomly transform hard examples to fill in batches. Experiments shows that this stategy lead training converge faster and a better performance. It is proved to be surprisingly effective in accelerating convergence and reducing training loss Figure 4(b).
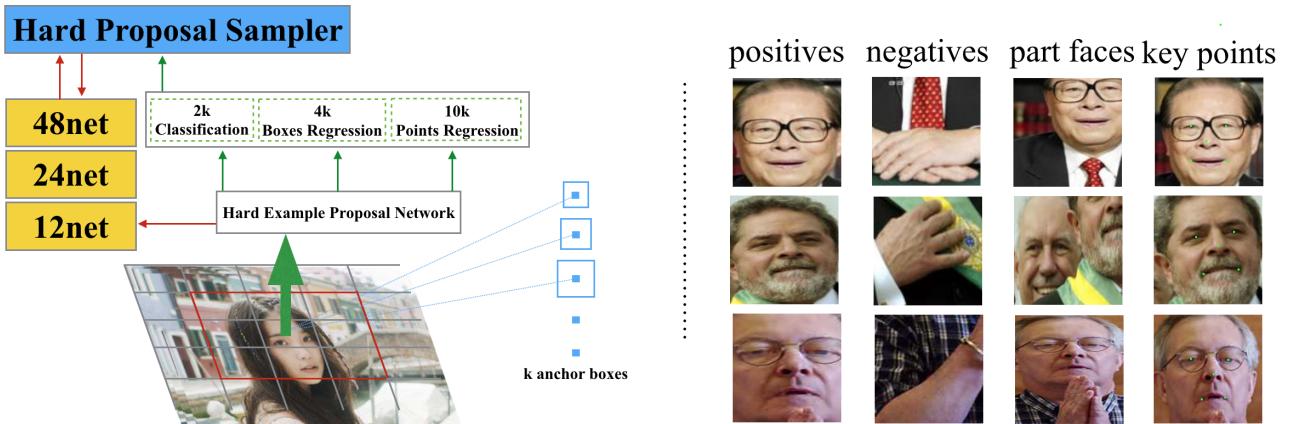
Figure 3: Hard Example Proposal Network generate hard proposals and labels during training process. According to loss, Hard Proposal Sampler choose only top 70% of training data backward propagate.

## 2.5 End To End Training.

Compared to other cascade convolutional neural network structure[11, 9, 3, 2, 4], we propose an *End-to-End* framework to train the cascade CNNs once a time. End-to-End training is to give the network model a raw input to obtain a final output, it can reduce the manual preprocessing and subsequent processing, the model can automatically adjust itself according to the data as much as possible, thereby increasing the overall fit of the model.

In ordinary cascade CNNs architecture[2], independent cascade CNN do not share weight and bias with other CNNs, it's impossible to update multi cascade CNNs during one forward propagate and backward propagate. In order to get rid of this limitation, as shown in Figure 2, we build a bridge between 12net and 24net, 24net and 48net. With this concatenate structure, the weight matrix and bias matrix is concatenated to higher cascade level CNN, the higher cascade level CNN is supplemented by the information from lower cascade level CNN which enable the detector to learn multi resolution information and offer a chance to backward propagate from 48net to 12net.

In this paper, an *Alternating Training* method is proposed to achieve End-to-End training. At the very beginning, 12net, 24net and 48net are trained independently, they modify their convolution layers in different ways. Then the convolution layer's weigh and bias parameters are extracted and duplicated to an End-to-End Cascade structure, and fine-tune them on same dataset with Hard Example Proposal Network. At this time, three seperate CNNs will share same weight, become more unified and robust.

# 3 Experiments

## 3.1 Effectiveness of Joint Training

To evaluate the contribution of joint training strategy on FDDB, we fix the 12net, 24net and modity the 48net to joint and non-joint versions. Figure 4(a) suggests that jointing face alignment with face detection task is beneficial for face detection.

