# A Jointly Learned Deep Architecture for Facial Attribute Analysis and Face Detection in the Wild

Keke He, Yanwei Fu, Xiangyang Xue
Fudan University
{kkhe15, yanweifu, xyxue}@fudan.edu.cn

## Abstract

*Facial attribute analysis in the real world scenario is very challenging mainly because of complex face variations. Existing works of analyzing face attributes are mostly based on the cropped and aligned face images. However, this result in the capability of attribute prediction heavily relies on the preprocessing of face detector. To address this problem, we present a novel jointly learned deep architecture for both facial attribute analysis and face detection. Our framework can process the natural images in the wild and our experiments on CelebA and LFWA datasets clearly show that the state-of-the-art performance is obtained.*

## 1. Introduction

The problem of automatically analyzing the facial attributes received increasing attention recently [24, 15, 37]. Face attributes may potentially benefit a number of real-world applications, such as face alignment[22, 26, 36, 38], head pose estimation [40, 32] and face verification [28, 27]. Analyzing facial attributes still remains challenging in real-world scenarios. For example, most existing works [24] predict the facial attributes on well-cropped face images. In an extension to real-world scenarios, those works have to utilize the face detector to localize the bounding box of human faces before facial attribute analysis.

Such a pipeline of face detection followed by facial attribute prediction is nevertheless undesirable in real-world facial attribute prediction due to two reasons. Firstly, with the preprocessing of face detector, the capability of facial attribute prediction has to heavily relies on the results of face detection. Unfortunately in an uncontrolled setting [15], the face is likely to have large out-of-plane tilting, occlusion and illumination variations, which may affect the algorithms of face detection and facial attribute analysis simultaneously. For example, the Viola-Jones face detector [29] works well for near-frontal faces, but less effective for faces in the pose far from frontal views. Second, face detection

is heterogeneous but subtly correlated with facial attribute prediction, and vice versa [33]. Efficiently and effectively exploiting this correlation may help in both detecting faces and predicting facial attributes.

We propose an end-to-end deep architecture to *jointly* learning to detect faces and analyze facial attributes. Essentially, the two tasks share the same bottom layers (feature map in Fig. 1) in the architecture to alleviate the problems of face variations aforementioned. For instance, one face region of one image is failed to be detected as a face due to the partial occlusion, and yet the network can detect the existence of some facial attributes such as "eyeglasses", and "smiling". The face region can still be used to optimize our network and in turn help supervise the learning of the face detection part, to improve the performance of face detection part.

In this work, the deep learning architecture is proposed for both face detection and facial attribute analysis. The network structure is illustrated in Fig. 1. Given an entire image, our architecture will firstly pass the image with convolutional layers (*conv1– conv5*) and max pooling layers to produce the *conv* feature map for face region proposals. The region of interest (RoI) pooling layer and face region proposal layer are two layers introduced to facilitate the jointly learning of our architecture. A region of interest (RoI) pooling layer pools each face proposal on the *conv5* feature map into a fixed-length feature vector which is further processed by two fully connected layers (*fc6* and *fc7*). Built on the *fc7* layer, each individual task of facial attribute analysis and face detection are jointly optimized. By virtue of the face detection subnet, our architecture can directly predict the facial attributes from the whole images, rather than using the well-cropped images as previous work [20, 24, 31]. Extensive experiments on benchmark CelebA [15] and LFWA [9, 15] datasets demonstrate that our method outperforms state-of-the-art alternatives.

**Contribution**. Our main contribution is to propose a novel deep architecture of jointly learning of face detection and facial attribute tasks, which is capable of utilizing both tasks

1

to better optimize the shared network and thus improve the performance of both tasks. Two layers – region of interest (RoI) pooling layer and face region proposal layer are introduced to enable learning these two tasks simultaneously. More importantly, we can predict the facial attributes in the wild and the input images do not need to be cropped and aligned as the standard practice in [24, 31]. Finally, our joint deep architecture achieves state-of-the-art result on the biggest benchmark dataset for facial attribute analysis–CelebA dataset [15] and LFWA [15, 9] datasets.

## 2. Related Work

**Facial attribute analysis.** It was first studied by Kumar et al. [13]. In term of different visual features and distinctive learning paradigm, the facial attribute analysis has been developed into three categories: (1) the methods [13] of using hand-crafted visual features, such as SIFT [16] and LBP [18]; (2) the methods of utilizing the recent deep features [15, 37, 31]; and (3) multi-task methods of learning facial attribute [1, 24]. We here highlight the differences between our architecture and these previous works. Liu *et al.* [15] cascaded three deep networks pre-trained for facial attribute prediction. In contrast, we show that the tasks of face detection and facial attribute prediction are highly correlated and our jointly learning architecture can improve both tasks. Rudd *et al.* [24] introduced a mixed objective optimization network which utilizes distribution of attribute labels to learn each task. Abdulnabi *et al.* [1] proposed a multi-task CNN model sharing of visual knowledge between tasks for facial attribute analysis. Comparing with [24, 1], we focus on jointly learning the face detection and facial attribute analysis; and our model can predict facial attributes on the images in the wild [24, 31] *i.e.* without cropped and aligned, as illustrated in Fig. 1 and Fig. 5.

