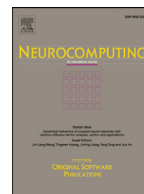




Contents lists available at ScienceDirect

Neurocomputing

journal homepage: www.elsevier.com/locate/neucom

A fast face detection method via convolutional neural network

Guanjun Guo^a, Hanzi Wang^{a,*}, Yan Yan^a, Jin Zheng^b, Bo Li^b^aFujian Key Laboratory of Sensing and Computing for Smart City, and School of Information Science and Engineering, Xiamen University, Xiamen 361005, PR China^bBeijing Key Laboratory of Digital Media, and School of Computer Science and Engineering, Beihang University, Beijing 100191, PR China

ARTICLE INFO

Article history:

Received 28 July 2017

Revised 9 December 2017

Accepted 4 February 2018

Available online xxx

Keywords:

Fast face detection

Convolutional neural network

Discriminative complete feature maps

ABSTRACT

Current face or object detection methods via convolutional neural network (such as OverFeat, R-CNN and DenseNet) explicitly extract multi-scale features based on an image pyramid. However, such a strategy increases the computational burden for face detection. In this paper, we propose a fast face detection method based on discriminative complete features (DCFs) extracted by an elaborately designed convolutional neural network, where face detection is directly performed on the complete feature maps. DCFs have shown the ability of scale invariance, which is beneficial for face detection with high speed and promising performance. Therefore, extracting multi-scale features on an image pyramid employed in the conventional methods is not required in the proposed method, which can greatly improve its efficiency for face detection. Experimental results on several popular face detection datasets show the efficiency and the effectiveness of the proposed method for face detection.

© 2019 Elsevier B.V. All rights reserved.

1. Introduction

In the Marr's theory of vision [1], a representational framework for vision was proposed, which consists of three levels: primal sketch, 1.5-D sketch and 3-D model representation. How to extract a proper image representation corresponding to the three levels is a critical problem in computer vision. So far, there are some milestone representations during the development of computer vision. For example, SIFT [2] and HOG [3] features have exhibited the good property of local invariance. Canny features [4] can capture the 1.5-D sketch, which represents the low-level image structure. The deconvolutional network [5] tries to characterize both mid-level and high-level representations in an unsupervised manner to understand images. However, finding a high-level representation for different computer vision tasks is still challenging.

Deep learning models have demonstrated impressive performance for different computer vision applications [6–11]. The deep convolutional neural network (CNN) [12] can map raw data from a manifold to the Euclidean space, in which features may be linearly separable. Generally speaking, there are two ways to extract features with deep CNN models. The first way is to extract features on an image pyramid obtained by using object proposals generation methods (such as SelectiveSearch [13], Edge-

boxes [14] and CM [15]) or the sliding-window strategy (such as R-CNN [6] and DenseNet [16]). However, this way suffers from high computational burden. The second way is to obtain features by using the sliding-window strategy on convolutional feature maps (such as OverFeat [17] and faster R-CNN [18]). The second way is usually more efficient to obtain features than the first way, when the sliding-window moves at each position. However, one problem for the second way is that the obtained features are less discriminative than those obtained by the first way. In addition, the computational complexity of both ways is high for object detection. Therefore, how to improve the efficiency of object detection remains a challenging problem due to the real-time requirement of object detection.

Face detection is a special object detection task, and it can be tackled by state-of-the-art object detection methods. For example, R-CNN, faster R-CNN and YOLO have been respectively extended to the task of face detection [19–22]. However, as mentioned above, these object detection methods are usually time-consuming. Moreover, these methods mainly rely on the input generic object proposals obtained by an object proposal generation method. Since the obtained object proposals may not cover small-sized faces, these detection methods are usually not effective for detecting small-size faces.

In this paper, we propose a fast face detection method based on Discriminative Complete Features (DCFs) extracted by an elaborately designed deep CNN. The extracted DCFs are sparse and insensitive to scale variations for face detection. As a result, direct classification on DCFs significantly improves the efficiency of

* Corresponding author.

E-mail addresses: gjguo@stu.xmu.edu.cn (G. Guo), hanzi.wang@xmu.edu.cn (H. Wang), yanyan@xmu.edu.cn (Y. Yan), jinzheng@buaa.edu.cn (J. Zheng), boli@buaa.edu.cn (B. Li).

face detection compared with several state-of-the-art face detection methods using CNN. The three key components of the proposed fast face detection method are summarized as follows: First, we generate a sparse feature space, which is implemented by using the rectified linear unit (ReLU) [24]. Second, the multi-scale features are resized to the desirable size (i.e., the same as the size of the input to a fully-connected layer) by using the nearest neighbor interpolation method [25], which can keep the sparseness and the topological structure of features. Third, for each window in an original input image, we propose a fast method to obtain the features based on the complete feature maps [26]. Based on the above three components, the desired features (i.e., DCFs) obtained from the feature maps before the fully-connected layer are extracted. As shown in Fig. 1, DCFs are linearly separable and sparse. Besides, by considering the fully-connected layer in CNN as a linear classifier, the generalization error bound of the linear classifier can be further calculated for face detection.

As the major contribution of this paper, a fast face detection method based on CNN is proposed. The sliding-window strategy can be used on the proposed discriminative complete features, which reduces the computational complexity of the proposed face detection method. Furthermore, the convolution operation in CNN is speeded up by using the Sparse Fast Fourier Transform algorithm. Experiments show the effectiveness and efficiency of the proposed face detection method.

The remainder of this paper is organized as follows. In Section 2, we review related work on object detection and face detection based on CNN. In Section 3, we propose a fast DCFs-based method by using a deep CNN for face detection. In Section 4, we conduct experiments to show the effectiveness of the proposed method and its performance on face detection. The main conclusions are summarized in Section 5.

2. Related work

From low-level vision representations to high-level features, the importance of feature invariance under different conditions (such as changes in illumination and scale) is obvious. DiCarlo et al. [27] review the evidences of neurons and computational models, and discover the computation process which can mimic powerful neuronal representations in the inferior temporal cortex. The authors point out that two neural pathways via a series of visual cortical areas try to untangle the neuronal representations for object recognition and object detection. These two neural pathways are widely known as the ventral pathway and the dorsal pathway, respectively. In some literatures [28,29], CNN plays a role similar to the ventral pathway, where the hierarchical features of CNN correspond to the neuronal representations from the visual area 1 (V1) to the visual area 4 (V4). CNN works well for many different vision tasks by deeply learning the features of data. For instance, Farabet et al. [30] show the impressive performance of CNN on scene labeling. Krizhevsky et al. [31] and He et al. [32] use CNN to classify the ImageNet dataset [8] and achieve the accuracy of 83% and 94.3%, respectively. R-CNN [6] improves mean average precision (mAP) by more than 30% relative to the previous best result on the VOC 2012 dataset [33], and it achieves a mAP of 53.3%.