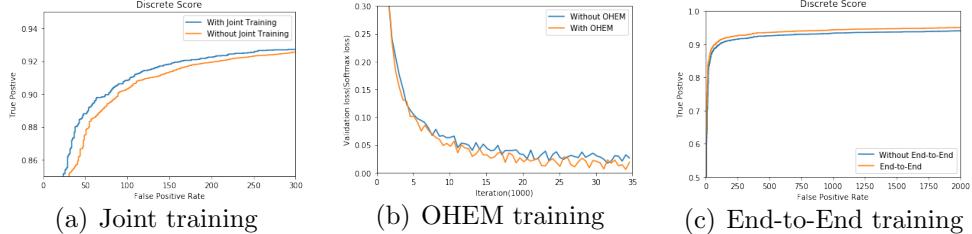(a) Joint training    (b) OHEM training    (c) End-to-End training

Figure 4: The figure shows the contribution of joint training strategy, online hard example mining strategy and end-to-end training strategy on FDDB[10] and AFLW[27] Benchmark

## 3.2 Effectiveness of OHEM in Hard Example Proposal Network

To evaluate the effectiveness of OHEM(online hard example mining) training strategy in Hard Example Proposal Network, we trained two 48net version(OHEM training version and Without-OHEM training version) on same dataset. Both of them have same learning rate($lr = 0.01$), batch size($batch\ size = 64$) and trained on Wider Face Training dataset[20]. Figure 4(b) represent the loss curves of two different training ways. According to the loss curves, OHEM can accelerate convergence and lead to a lower loss value. It is easy to draw a conclusion that Hard Example Proposal Network, especially OHEM, is beneficial to performance improvement.

## 3.3 Effectiveness of End-to-End

To evaluate the effectiveness of End-to-End training strategy, the same 12net, 24net and 48net initialize models are used to generate End-to-End and non End-to-End training versions. Figure 4(c) represent the face detection result of two different training strategies on FDDB. End-to-End structure can lead to a more unified and robust system. The ROC curve suggests that End-to-End is beneficial to cascade CNNs structure.
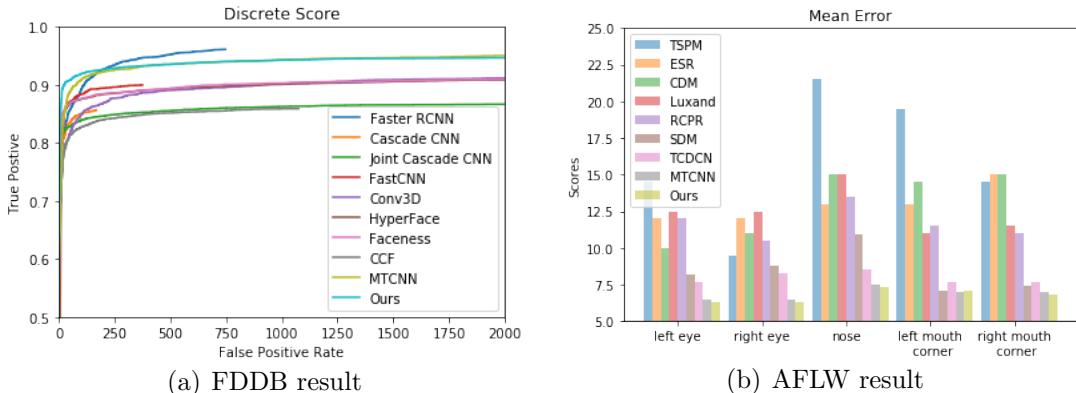


(a) FDDB result    (b) AFLW result

Figure 5: The figure shows the comparison of our model with other state-of-art methods [7, 2, 9, 11, 22, 23, 25, 24, 26] on FDDB and [13, 15, 16, 14, 12, 17, 18] on AFLW.

## 3.4 Evaluation on FDDB

FDDB is a face detection dataset and benchmark, it contains 5,171 annotated faces in 2,845 images. We compare our model with other state-of-the-art methods[7, 2, 9, 11, 22, 23, 25, 24, 26] on FDDB, the ROC result is shown in Figure 5(a).

## 3.5 Evaluation on AFLW

AFLW is a facial key points detection benchmark which contains 24,386 faces with facial landmarks. We compare our model with other state-of-the-art methods[13, 15, 16, 14, 12, 17, 18] on AFLW, the mean error result is shown in Figure 5(b).

# 4 Conclusion

In this paper, we have proposed an End-to-End Multi-Output Cascade CNNs and introduce details on how to achieve. Experiments demonstrate that our framework can cope with the complex environment in reality and perform well on several challenging benchmarks(including FDDB[10] for face detection, and AFLW[27] for face alignment). In the future, we will make more efforts on the correlation between face detection and object detection or other aspects to further improve the performance. The further improvement will be published in the future.