**Face detection.** The work of Viola and Jones [29, 30] made face detection usable in many real world applications. The cascaded classifiers were built on Harr-like features to detect human faces. The deformable part model (DPM) [6] on top of HOG feature is a general object detector and can also be used for face detection [4]. Recent advances of deep learning architectures also inspire another category of methods for face detection. Yang [34] used the fully convolutional networks (FCN) to generate the heat map of facial parts for producing face proposals. In contrast to these works, our face detection task gets benefit from not only the well-designed architecture (in Fig. 1), but also the jointly learning process with facial attribute prediction. Further inspired by Fast-RCNN [7, 23], we take the face detection as a special case of the general semi-rigid object detection. More specifically, given an image, our face detection will try to answer two questions: (1) whether this patch contains faces or not? *i.e.* face score task in Fig. 1. (2) can we de-

tect the bounding box of faces if there is any face? *i.e.* face bounding box task in Fig. 1. The jointly optimizing these two sub-tasks will better solve face detection.

**Multi-task learning.** Our framework can be categorized as multi-task learning [35, 19], which shares the information and explores the similarity of related tasks on the same data. The multi-task learning can facilitate a wide range of tasks and applications, including but not limited to action recognition [39], information retrieval [14], facial landmark detection [36, 20], and facial attribute prediction [25, 5, 1, 21, 17]. The recent technical report [20] also combine face detection with the tasks of locating face landmarks and recognizing gender. However, unlike our facial attribute prediction on the images in the wild, their work still has to firstly detect face regions for the facial landmarks and predicting gender.

## 3. Our Deep Architecture

In this section, we firstly overview our network in Sec. 3.1, and then the tasks are defined in Sec. 3.2. The Face region proposal and RoI pooling layers are explained accordingly in Sec. 3.3 and Sec. 3.4. Finally, we utilize the network to solve the facial attribute prediction and face detection in Sec. 3.5.

### 3.1. Overview

Figure 1 shows our framework for jointly face detection and attribute analysis. Our architecture takes an entire image as input and a set of face bounding boxes as labels for training. The whole network firstly processes the image with several convolutional layers (*conv1– conv5*), and max pooling layers to produce a *conv* feature map for face region proposal. For feature map of each proposed face, we employ a region of interest (RoI) pooling layer to pool it into a fixed-length feature vector. Each feature vector is further processed by two fully connected layers (*fc6* and *fc7*) and thus used for the tasks of facial attribute analysis and face detection. On top layers, our architecture has *face detection branch* and *facial attribute branch*.

Our network follows the art and design of VGG-16 [3]. Particularly, the kernel size, stride and the number of filters in convolutional layers (*conv1– conv5*) and the two fully connected layers (*fc6* and *fc7*) are exactly the same as the corresponding layers in VGG-16 architecture.

### 3.2. Task formulation

Suppose we have the labelled source training dataset $\mathcal{D}_s = \{\mathbf{I}, \mathbf{a}, \mathbf{L}\}$ with $N$ training instances and $M$ attributes. $\mathbf{I}$ denotes the patches of training images and $\mathbf{L}$ denotes the labels. We use the $\mathbf{L}$ matrix to both denote whether an image patch contains the face, and whether a facial attribute exists in the image patch. Particularly, for the $i$-th image
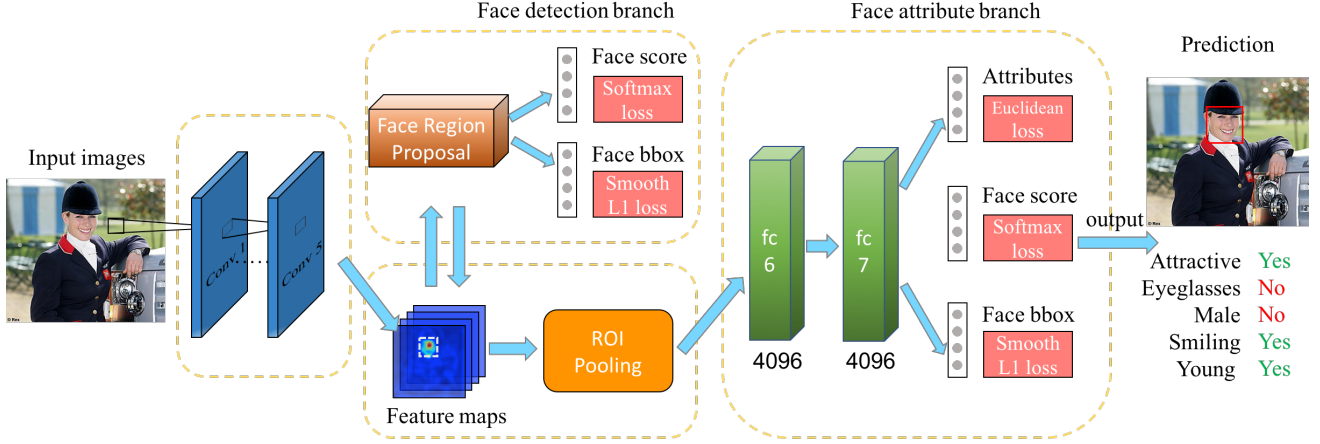
Figure 1. Overview of the proposed architecture.