Although CNN has shown good performance for object detection, its computational efficiency is not high. OverFeat [17] obtains the multi-scale dense features by using the sliding window strategy, and it achieves mAP of 24.3% in the ILSVRC2013 competition [8]. However, the detection process of OverFeat is very time-consuming, because it requires running the classifier and regressor networks across all the possible locations and scales. To improve the efficiency, R-CNN first generates many class-independent proposal windows, and then extracts features on the proposal windows with the trained CNN on a multi-scale images pyra-

mid. Afterwards, a score is assigned to each proposal window by applying a linear SVM to the features. R-CNN takes about 15 seconds to run a detector on a 500×375 image. Another recently proposed method, called SPP-net [7], removes the constraint of the fixed-sized input in the conventional CNN by replacing the pooling layer with a spatial pyramid pooling layer. Instead of extracting features from image regions, SPP-net directly obtains features from the feature maps before the fully-connected layer. However, SPP-net can not deal with flexible CNN architectures for different vision tasks. To reduce the computational complexity of the CNN-based object detection methods and design flexible CNN architectures, Giusti et al. [26] fragment the extended maps generated from each max-pooling layer. When a sliding-window moves in a test image, the corresponding features are obtained by looking up all the fragments in each layer for classification. This object detection method is theoretically faster than the other object detection methods in which features are obtained by using a CNN model on each patch in an image. However, considering the image size and the number of layers, this method is still not computationally efficient. Based on the above reviews, we can see that most current object detection methods based on deep learning still consume much time during the object detection process.

Several object detection methods based on CNN are applied to the task of face detection. For example, Lin et al. [34] propose the MLetNet method, where the number of the output nodes of the LeNet CNN structure [12] is modified for face detection, to detect the masked faces of possible terrorists. FaceCraft [35] jointly trains a region proposal network (RPN) [18] and a fast R-CNN model [36] to improve the detection rate of the fast R-CNN method for face detection. However, as mentioned above, fast R-CNN is time-consuming since it extracts features on an image pyramid. MTCNN [37] utilizes a cascaded multi-task framework, which exploits the inherent correlation between the tasks of detection and alignment, to boost up the performance of both face detection and alignment. Since MTCNN uses multiple CNN networks to detect faces, it is still not computationally efficient. Several other face detection methods also apply fast R-CNN or faster R-CNN to the task of face detection, such as [19,20]. However, these two methods are not effective at detecting small-sized faces since the object proposal generation methods are designed for generic objects instead of faces. Recently, SAFD [38] uses a scale proposal network to predict the possible scales of the faces in a test image, which can effectively boost the performance of CNN based face detectors.

Next, we propose a fast face detection method by directly using the sliding window strategy on discriminative complete feature maps.

3. A fast face detection method based on discriminative complete features

In this section, we describe the details of the proposed fast face detection method via CNN. In Section 3.1, we give the overview of the proposed method. In Section 3.2, we show the details of the proposed DCFs. In Section 3.3, a fast face detection method is given by performing direct classification on DCFs. In Section 3.4, an ensemble model is used to further reduce the classification error on DCFs.

3.1. Overview of the proposed method

In this paper, we propose a fast method for face detection by directly using the sliding window strategy on the feature maps before the fully-connected layer, where a specific CNN is elaborately designed, as shown in Fig. 2. We expect the trained CNN model can map the whole input image to a feature space, in which the new features in each sliding window are linearly separable. Then, the

computational complexity of the proposed face detection method is significantly reduced by performing classification on the feature space. Since the classification layer in a CNN model is usually a linear classifier, we choose the output of the layer before the classification layer as the new features. All the layers before the classification layer act as a nonlinear mapping function (see Fig. 2).

The proposed method includes two key elements, which are the nonlinear mapping function and the sparse discriminative features for face detection. The nonlinear mapping function projects raw data onto the Euclidean space, and the sparseness of features is beneficial to linear separability. As shown in Fig. 2, the convolution layers and the max-pooling layers act as the nonlinear mapping function, where ReLU constrains the output features to be sparse. At the same time, the local contrast normalization (LCN) layers generate compact and competitive features. The LCN layers are inspired by computational neuroscience models [39], where local competition between adjacent features in a feature map is enforced. It has been empirically proven that the LCN layers can reduce the error rates and make supervised learning considerably faster [40].

3.2. Discriminative complete features (DCFs)

To improve the efficiency of the proposed face detection method, we directly perform classification on the feature maps. However, in the process of forward propagation of CNN on a test image, one potential problem is that some discriminative information for face detection is lost when the operation of max-pooling is applied to the whole test image. As illustrated in Fig. 3, the max-pooling layer partitions the input feature map (E) into a set of 2×2 non-overlapping kernels. For each kernel, it outputs the maximum pixel value. There are four ways to tile max-pooling kernels over the feature map (E) corresponding to the four different offsets in the max-pooling kernel. For convenience, the four different offsets are denoted as A, B, C, D, respectively. If the 2×2 non-overlapping max-pooling kernels are tiled over the input feature map (E) beginning with the offset A, the output feature map is shown in Fig. 3(A). For convenience, Fig. 3(A) is called as fragment (A). Similarly, fragments (B), (C) and (D) are obtained when the 2×2 non-overlapping max-pooling kernels are tiled over the input feature map (E) beginning with the offset B, C and D, respectively. The general operation of the max-pooling step on the test image only obtains a fragment (A) after tiling non-overlapping max-pooling kernels over the feature map (E) starting with the offset A, and the other fragments with the other offsets are not obtained. If a sliding window moves at the other offsets (i.e., B, C, D), the obtained features may not be discriminative. Therefore, we adopt the strategy of [26] to keep all the possible fragments. The corresponding fragment is created when we tile non-overlapping max-pooling kernels over the feature map (E) in Fig. 3 starting with different offsets (i.e., A, B, C, D), and each fragment is independent to the other fragments at the same layer. Based on the fragments in the extended feature maps, we can extract DCFs and directly perform classification on DCFs. The definition of each layer, which is used to extract DCFs, is given next.

Given an input image x , the trained filter bank k and bias b , the output of the convolution layer can be written as:

$$x_{oj}^l = f \left(\sum_{i \in M_j} x_{oi}^{l-1} * k_{ij}^l + b_j^l \right), \quad (1)$$

where o denotes the index of fragments; l denotes the current layer; i and j denote the indexes of the input and output feature maps, respectively; M_j represents a selection of the input feature maps, which are the output features from the previous layer; $*$ denotes the operator of convolution; f denotes the activation func-

tion, where $f(x) = \max(0, x)$ is the ReLU [24] activation function. The output of ReLU is 0 if the input is less than 0, and raw output otherwise. Thus, ReLU constrains the output features to be sparse.