# References

[1] Viola P, Jones M J.: Robust real-time face detection. International journal of computer vision. 57, 137-154 (2004)

[2] Li H, Lin Z, Shen X, et al.: A convolutional neural network cascade for face detection. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 5325-5334 (2015)

[3] Sun Y, Wang X, Tang X.: Deep convolutional network cascade for facial point detection. Proceedings of the IEEE conference on computer vision and pattern recognition. 3476-3483 (2013)

[4] Zhou E, Fan H, Cao Z, et al.: Extensive facial landmark localization with coarse-to-fine convolutional network cascade. Proceedings of the IEEE International Conference on Computer Vision Workshops. 386-391 (2013)

[5] Yang S, Luo P, Loy C C, et al.: From facial parts responses to face detection: A deep learning approach. Proceedings of the IEEE International Conference on Computer Vision. 3676-3684 (2015)

[6] Long J, Shelhamer E, Darrell T.: Fully convolutional networks for semantic segmentation. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 3431-3440 (2015)

[7] Ren, Shaoqing, et al.: Faster r-cnn: Towards real-time object detection with region proposal networks. Advances in neural information processing systems. 91-99 (2015)

[8] Shrivastava, Abhinav, Abhinav Gupta, and Ross Girshick.: Training region-based object detectors with online hard example mining. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 761-769 (2016)

[9] Zhang K, Zhang Z, Li Z, et al.: Joint Face Detection and Alignment Using Multitask Cascaded Convolutional Networks. IEEE Signal Processing Letters. 23(10), 1499-1503 (2016)

[10] V. Jain, and E. G. Learned-Miller, FDDB: A benchmark for face detection in unconstrained settings, Technical Report UMCS-2010-009, University of Massachusetts, Amherst. (2010)

[11] Chen, Dong, et al.: Joint cascade face detection and alignment. European Conference on Computer Vision. Springer International Publishing. 109-122 (2014)

[12] Burgos-Artizzu, Xavier P., Pietro Perona, and Piotr Dollr.: Robust face landmark estimation under occlusion. Proceedings of the IEEE International Conference on Computer Vision. 1513-1520 (2013)

[13] Zhu, Xiangxin, and Deva Ramanan.: Face detection, pose estimation, and landmark localization in the wild. Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on. 2879-2886 (2012)

[14] Documentation, Luxand FaceSDK.: Luxand FaceSDK 4.0 Face Detection and Recognition Library. Developer s Guide, Copyright 2011. (2005).

[15] Cao, Xudong, et al.: Face alignment by explicit shape regression. International Journal of Computer Vision. 107.2, 177-190 (2014)

[16] Yu X, Huang J, Zhang S, et al.: Pose-free facial landmark fitting via optimized part mixtures and cascaded deformable shape model. Proceedings of the IEEE International Conference on Computer Vision. 1944-1951 (2013)

[17] Xiong, Xuehan, and Fernando De la Torre.: Supervised descent method and its applications to face alignment. Proceedings of the IEEE conference on computer vision and pattern recognition. 532-539 (2013)

[18] Zhang, Zhanpeng, et al.: Facial landmark detection by deep multi-task learning. European Conference on Computer Vision. Springer International Publishing. 94-108 (2014)

[19] Ren, Shaoqing, et al.: Face alignment at 3000 fps via regressing local binary features. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 1685-1692 (2014)

[20] Yang, Shuo, et al.: Wider face: A face detection benchmark. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 5525-5533 (2016)

[21] Liu, Ziwei, et al.: Deep learning face attributes in the wild. Proceedings of the IEEE International Conference on Computer Vision. 3730-3738 (2015)

[22] Triantafyllidou D, Tefas A.: A Fast Deep Convolutional Neural Network for Face Detection in Big Visual Data. Advances in Big Data. Springer International Publishing. 61-70 (2016)

[23] Li Y, Sun B, Wu T, et al.: Face Detection with End-to-End Integration of a ConvNet and a 3D Model. European Conference on Computer Vision. Springer International Publishing. 420-436 (2016)

[24] Ranjan R, Patel V M, Chellappa R. Hyperface: A deep multi-task learning framework for face detection, landmark localization, pose estimation, and gender recognition. arXiv preprint arXiv. 1603.01249 (2016)

[25] Yang B, Yan J, Lei Z, et al.: Convolutional Channel Features. IEEE International Conference on Computer Vision. 82-90 (2016)

[26] Yang S, Luo P, Loy C C, et al.: From Facial Parts Responses to Face Detection: A Deep Learning Approach. IEEE International Conference on Computer Vision. 3676-3684 (2016)

[27] Kstinger M, Wohlhart P, Roth P M, et al.: Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization. Computer Vision Workshops (ICCV Workshops), 2011 IEEE International Conference on. 2144-2151 (2011)