patch $\mathbf{I}_i$ ($i = 1, \cdots, N$), we use $\mathbf{L}_{i\star} = \mathbf{0}$ to indicate that this image patch does not contain a human face. If $\mathbf{I}_i$ is a face image patch, we use $\mathbf{L}_{ij} = +1$ to denote the existence of $j-$th facial attribute $\mathbf{a}_j$ ($j = 1, \cdots, M$); $\mathbf{L}_{ij} = -1$ otherwise.

As illustrated in Fig. 1, our network extracts the image patch to a $4096 - dim$ feature vector, and it is denoted as $f(\mathbf{I}; \Theta)$, where $\Theta$ is the parameter set of the deep architecture. We have the prediction tasks of facial attribute, face score and face bounding box. Each task has their own parameters on the last layer.

**Facial attribute prediction.** To predict the attributes of the image $\mathbf{I}^*$, we need to learn a function $\mathbf{L}^*_{i\mathbf{a}} = \Psi(\mathbf{I}^*)$ to predict the facial attribute $\mathbf{a}$. We thus have the predicting function $\Psi = [\psi_i]_{i=1,\cdots,M}$, and $\psi_i(\mathbf{I}) : \mathbb{R}^{4096} \to \{+1, -1\}$; specifically, we consider $\Psi(x) = W_a^T x$, where $W_a \subseteq \mathbb{R}^{4096 \times 40}$ for all 40 attributes.

**Face score.** This task aims to predict the score whether an image patch is face. Essentially, we need to learn $\Phi(x) : \mathbb{R}^{4096} \to \{+1, -1\}$, and $\Phi(x) = W_s^T x$, where $W_a \subseteq \mathbb{R}^{4096 \times 2}$ for binary prediction task.

**Face bounding box.** We regress the bounding box of faces by $\Omega(x) : \mathbb{R}^{4096} \to \mathbb{R}^4$, and we configure the form as $\Omega(x) = W_b^T x$, and $W_b \subseteq \mathbb{R}^{4096 \times 4}$ for the new bounding box.

### 3.3. Face region proposal layer

Inspired by the work of object detection [23], our network has the branch of face region proposal. Specifically, this branch takes an image of any size as input and outputs a set of rectangular face proposals. To generate the facial region proposal, we slide the branch over the convolutional feature map after the five convolutional layers (*conv1– conv5*) in Fig. 1. The extracted *conv5* features are further fed into two sibling fully connected layers, i.e. *face*

*bounding box* (bbox) and *face score* layer. Face bounding box layer employs the smooth $L_1$ loss (defined in [7]) enables a regressor to predict the facial bounding box. Face score layer utilizes the softmax loss to indicate whether the bounding box is a face. For each image, we selected top-300 face region proposals in term of the face scores computed.

### 3.4. The RoI pooling layer

The RoI pooling layer can convert the feature maps of face region proposal (with the size of $h \times w$) to a fixed spatial extent of $H \times W$, which facilitates the further process. Here, $h, w, H, W$ are heights and widths of each rectangle region of feature map in Fig. 1. For the varying size of input feature patches, we vary the size of filters of pooling layer with the sub-window of approximate size $h/H \times w/W$. The RoI pooling technique is adopted from the work of object detection [7] which inspired by the SPPnets [8].

### 3.5. Facial attribute analysis and face detection

The fixed-length feature vector extracted from RoI Pooling layer is further fed into two fully connected layers (*fc6* and *fc7*) in order to enable the tasks of facial attribute analysis. Particularly, built on the layers of *fc6* and *fc7*, the multi-task network structure is employed to analyze each attribute individually. This can be modeled as the minimization of the expected loss over all the training instances which is

$$\{\Theta, W_a, W_s, W_b\} = \operatorname{argmin} \mathcal{L}(\mathbf{I}; \Theta, W_a, W_s, W_b) \quad (1)$$

where $\mathcal{L}(\mathbf{I}; \Theta, W_a, W_s, W_b)$ is the loss function of jointly learning; and we have

$$\mathcal{L}(\mathbf{I}; \Theta, W_a, W_s, W_b) = \lambda_1 \mathcal{L}_a(\Phi(f(\mathbf{I}))) + \lambda_2 \mathcal{L}_s(\Psi(f(\mathbf{I}))) \quad (2)$$

$$+ \lambda_3 \mathcal{L}_b(\Omega(f(\mathbf{I})))$$

3

here the loss function $\mathcal{L}_a \left( \Phi \left( f \left( \mathbf{I} \right) \right) \right)$ on *facial attribute prediction* is the mean square error loss of all attributes; we use softmax loss $\mathcal{L}_s \left( \Psi \left( f \left( \mathbf{I} \right) \right) \right)$ for *face score task,* and finally the smooth $L_1-$ loss $\mathcal{L}_b \left( \Omega \left( f \left( \mathbf{I} \right) \right) \right)$ [7] is employed to regress the *face bounding box.*