For the max-pooling layer, we have

$$x_{oj}^l(m, n) = \max(x_{oj}^{l-1}(p : (p + s), q : (q + s))), \quad (2)$$

where m and n respectively denote the row and column index of a feature map in the current layer, respectively; s denotes the kernel size of the max-pooling layer; $p = s \times (m - 1) + \kappa + 1$ and $q = s \times (n - 1) + \kappa + 1$, where κ ($0 < \kappa < s$) denotes the offset, p and q denote the starting positions of the row and column of each tiled max-pooling kernel, respectively; the colon is used to pick out the selected rows or columns of the feature maps, and here $o \leq s^2$. The max operation has the property of the local invariance [41].

Inspired by the computational neuroscience, the local contrast normalization (LCN) layer [39] mimics the cells in V1 to enforce local competition in the feature maps. In order to extract features from an input image with an arbitrary size, we set the local contrast normalization layer as:

$$x_{oj}^l(m, n) = \frac{x_{oj}^l(m, n)}{(\kappa + \alpha \sum_{i=\max(0, j-\frac{r}{2})}^{\min(N-1, j+\frac{r}{2})} x_{oi}^l(m, n)^2)^\beta}, \quad (3)$$

where r denotes the number of the adjacent feature maps (usually set to be a positive integer in the training process); N denotes the total number of the feature maps at the current layer. κ , α , β are the hyperparameters which are set to be appropriate float point values in the training process. Eqs. (1)–(3) can be used to extract the DCFs, where all the discriminative information for face detection is kept. Thus, we call the features obtained from the feature maps before the fully-connected layer as the discriminative complete features (DCFs). As a result, direct classification on the feature maps before the fully connected layer can improve the efficiency of the detection procedure.

3.3. Speedup for face detection based on DCFs

Based on the above-mentioned three layers (i.e., the convolution layer, the max-pooling layer, the local contrast layer), the input images are mapped into a feature space, where the features for face detection are linear separable. Therefore, we can directly perform classification on DCFs extracted from each patch corresponding to a sliding window by using the weight vector of the fully-connected layer. Instead of using the sliding window technique on the input images, direct classification on DCFs can significantly improve the efficiency for face detection. As a matter of fact, the weight vector of the fully-connected layer can be reshaped into a kernel matrix so that the dot product between DCFs and the weight vector is transformed into the convolution between DCFs and the kernel matrix [42]. In addition, the convolution can be implemented by using the Sparse Fast Fourier Transform algorithm [43], which also improves the speed for face detection since the DCFs of an image are sparse.

To show the efficiency of the proposed method for face detection, the theoretical analysis on the computational complexity of the floating point operations (FLOPS) required by the proposed method and a couple of other methods (including the patch-based method [6] and the image-based method [26]) is presented next. The patch-based method applies a trained CNN model to each overlapping patch in the test image. The image-based method performs the convolution only once for the test image. The main difference between the image-based method and the proposed DCFs-based method is that the image-based method applies the sliding window technique on the original test image, while the proposed DCFs-based method directly uses the sliding window technique on the feature maps. Since the size of the feature maps

is far smaller than that of the original image, the proposed method is much more efficient. In addition, the convolution operation in CNN is speeded up by using the Sparse Fast Fourier Transform algorithm in the proposed face detection method. The FLOPS required by the patch-based method $FLOPS_I^{patch}$, the image-based method $FLOPS_I^{Image}$ and the proposed DCFs-based method $FLOPS_I^{DCF}$ are given in detail next. Here, we only analyze the computational complexity of the FLOPS required by the three methods on the convolutional layer and the fully-connected layer, because these two layers are the most time-consuming.

For the convolutional layer, the FLOPS required by the patch-based method is written as:

$$FLOPS_I^{patch} = 2A_l |P_{l-1}| |P_l| w_l^2 s_l^2, \quad (4)$$

where A_l , w_l , s_l denote the number of the pixels of an input image, the number of the pixels of the feature map in the l th layer, and the kernel size at the l th layer, respectively; $|P_{l-1}|$ and $|P_l|$ denote the numbers of the feature maps in the previous layer and in the current layer, respectively. The patch-based method is quite time-consuming since the number of convolution operations on each feature map is equal to the number of pixels of an input image, which causes a large number of redundant computations.

In contrast, the FLOPS required by the image-based method is written as:

$$FLOPS_I^{Image} = 2A_l |P_{l-1}| |P_l| F_l s_l^2, \quad (5)$$

where A_l denotes the number of the pixels in the feature map; F_l denotes the number of the fragments in the current layer. As the size of feature maps decreases, the FLOPS required by the image-based method on the convolutional layer reduces.

The DCFs-based method utilizes the Sparse Fast Fourier Transform algorithm to improve the efficiency of convolution. Thus, the FLOPS required by the DCFs-based method can be written as:

$$FLOPS_I^{DCF} = 2A_l^* \ln(A_l) |P_{l-1}| |P_l| F_l, \quad (6)$$

where A_l^* denotes the number of the non-zeros pixels in the feature map.

For the fully-connected layer, the FLOPS required by the patch-based method, the image-based method and the DCFs-based method are respectively written as Eqs. (7)–(9).

$$FLOPS_F^{patch} = 2A |P_{l-1}| |P_l| s_l. \quad (7)$$

$$FLOPS_F^{Image} = 2A F_l s_l. \quad (8)$$

$$FLOPS_F^{DCF} = 2A^* \ln(A_l) 2F_l. \quad (9)$$

Since the proposed method uses the sliding window technique on DCFs instead of the original test image, the number of candidate windows used for classification in the proposed method is significantly reduced compared with that used in the other two methods (i.e., the patch-based method and the image-based method).

As the sliding window moves at each position of the feature map in a fragment of DCFs, a resulting response map is obtained by using the classifier layer on the feature vector obtained from each sliding window. Multi-scale DCFs, which are obtained by resizing DCFs with the nearest neighbor interpolation method, are also used for classification. The response maps of the multi-scale DCFs are shown in Fig. 4. The estimated scale of a face candidate is obtained by using the non-maximum suppression technique, and the response maps of all fragments are merged into one response map using non-maximum suppression. Finally, by backprojecting the estimated scale and the position of the centroids of the connected regions in the response map onto the original image, the locations of the faces can be detected with the estimated scale in the original image.

3.4. An ensemble model for reducing generalization error

In this section, an ensemble model [44] is used to improve the performance of the proposed method. Based on the PAC theory [45], which shows the theoretical analysis of a linear classifier, we can analyze the generalization error of the linear classifier at the fully-connected layer. In the framework of PAC, the PAC-bayesian margin bound is the upper bound of generalization error $R[w]$ (see Theorem 3 in [46]). Then, the number of model instances can be computed. Specifically, given the upper bound of generalization error $R[w]$ and a generalization error threshold of the linear classifier $H[w]$, the number of model instances v can be calculated as:

$$v = \left\lceil \frac{\ln(H[w])}{\ln(R[w])} \right\rceil, \quad (10)$$

where $\lceil \cdot \rceil$ denotes the ceiling function.