We denote the shared parameters of deep architectures as $\Theta$ which are jointly optimized by all these three tasks. Since these three tasks are highly related, learning in such a way can not only greatly reduce the prediction error of each individual task, but also accelerate the convergence rate of learning the whole network.

For the testing image $\mathbf{I}_k$, the predicted result $\hat{\mathbf{L}}_{kj}$ of the $j-$th facial attribute $\mathbf{a}_j$ $(j = 1, \cdots, M)$ is thresholded by

$$\hat{\mathbf{L}}_{kj} = \begin{cases} 1 & \psi \left( \mathbf{I}_k \right) > \tau \\ -1 & \psi \left( \mathbf{I}_k \right) \leq \tau \end{cases} \qquad (3)$$

where $\tau$ is the threshold parameter.

## 4. Experiments

### 4.1. Datasets and settings

We conduct the experiments on the CelebA dataset [15] and LFWA dataset [9, 15].

**CelebA [15]** contains approximately 200k images of 10k identities. Each image is annotated with 5 landmarks (two eyes, the nose tips, the mouth corners) and binary labels of 40 attributes. To make a fair comparison with the other facial attribute methods, the standard split is used here: the first $160k$ images are used for training, $20k$ images for validation and remaining $20k$ for testing. CelebA provides two types of training images, *i.e.*, aligned and cropped face images and raw images as shown in the first and second row of Fig. 4 respectively. We use the raw images for training our joint deep architecture. Since there is no ground-truth bounding box label for faces in raw images, the face detector [4, 11] is employed here to help generate the face bounding box for training images. We use the implementation of dlib toolbox [11, 4] for generating labels. Note that face bounding box generation is only required in the training stages for synthesizing labels. At testing phrase, raw images are directly input for both facial attribute analysis and face detection.

**LFWA [15]** is constructed based on face datasets LFW [9]. It contains approximately 13143 images of $10k$ identities. Following [15], 50% of the images for training, and the other 50% are used for testing. LFWA has 40 binary facial attributes, the same as CelebA. We also generate face bounding box for LFWA to train our joint learning network.

**Evaluation metrics.** We take the attribute prediction as classification tasks and thus mean accuracy can be computed. Particularly, we evaluate the performance as comparable to [24, 15] by the mean error which is defined as

$mean\,error = 1 - mean\,accuracy.$

**Implementation and Parameter settings.** The $\tau$ is set as 0 in Eq (3); and square error loss is used in Eq (1). We empirically set $\lambda_1 = \lambda_2 = 1, \lambda_3 = 2$ in Eq (2). The convolutional layers and fully connected layers are initialized by the 16 layer VGG network [3] individually. We use the open source deep learning framework Caffe [10] to implement our structure. A single end-to-end model is used for all the testing. We employ the stochastic gradient descent to train our network. Dropout is used for fully connected layers and the ratio is set to 0.5. For training CelebA dataset, with initial learning rate 0.001, and gradually decreased by $1/10$ at $100k$, $150k$ iterations, the total training iterations are $180k$. For training LFWA dataset, as there are only limited training images - 6263, We fine-tune from pre-trained CelebA model, with initial learning rate 0.0001, and decreased by $1/10$ at $40k$ iterations, the total training iterations is $60k$. Once trained, our framework can predict the facial attributes and detect faces on the images in the wild for any testing images.

**Running cost.** Our facial model get converged with $180k$ iterations and it takes 29 hours on CelebA with one NVIDIA TITANX GPU. On LFWA dataset, Our facial model get converged with $60k$ iterations and takes 9 hours. For training all the model, and it takes around 4 GB GPU memory.

### 4.2. Competitors.

Our model is compared against state-of-the-art methods and several baselines. Particularly, (1) **FaceTracer** [12] is one of the best methods with the hand-crafted features. The features used including HOG [2] and color histograms of facial regions of interested to train SVM classifier for predicting facial attributes. (2) **LNets+ANet** [15] migrates two deep CNN face localization networks to one deep CNN network for facial attribute classification. (3) **Walk and Learn** [31] learns good representations for facial attributes by exploiting videos and contextual data (geo-location and weather) as the person walks. (4) **Moon** [24] is a mixed objective optimization multi-task network to learn all facial attributes, achieves the best result on CelebA dataset. (5) **Cropped model** is a variant of our model without using the face detection branch. The raw image is also used as the input and we crop the faces from the images by Dlib toolkit to train facial attribute models. The processed face images are utilized to train the model. (6) **Aligned model** uses the aligned and cropped images provided by CelebA and LFWA to train facial attribute model, which is used by most of the state-of-the-art methods. To make a fair comparison, these two baseline models are initialized by 16 layer VGG network; same as the proposed structure.
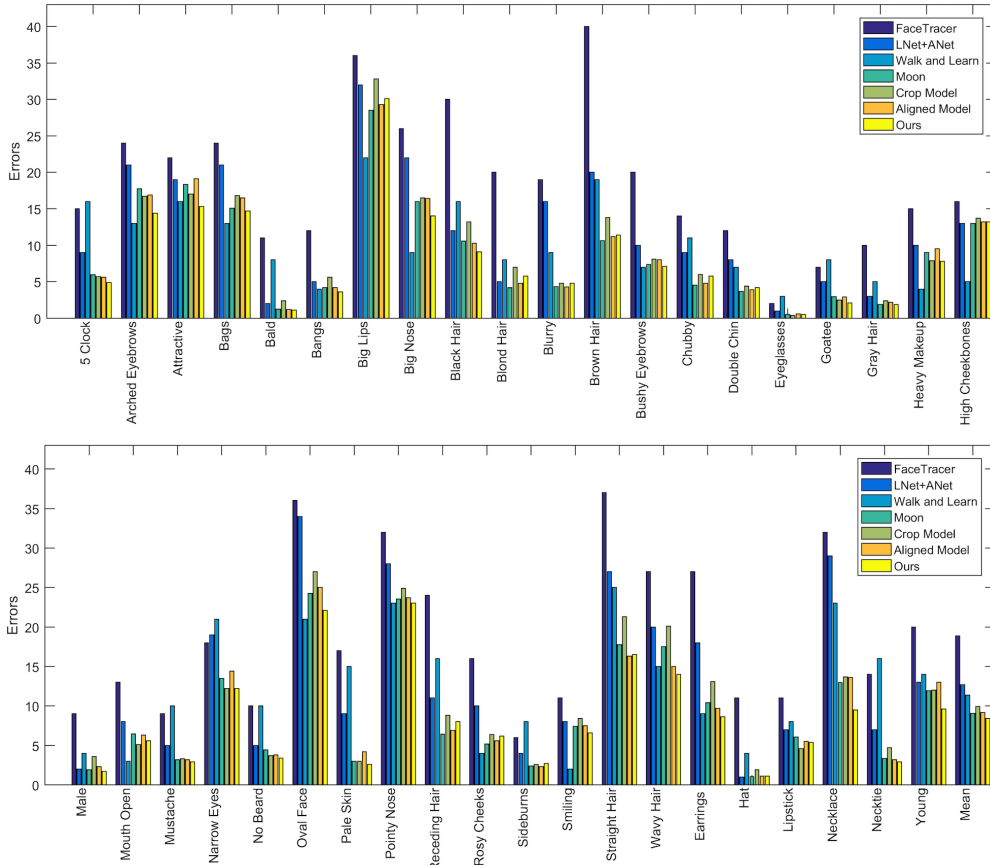
Figure 2. Performance comparison with state-of-the-art methods on CelebA on all 40 attributes. X-axis indicates each attribute, while y-axis is the error rate of attribute prediction (the lower value, the better performance of one method). The averaged error rate of FaceTracer, LNets+ANet, Walk and Learn, Moon, Baseline (Cropped images), Baseline (Aligned images), and ours are 18.88%, 12.70%, 11.35%, 9.06%, 9.95%, 9.19% and **8.41%** respectively.

| Methods | CelebA (%) | LFWA (%) |
|---|---|---|
| FaceTracer [12] | 18.88 | 26.07 |
| LNets+ANet [15] | 12.70 | 16.15 |
| Walk-and-Learn [31] | 11.35 | 13.40 |
| Moon [25] | 9.06 | - |
| Baseline: Cropped Model | 9.95 | 15.99 |
| Baseline: Aligned Model | 9.19 | 15.06 |
| Ours | **8.41** | **13.13** |

Table 1. Comparison of mean error on CelebA and LFWA datasets.

## 4.3. Comparison with the State of the Art

We compared our model with state-of-the-art methods: FaceTracer [12], LNets+ANet [15], Walk and Learn [31], Moon [24] and two baselines. The results on the testing split of CelebA and LFWA are reported in Tab. 1, Fig. 2 and Fig. 3, respectively. Note that please refer to the supplementary material for the full comparison results on each attribute. Comparing with all the competitors, we draw the following conclusions.