3.5. Face bounding box regression

Inspired by R-CNN [47] and MDNet [48], we also train a regressor, which is used to regress the face bounding boxes. The regressor consists of two linear layers, where the first layer has one hundred output nodes and the second layer has four output nodes. The learning objective of the regressor is formulated as the mean squared error between the four parameter values of a candidate face window (i.e., the horizontal and vertical coordinate values of the top-left corner, the width and height of the face bounding box) and those of the ground truth face bounding box. For each candidate face window, its DCFs are used to train the regressor. We expect that the regressor can refine the face bounding box. Usually, the candidate face window is obtained by using the sliding window strategy on DCFs with several fixed scales. However, the resulting face bounding boxes may not be accurate enough, since the width and height of a face bounding box are continuous values instead of the values in several fixed scales. The regressor can overcome the above drawback by further refining the bounding box.

4. Experiments

In this section, the detailed evaluation of the performance of the proposed method on face detection is shown. Section 4.1 gives the parameter settings used in training the CNN model. Section 4.2 visualizes the trained CNN network. Section 4.3 presents the comparison of efficiency and performance between the proposed method and several state-of-the-art face detection methods on two public face datasets.

4.1. Parameter settings

We use the GPU implementation of the *cuda-convnet* code [31] to train the CNN model on various datasets. In order to adapt the parameters of the trained model to test images of arbitrary size, the batch normalization step in *cuda-convnet* is not performed during the training process.

Data pre-processing: We collect 13,466 face images from the dataset provided by Sun et al. [49], 29,821 face images from the CAS-PEAL-R1 dataset [50] and 47,220 face images from the web. The collected face dataset totally includes 90,507 face images with different lighting conditions, poses and ages. Moreover, 100,000 non-face images from the VOC2012 dataset [33] are collected for training. All the face and non-face images are converted to the gray images and resized to the size of 32×32 . All images are split into 19 batches and 10,000 images per batch are used for training.

Hyperparameters for training: To balance the efficiency and effectiveness of the proposed face detection method, we design a lightweight CNN structure, which is shown in Table 1. As shown

Table 1
The CNN structure for face detection.

Layer	Type	Input	#Channels	#Filters	Filter Size	Activation
Layer1	conv	Data	3	16	5×5	ReLU
Layer2	pool	Layer1	16	—	—	—
Layer3	cmrnorm	Layer2	16	—	—	—
Layer4	conv	Layer3	16	16	5×5	ReLU
Layer5	cmrnorm	Layer4	16	—	—	—
Layer6	pool	Layer5	16	—	—	—
Layer7	FC	Layer6	—	—	—	—
Layer8	softmax	Layer7	—	—	—	—

Table 2
The hyperparameters used for training the CNN model of the proposed method.

LayerType	epsW	epsB	momW	momB	wc
Layer1	0.0030	0.0040	0.9000	0.9000	0.0000
Layer4	0.0001	0.0002	0.9000	0.9000	0.0000
Layer7	0.0002	0.0003	0.9000	0.9000	0.0100
scale	—	—	—	—	—
Layer3	0.0010	0.7500	—	—	—
Layer5	0.0010	0.7500	—	—	—

in Table 1, Layer1 and Layer4 are the convolutional layers whose filter size is 5×5 and the number of output feature maps is 16; Layer2 and Layer6 are the max-pooling layers whose filter size is 2×2 ; Layer3 and Layer5 are the local contrast normalization layers defined in Eq. (3); Layer7 is a fully connected layer. The hyperparameters used for training are given in Table 2, where *epsW* and *epsB* respectively denote the learning rate for the weight vector and the bias in the CNN model; *momW* and *momB* respectively denote the momentum of the weight vector and that of the bias in the CNN model; *wc* denotes the weight decay parameter. In the training process, we use 16 batches of images for training and 3 batches of images for validation. The whole training process takes around one hour on a computer equipped with a GTX-780Ti GPU. The trained CNN achieves 99.74% of test precision on the two validation batches. Given $H[w] = 0.0008$, two CNN models (calculated by Eq. (10)) are trained, and the filters at the convolution layers in one model are shown in Fig. 5.

Test environment: For face detection, our method is programmed in Matconvnet [51] and runs on a computer equipped with a i7-4.0 GHZ CPU, a GTX-780Ti GPU and 32 GB RAM. We evaluate the proposed face detection method on the Fddb [52] and the AFW [53] datasets. The Fddb dataset contains 5171 faces in 2845 images and the AFW dataset contains 468 faces in 205 images.

4.2. Visualization of the trained CNN

After training the parameters of the proposed method using the *cuda-convnet* framework, the trained model is applied to detect the faces in the test images. One fragment of the DCFs is shown in Fig. 6. We can see that the intensities of most pixels in some feature maps are zero due to the usage of ReLU, and the sparse ratio is approximately equal to 0.5.

In order to analyze the scale invariance of DCFs, multi-scale DCFs are obtained by resizing the DCFs from the original images using the nearest neighbor interpolation method. The multi-scale DCFs corresponding to the faces (represented by the bounding boxes) in the test images are resized to the size of the fully-connected layer. For fair comparison, the general features before the fully-connected layer in the traditional CNN [12] and the dense features obtained by the proposed method with the sigmoid activation function are also evaluated. All the CNN models use the same structure as shown in Table 1. Using the t-SNE [23] visualization algorithm, 5000 randomly selected resized DCFs and

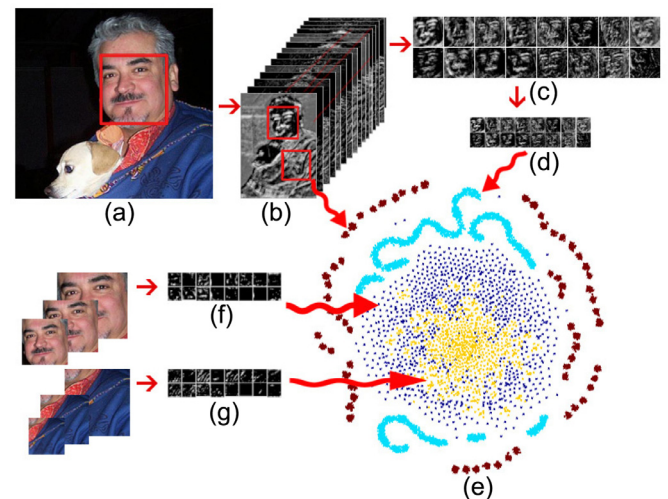


Fig. 1. The visualization of DCFs and general CNN features by using the t-SNE [23] visualization algorithm. (a) A test image. (b) The corresponding DCFs of the test image. (c) The features obtained from a set of the sliding windows. (d) Resizing the obtained features into the size of the fully-connected layer of CNN. (e) The visualization of the distribution of DCFs (shown in the brown and cyan dots) and general features of CNN (shown in the blue and yellow dots). (f) and (g) The general features extracted based on an image pyramid via CNN.

general features are shown in Fig. 7. From the figure, we can see that the resized DCFs are still robust for linear classification since the DCFs are sparse. However, the dense features extracted by the traditional CNN are not robust for linear classification.