**Our model beats all the other methods by the mean error on CelebA and LFWA dataset.** As we can see from Fig. 2, our approach obtains the mean error **8.41%** which is the lowest among all the competitors, and it outperforms all the other methods on CelebA dataset. And on LFWA dataset as shown in Fig. 3, our mean error is only **13.13%** which is the lowest among all the competitors. This validates the efficacy of our joint learning architecture. Particularly, We compare the classification error in each individual attribute in Fig. 2 and Fig. 3. We note that on more than half among the total 40 attributes, our framework is significantly better than the other competitors, since our jointly learning architecture can efficiently leverage the information between face detection and facial attribute prediction. Particularly, comparing with the other works, all our attribute prediction tasks share the same deep learning architecture, *i.e. conv1-conv5*, feature map and *fc6, fc7* layers as illustrated in Fig. 1). These shared structures implicitly model the correlations between each attribute task. Furthermore,
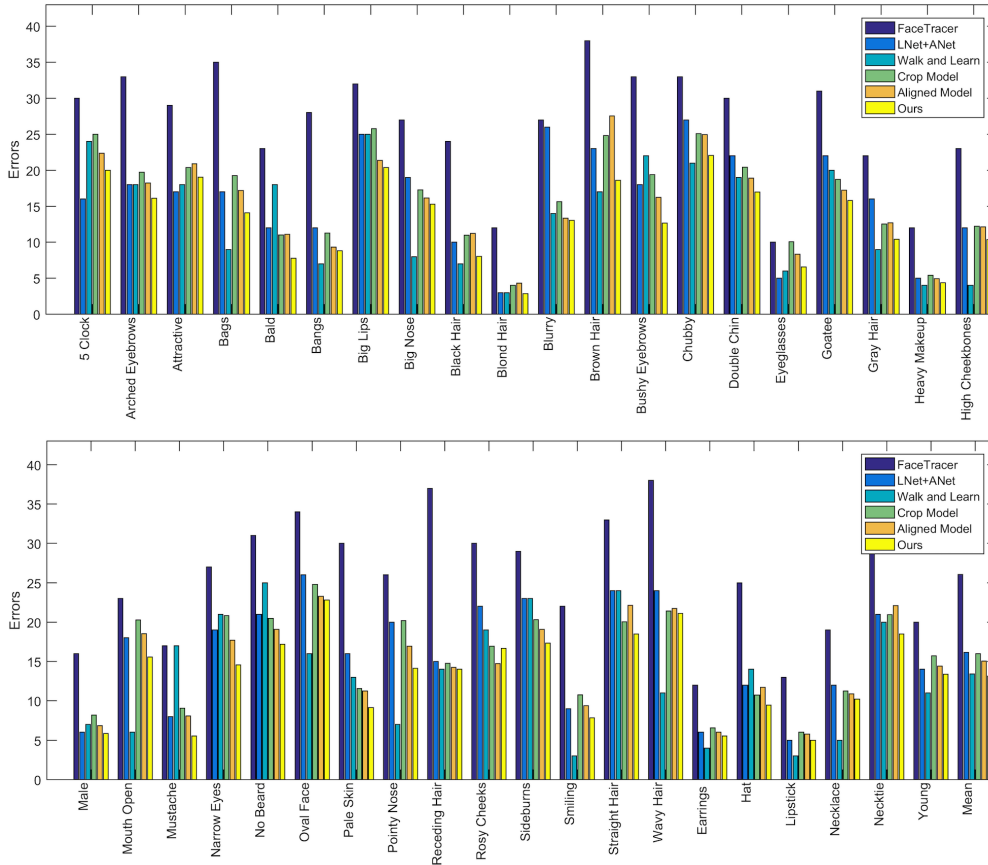
Figure 3. Performance comparison with state-of-the-art methods on LFWA on all 40 attributes. The averaged error rate of FaceTracer, LNets+ANet, Walk and Learn, Baseline (Cropped images), Baseline (Aligned images), and ours are 26.07%, 16.15%, 13.40%, 15.99%, 15.06% and **13.13%** respectively.



Figure 4. The first row is aligned images provided by the CelebA. The second row is CelebA raw images. We train our end-to-end model on raw images.

on gender attribute prediction, our framework can achieve 1.7% error on CelebA; in contrast, the error of HyperFace [20] is 3.0%. That indicates that our results can get 1.3% improvement over that of HyperFace [20].

**Our model can process the images in the wild.** Our face region proposal and RoI pooling layers are flexible enough to directly process the images in the wild. This thus better

demonstrates the effectiveness of our models. Specifically, our model not only achieves the best performance on the benchmark dataset – CelebA and LFWA, but also our tasks of facial attribute prediction do not need to align and crop the facial images as have done in many previous work [24, 12, 15]. Additionally, unlike the work of Walk and Learn [31], our model is not trained by the external data; and yet still obtains better results than Walk and Learn [31].