4.3. The performance comparison for face detection

In this section, we compare the computational complexity between the proposed method and several other methods (including the image-based and patch-based methods) for face detection. We assume that five scales are used for detection and we use a test image with the size of 800×600 for evaluation. The FLOPS required the three methods at different layers (i.e., 1, 4, 7) in CNN is shown in Table 3. Moreover, Table 4 gives the CPU and GPU time used by the four CNN-based face detection methods on a 800×600 test image. In Table 4, for the patch-based method, we use the representative R-CNN, while we adopt the method of Giusti et al. [26] for the image-based method. In addition, SPP-net, Fast R-CNN, Faster R-CNN, DeepIR and YOLO are compared since these face detection methods are closely related to the proposed face detection method. SPP-net is similar to the proposed DCFs-based face detection method, because both methods directly perform detection on feature maps. However, compared with the proposed method, SPP-net uses a more complicated CNN structure, which has more computational burden for face detection. From Tables 3 and 4, we can see that the proposed method is much more efficient than the patch-based method, the image-based method and the SPP-net method considering both theoretically

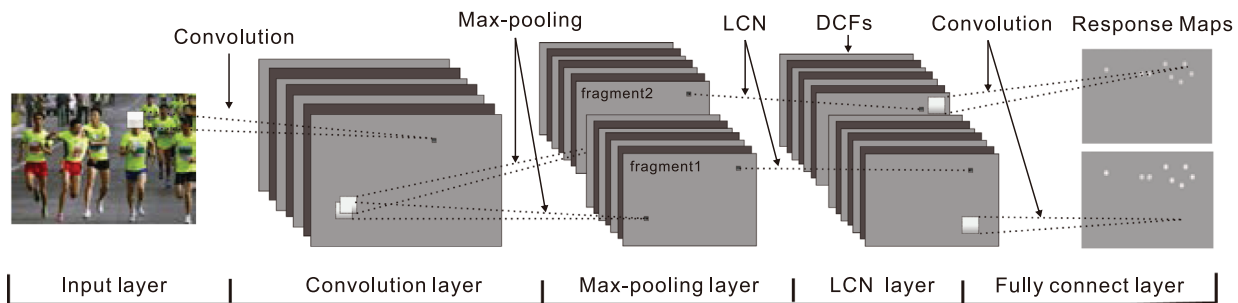


Fig. 2. The proposed DCFs-based method for object detection. Different fragments are created corresponding to different offsets in a max-pooling layer, which preserves all the possible discriminative features. LCN is the abbreviation of Local Contrast Normalization.

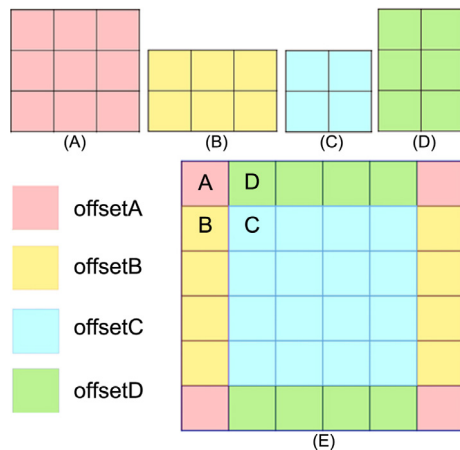


Fig. 3. Illustration of the fragments of the extended feature maps based on different offsets of the max-pooling kernel. (A), (B), (C), (D), respectively denote different fragments of the extended maps after max-pooling corresponding to different offsets in (E). (E) denotes the input feature map for a max-pooling layer. For the kernel size of 2×2 , A, B, C, D in (E) denote different offsets respectively.

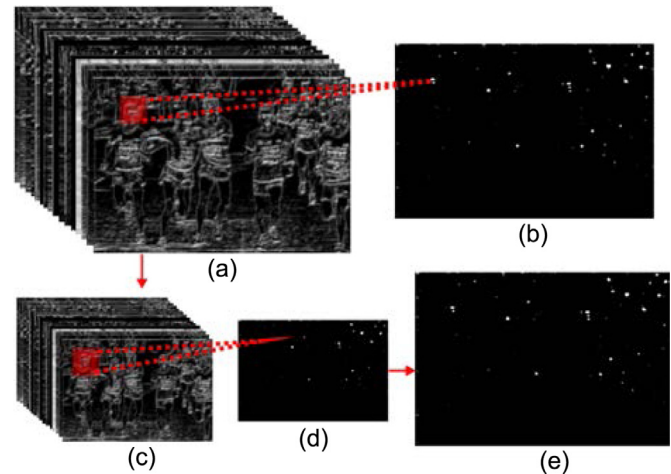


Fig. 4. Face detection on the DCFs feature maps. (a) The DCFs obtained by the proposed method. (b) The response map corresponding to the DCFs at the last layer. (c) The resized feature maps obtained by the nearest neighbor interpolation method. (d) The output response map of the resized DCFs. (e) The resized response map (the same size as (b)) for further processing by using non-maximum suppression.

computational complexity (i.e., the required FLOPS in Table 3) and running time (in Table 4). The proposed face detection method also runs faster than the fast R-CNN, faster R-CNN and DeepIR methods. The reason is that the CNN structures used in the three competing methods are more complex than that of the proposed method, which results in more computation. Moreover, the object proposal generation methods employed in the three competing methods consume more CPU time (about 2 s) than the proposed method (about 0.1 s). In terms of GPU time, the running time of the proposed method is 0.09, which is faster than that of the other six competing methods (about 3.3 to 278.9 times faster) except



Fig. 5. The trained filters of the first convolutional layer in the CNN model used by the proposed method for face detection.

for YOLO. The YOLO method, which is implemented in C/C++, runs faster than the proposed face detection method, which is implemented in Matlab and C/C++. However, the proposed face detection method obtains higher true positive rate than YOLO (see Fig. 9).

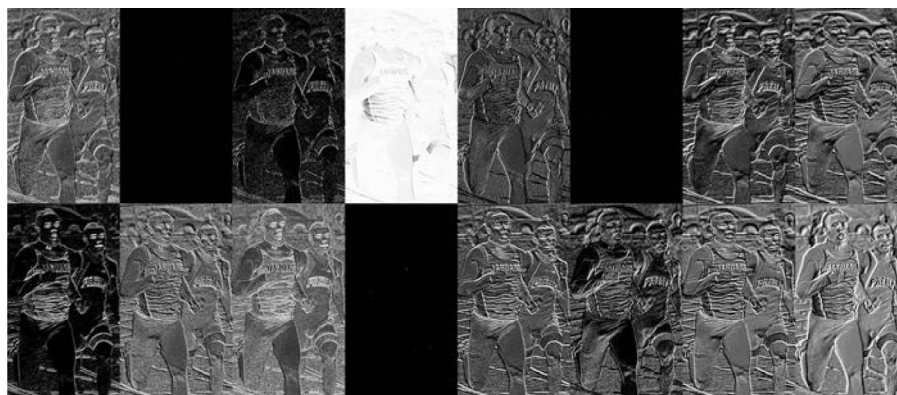


Fig. 6. One fragment of DCFs at the sixth layer. The number of the feature maps in one fragment is equal to the number of channels at the sixth layer in Table 1.

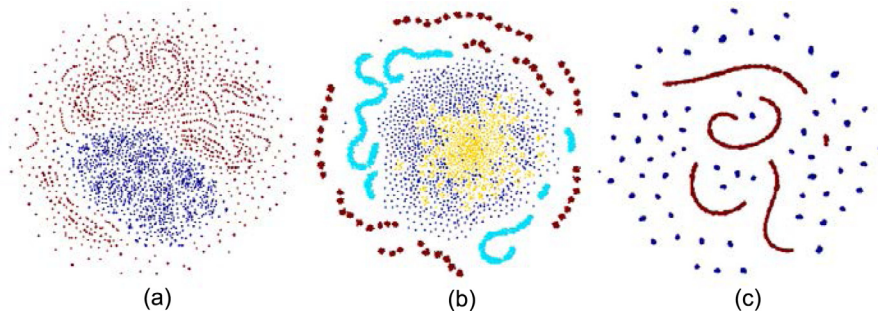


Fig. 7. The visualization of the separability of the features extracted by CNN with different architectures. (a) General features of 32×32 images extracted by CNN (the blue and red points denote the face features and the nonface features, respectively). (b) The resized DCFs in accompany with the general features (also see Fig. 1 for details). (c) The features of 32×32 images extracted by CNN with the activation function of sigmoid (the red and blue points denote the face features and the nonface features, respectively).

Table 3

The FLOPS [$\cdot 10^{10}$] required by the three methods at three layers with a given parameter setting.

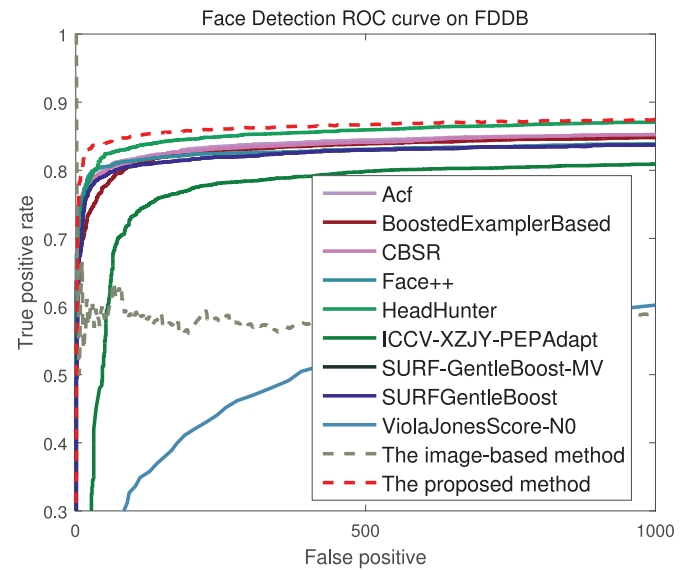
	Layer			Total
	1	4	7	
A	800×600	800×600	800×600	—
$ P_{l-1} $	1	16	16	—
$ P_l $	16	16	16	—
A_l	800×600	400×300	200×150	—
S_l	25	25	1024	—
F_l	1	4	16	—
$FLOPS^{Patch}$	34.56×5	157.2864×5	0.98304×5	964.1472
$FLOPS^{Image}$	0.002×5	0.6144×5	0.98304×5	7.9972
$FLOPS^{DCF}$	0.002	0.135168	0.00096×5	0.141968

Table 4

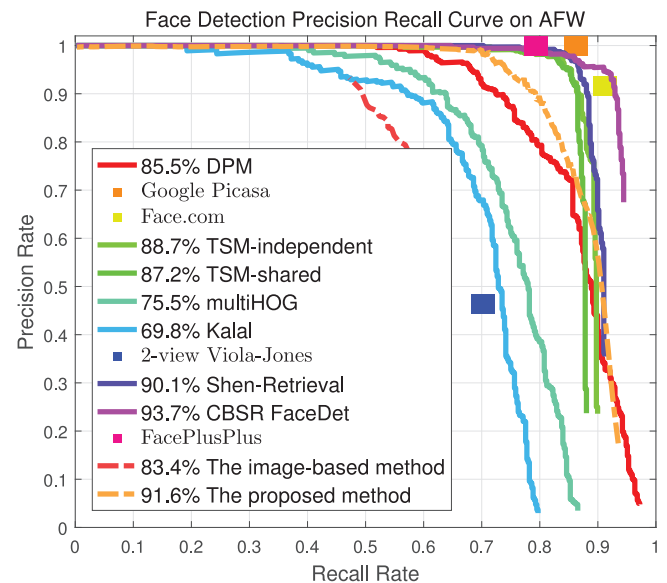
The CPU and GPU time (in seconds) used by the seven CNN based face detection methods and the proposed face detection method on a 800×600 test image. Note that YOLO is implemented in C/C++, and the other competing methods are implemented in Matlab and C/C++.

Method	CPU time [s]	GPU time [s]	Total time
Patch-based	2.3	25.1	28.1
Image-based	0	2.7	2.7
SPP-net	2.3	0.3	2.6
Fast R-CNN	2.7	10.1	12.8
Faster R-CNN	0.3	0.4	0.7
DeepIR	0.3	0.4	0.7
YOLO	0	0.04	0.04
DCF-based	0.1	0.09	0.19

Next, we compare the performance obtained by different methods for face detection on the Fddb and AFW datasets. On the Fddb dataset, we use the criterion given in [52,54], which measures the degree of match between a detected and a manually annotated face region. The cutoff ratio of the intersected areas to the joined areas of the two regions is set to 0.5 in our case. We compare the proposed method with the following detection methods: (1) ACF-multiscale [55], (2) HeadHunter [56], (3) CBSR [57], (4) Boosted Exemplar [58], (5) SURF frontal/multiview [59], (6) SURF GentleBoost [59], (7) Face++ [60], (8) PEP-Adapt [61], (9) Viola-Jones [62] and (10) the image-based method [26]. We select these methods since most of these methods are the state-of-the-art face detection methods based on hand-crafted features, and their results have been reported in the Fddb website [52]. In addition, we apply the image-based method to face detection as a comparison. The methods in [38,63], which also use CNN for the task of face detection, are not compared since their source codes are not available in public. The comparison of the discrete ROC curves obtained by the competing methods is shown in Fig. 8(a), from which we can see that the proposed method achieves better



(a)



(b)

Fig. 8. The comparison of the discrete ROC and PR curves for face detection obtained by the competing methods on the Fddb and AFW datasets. (a) The results of ROC curves obtained by the competing methods on the Fddb dataset. (b) The results of PR curves obtained by the competing methods on the AFW dataset.