**Our model is very efficient in term of jointly learning to detect faces and predict facial attributes.** Our joint learning results have greatly improved over those of two baseline models – Cropped, Aligned models. These two methods are yet another two naive baselines of directly learning the facial attribute tasks. Particularly, on CelebA dataset, our model beats the two methods on the classification error of 27 attributes (totally 40 facial attributes) as compared in Fig. 2. Figure 3 shows that our joint model hit better results than the two baseline models on all the 40 attributes. This reveals attribute prediction tasks get benefit from the face detection task. Also note that unlike LNets+ANet [15] using two networks to localize the face and another one network
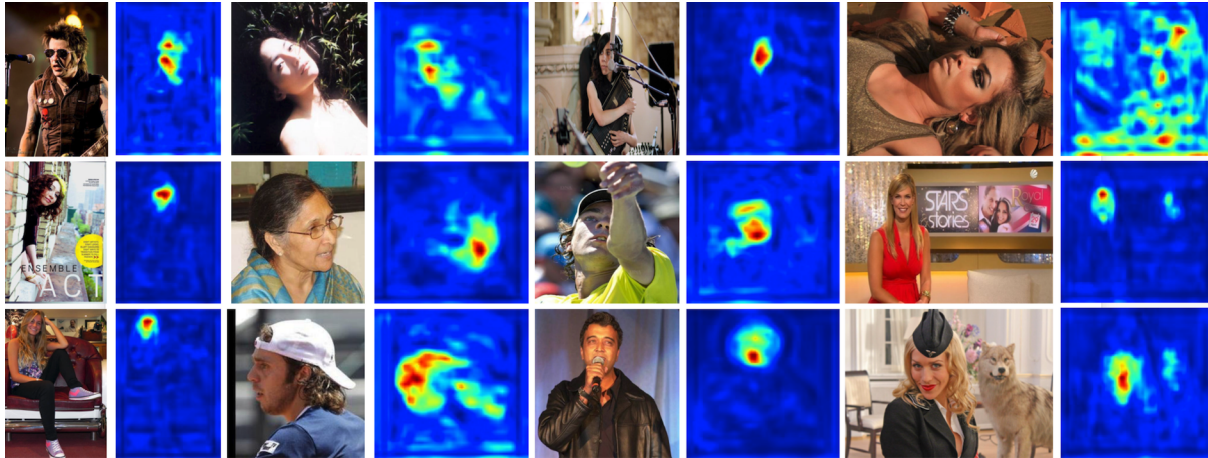
Figure 5. The visualization of feature map for *conv5* layer. The images are from CelebA test set. The different types of images are shown in the first three columns. These images have large pose variation of heads and high occlusion. Some failed examples are shown in the last column.

to extract features, we can train our model by a single end-to-end network here since the shared architectures (*conv1* – *conv5*) can explicitly model the important facial parts, again thanks to the face detection subnet.

**Qualitative results.** The important facial parts extracted by the shared feature map are visualized in Fig. 5. In particular, four groups of images are shown; and the feature map has higher activation on the regions of human faces. The first 3 columns are success cases, even thought the images have large pose (Column 2) or very high occlusion (Column 3). We also list some failure examples in Column 4, which are some extreme cases. Some of them are caused by ambiguous or too small view (e.g. the cat face is also similar to human's). The feature map reveals our model pay attention on the important facial parts which helps to analyze facial attributes.

### 4.4. Results of Face Detection

Our face detection is compared against the face detector of dlib toolbox [11], which is an implementation of the generic object detector [4] on face images. On CelebA dataset, there are only 5 landmarks of faces, but no labeled ground-truth bounding box of faces for raw images. We generate the ground truth bounding box by 5 landmark coordinates for the test set. If the IOU (Intersection over Union) of predicted bounding box and the ground truth bounding box is larger than 0.5, we assume the prediction is correct. We compare the two detectors on the testing split of CelebA.

The results are shown in Fig. 6. We compare the Precision-Recall curve for two methods. We find that both the Dlib detector and our face detector can achieve very high face detection on the CelebA test set as shown in the left subfigure of Fig. 6, which shows the efficacy of both de-
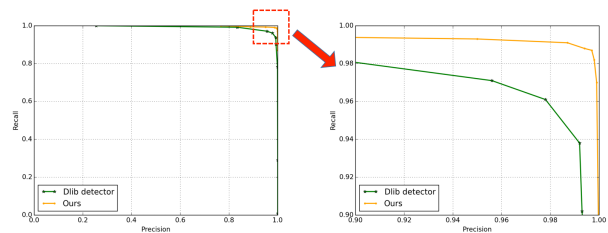


Figure 6. Precision-Recall (PR) results of face detection.

tectors of solving the task of face detection. Nevertheless, our face detector still beats the Dlib detector by a relatively large margin. The AUC values of Dlib detector and ours are 0.938 and 0.982 respectively. To better show the difference, we highlight the up-right corner of PR curve in the right subfigure as shown in Fig. 6. This result further proves that our jointly learned architecture can make the face detection and face attribute tasks help each other via sharing the same parameters of the deep network. This shows that attribute learning does help for face detection also it provides more detailed information about the face.

## 5. Conclusion

In this paper, we propose a novel joint deep architecture for facial attribute prediction and face detection. Different from the previous pipeline of face detection followed by facial attribute prediction, our architecture takes an entire image as input, enables both face detection and facial attribute analysis. The proposed architecture can not only exploit the correlation of face detection and face attribute prediction, but also boost both tasks. The experimental results on CelebA and LFWA datasets show the efficacy of proposed methods over the other state-of-the-art methods.