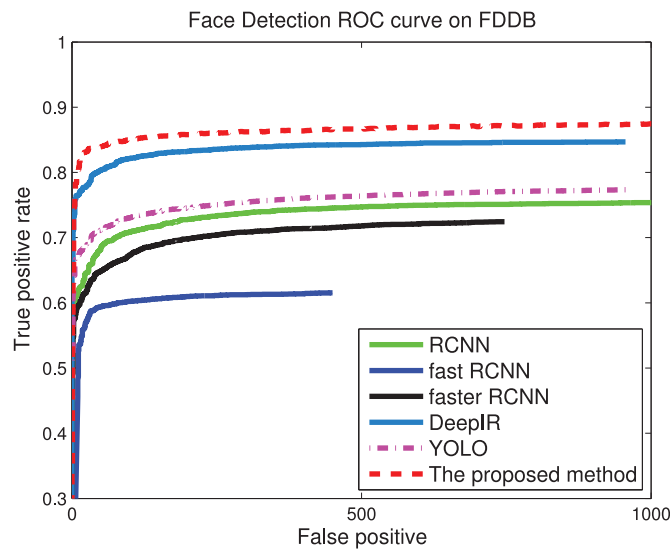


Fig. 9. The comparison of the discrete ROC curves obtained by the six face detection methods based on CNN on the FDDB dataset.

performance than most of the competing methods except for HeadHunter. Due to the fact that the proposed method uses CNN to learn effective features of faces, the proposed method obtains better performance than that obtained by the most other competing methods. Although the image-based method also learns features for faces by using CNN, it obtains worse performance than that obtained by the other competing methods except for Viola-Jones. The main reason is that the features obtained the image-based method are less discriminative and not robust to scale variations. Moreover, the image-based method does not use an extra CNN model to eliminate false detection.

On the AFW dataset, we report the Precision-Recall curves for evaluation. Using the same settings as in [57], we compare the proposed method with the following eleven methods: (1) DPM [53], (2) CBSR [57], (3) TSM [53], (4) Kalal [64], (5) Face++ [60], (6) Google Picasa, (7) Face.com, (8) multiHOG [53], (9) Viola-Jones [62], (10) Shen-Retrieval [65], and (11) Image-based method [26].

We select these methods since they are state-of-the-art face detection methods, which have been tested on the AFW dataset. The comparison of the Precision-Recall curves is shown in Fig. 8(b). The average precision obtained by the proposed method on the AFW dataset is 91.6%, which is worse than CBSR, Google Picasa and Face.com due to the lack of a large number of training faces in the wild. However, the proposed method shows excellent classification performance when the recall ratio is high. The recall ratio obtained by the proposed method is higher than that obtained by DPM, Kalal, multiHOG, Viola-Jones and the image-based method when the precision ratio is greater than 0.7.

Several object detection methods based on CNN (such as R-CNN, faster R-CNN, YOLO) can be used for face detection, and they can achieve comparative performance on the FDDB dataset. To show the effectiveness of the proposed method, we implement and evaluate several state-of-the-art object detection methods on the FDDB dataset. We compare the proposed method against the following five face detection methods based on CNN: (1) R-CNN [19], (2) fast R-CNN [20], (3) faster R-CNN [19], (4) DeepIR [21] and (5) YOLO [22]. The five competing methods perform face detection by using CNN classifiers on the face proposals generated by different object proposal generation methods. The YOLO method applies a single neural network to a whole image. The network divides the image into regions and predicts bounding boxes and probabilities of each region being an object. Then, these bounding boxes are weighed by the predicted probabilities. Fig. 9 shows the ROC curves obtained by the six competing methods on the FDDB dataset. As can be seen, the proposed method obtains the best performance. The reason is that the proposed method can obtain effective features when a sliding window strategy is applied on DCFs. The R-CNN, fast R-CNN and faster R-CNN methods obtain worse performance than the other competing methods since the performance of the three methods for face detection is greatly affected by the quality of object proposals generated by the object proposal generation methods. More specifically, both R-CNN and fast R-CNN use the Selective Search object proposal generation method (SS) to generate object proposals, and faster R-CNN uses the Region Proposal Network (RPN) to generate object proposals. However, both SS and RPN are not effective at generating object proposals for faces, especially for small-sized faces. DeepIR improves RPN by re-

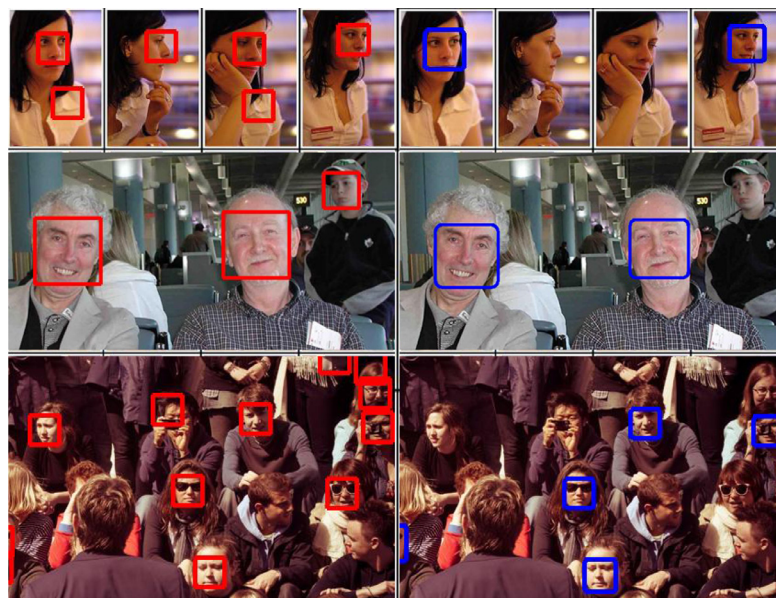


Fig. 10. The comparison results of face detection. Red bounding boxes (left) denote the result obtained by the proposed method. Blue bounding boxes (right) denote the result obtained by Face++.

training the PRN model on the face dataset, and it increases the discriminative ability of the classifier by using the hard negative mining strategy. Thus, DeepIR obtains better performance than R-CNN, fast R-CNN and faster R-CNN. However, DeepIR obtains worse performance than the proposed method since the latter directly uses a sliding window strategy on the DCFs, which can achieve high face recall than the object proposal generation method used in DeepIR.

Some detection results are given in Fig. 10. We can see that the proposed method can find more correct faces with occlusions and different poses, where the false detection can be eliminated by using an extra CNN model during the classification process.

5. Conclusion

In this paper, we propose a fast face detection method based on DCFs extracted by an elaborately designed CNN. Compared with the state-of-the-art face detection methods using CNN, the proposed method performs direct classification on DCFs, which can significantly improve the efficiency during the face detection process. Moreover, the proposed method can effectively detect small-size faces by using the sliding window strategy on DCFs. In contrast, current state-of-the-art face detection methods based on CNN are hard to detect small-sized faces since the performance of these methods is greatly affected by the quality of object proposals generated by the object proposal generation methods. Experimental results have shown that the proposed DCFs-based face detection method can achieve promising performance on several popular face detection datasets.