# References

[1] A. H. Abdulnabi, G. Wang, J. Lu, and K. Jia. Multi-task cnn model for attribute prediction. *IEEE TMM*, 2015.

[2] A. Bosch, A. Zisserman, and X. Munoz. Representing shape with a spatial pyramid kernel. In *CVIR*, 2007.

[3] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the devil in the details: Delving deep into convolutional nets. In *BMVC*, 2014.

[4] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *CVPR*, 2005.

[5] M. Ehrlich, T. J. Shields, T. Almaev, and M. R. Amer. Facial attributes classification using multi-task representation learning. In *CVPRW*, 2016.

[6] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *IEEE TPAMI*, 2010.

[7] R. Girshick. Fast r-cnn. In *ICCV*, 2015.

[8] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In *ECCV*, 2014.

[9] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical report.

[10] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.

[11] D. E. King. Dlib-ml: A machine learning toolkit. *JMLR*, 2009.

[12] N. Kumar, P. Belhumeur, and S. Nayar. Facetracer: A search engine for large collections of images with faces. In *ECCV*, 2008.

[13] N. Kumar, A. C. Berg, P. N. Belhumeur, and S. K. Nayar. Attribute and simile classifiers for face verification. In *ICCV*, 2009.

[14] X. Liu, J. Gao, X. He, L. Deng, K. Duh, and Y.-Y. Wang. Representation learning using multi-task deep neural networks for semantic classification and information retrieval.

[15] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *ICCV*, 2015.

[16] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 2004.

[17] Y. Lu, A. Kumar, S. Zhai, Y. Cheng, T. Javidi, and R. Feris. Fully-adaptive feature sharing in multi-task networks with applications in person attribute classification. In *CVPR*, 2017.

[18] T. Ojala, M. Pietikäinen, and T. Mäenpää. Gray scale and rotation invariant texture classification with local binary patterns. *IEEE TPAMI*, 2002.

[19] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010.

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

[21] R. Ranjan, S. Sankaranarayanan, C. D. Castillo, and R. Chellappa. An all-in-one convolutional neural network for face analysis. In *arxiv*, 2016.

[22] S. Ren, X. Cao, Y. Wei, and J. Sun. Face alignment at 3000 fps via regressing local binary features. In *CVPR*, 2014.

[23] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015.

[24] E. M. Rudd, M. Günther, and T. E. Boult. Moon: A mixed objective optimization network for the recognition of facial attributes. In *ECCV*, 2016.

[25] E. M. Rudd, M. Gunther, and T. E. Boult. Moon:a mixed objective optimization network for the recognition of facial attributes. In *ECCV*, 2016.

[26] C. Sagonas, G. Tzimiropoulos, S. Zafeiriou, and M. Pantic. 300 faces in-the-wild challenge: The first facial landmark localization challenge. In *ICCVW*, 2013.

[27] Y. Sun, Y. Chen, X. Wang, and X. Tang. Deep learning face representation by joint identification-verification. In *NIPS*, 2014.

[28] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *CVPR*, 2014.

[29] P. Viola and M. Jones. Rapid object detection using a boosted cascade of simple features. In *CVPR*, 2001.

[30] P. Viola, J. C. Platt, C. Zhang, et al. Multiple instance boosting for object detection. In *NIPS*, 2005.

[31] J. Wang, Y. Cheng, and R. Schmidt Feris. Walk and learn: Facial attribute representation learning from egocentric video and contextual data. In *CVPR*, 2016.

[32] H. Yang, W. Mou, Y. Zhang, I. Patras, H. Gunes, and P. Robinson. Face alignment assisted by head pose estimation. *arXiv*, 2015.

[33] S. Yang, P. Luo, C.-C. Loy, and X. Tang. From facial parts responses to face detection: A deep learning approach. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3676–3684, 2015.

[34] S. Yang, P. Luo, C. C. Loy, and X. Tang. From facial parts responses to face detection: A deep learning approach. In *ICCV*, 2015.

[35] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *Advances in neural information processing systems*, pages 3320–3328, 2014.

[36] Z. Zhang, P. Luo, C. C. Loy, and X. Tang. Facial landmark detection by deep multi-task learning. In *ECCV*, 2014.

[37] Y. Zhong, J. Sullivan, and H. Li. Face attribute prediction using off-the-shelf cnn features. In *arxiv*, 2016.

[38] E. Zhou, H. Fan, Z. Cao, Y. Jiang, and Q. Yin. Extensive facial landmark localization with coarse-to-fine convolutional network cascade. In *ICCVW*, 2013.

[39] Q. Zhou, G. Wang, K. Jia, and Q. Zhao. Learning to share latent tasks for action recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2264–2271, 2013.

[40] X. Zhu and D. Ramanan. Face detection, pose estimation, and landmark localization in the wild. In *CVPR*, 2012.