Declaration of competing interest

None

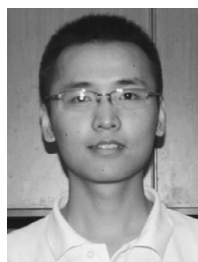
Acknowledgments

This work was supported by the National Natural Science Foundation of China under Grants U1605252, 61872307, 61571379 and 61876014, and by the Natural Science Foundation of Fujian Province of China under Grant 2017J01127.

References

- [1] D. Marr, W.H. Freeman, Company, Vision, MIT Press, 1982.
- [2] D.G. Lowe, Distinctive image features from scale-invariant keypoints, *Int. J. Comput. Vis.* 60 (2) (2004) 91–110.
- [3] N. Dalal, B. Triggs, Histograms of oriented gradients for human detection, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2005, pp. 886–893.
- [4] J. Canny, A computational approach to edge detection, *IEEE Trans. Pattern Anal. Mach. Intell.* 8 (6) (1986) 679–698.
- [5] M. Zeiler, G. Taylor, R. Fergus, Adaptive deconvolutional networks for mid and high level feature learning, in: *Proceedings of the IEEE International Conference Computer Vision*, 2011, pp. 2018–2025.
- [6] R. Girshick, J. Donahue, T. Darrell, J. Malik, Rich feature hierarchies for accurate object detection and semantic segmentation, in: *Proceedings of the IEEE Conference Computer Vision Pattern Recognition*, 2014, pp. 580–587.
- [7] K. He, X. Zhang, S. Ren, J. Sun, Spatial pyramid pooling in deep convolutional networks for visual recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* 37 (9) (2015) 1904–1916.
- [8] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M.S. Bernstein, A.C. Berg, L. Fei-Fei, Imagenet large scale visual recognition challenge, *Int. J. Comput. Vis.* 115 (3) (2015) 211–252.
- [9] Q.-Q. Tao, S. Zhan, X.-H. Li, T. Kurihara, Robust face detection using local CNN and SVM based on kernel combination, *Neurocomputing* 211 (2016) 98–105.
- [10] F.H.C. Tivive, A. Bouzerdoum, A hierarchical learning network for face detection with in-plane rotation, *Neurocomputing* 71 (16) (2008) 3253–3263.
- [11] K. Lu, X. An, J. Li, H. He, Efficient deep network for vision-based object detection in robotic applications, *Neurocomputing* 245 (2017) 31–45.
- [12] Y. LeCun, L. Bottou, G. Orr, K.-R. Müller, Efficient backprop, in: *Proceedings of the Advances in Neural Information Processing Systems*, 1998, pp. 9–50.
- [13] J. Uijlings, K. van de Sande, T. Gevers, A. Smeulders, Selective search for object recognition, *Int. J. Comput. Vis.* 104 (2) (2013) 154–171.
- [14] C.L. Zitnick, P. Dollár, Edge boxes: locating object proposals from edges, in: *Proceedings of the European Conference on Computer Vision Conference*, 2014, pp. 391–405.
- [15] G. Guo, H. Wang, W.L. Zhao, Y. Yan, X. Li, Object discovery via cohesion measurement, *IEEE Trans. Cybern.* 40 (3) (2018) 862–875.
- [16] F.N. Iandola, M.W. Moskewicz, S. Karayev, R.B. Girshick, T. Darrell, K. Keutzer, Densenet: implementing efficient convnet descriptor pyramids, *CoRR* (2014).
- [17] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, Y. LeCun, Overfeat: integrated recognition, localization and detection using convolutional networks, *Proceedings of the International Conference on Learning Representations* (2014).
- [18] S. Ren, K. He, R. Girshick, J. Sun, Faster R-CNN: Towards real-time object detection with region proposal networks, in: *Proceedings of the Advances in Neural Information Processing Systems*, 2015, pp. 91–99.
- [19] H. Jiang, E.G. Learned-Miller, Face detection with the faster R-CNN, *CoRR abs/1606.03473* (2016).
- [20] D. Triantafyllidou, A. Tefas, A fast deep convolutional neural network for face detection in big visual data, in: *Proceedings of the INNS Conference on Big Data Advances in Big Data*, 2017, pp. 61–70.
- [21] X. Sun, P. Wu, S.C.H. Hoi, Face detection using deep learning: an improved faster RCNN approach, *CoRR abs/1701.08289* (2017).
- [22] J. Redmon, S. Divvala, R. Girshick, A. Farhadi, You only look once: Unified, real-time object detection, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 779–788.
- [23] L. van der Maaten, G. Hinton, Visualizing high-dimensional data using T-SNE, *J. Mach. Learn. Res.* 9 (2008) 2579–2605.
- [24] G.E. Hinton, Rectified linear units improve restricted Boltzmann machines, in: *Proceedings of the International Conference on Machine Learning*, 2010, pp. 807–814.
- [25] R.C. Gonzalez, R.E. Woods, *Digital Image Processing*, (3rd ed.), Prentice-Hall, Inc., 2006.
- [26] A. Giusti, D.C. Ciresan, J. Masci, L.M. Gambardella, J. Schmidhuber, Fast image scanning with deep max-pooling convolutional neural networks, in: *Proceedings of the IEEE International Conference on Image Processing*, 2013, pp. 4034–4038.
- [27] J.J. DiCarlo, D. Zoccolan, N.C. Rust, How does the brain solve visual object recognition? *Neuron* 73 (2012) 415–434.
- [28] U. Güçlü, M.A.J. van Gerven, Deep neural networks reveal a gradient in the complexity of neural representations across the ventral stream, *J. Neurosci.* 35 (27) (2015) 10005–10014.
- [29] K. Simonyan, A. Zisserman, Two-stream convolutional networks for action recognition in videos, in: *Proceedings of the Advances in Neural Information Processing Systems*, 2014, pp. 568–576.
- [30] C. Farabet, C. Couprie, L. Najman, Y. LeCun, Learning hierarchical features for scene labeling, *IEEE Trans. Pattern Anal. Mach. Intell.* 35 (8) (2013) 1915–1929.
- [31] A. Krizhevsky, I. Sutskever, G.E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Proceedings of the Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.
- [32] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [33] M. Everingham, L. Van Gool, C.K.I. Williams, J. Winn, A. Zisserman, The pascal visual object classes (VOC) challenge, *Int. J. Comput. Vis.* 88 (2) (2010) 303–338.
- [34] S. Lin, L. Cai, X. Lin, R. Ji, Masked face detection via a modified Lenet, *Neurocomputing* 218 (2016) 197–202.
- [35] H. Qin, J. X. Li, X. Hu, Joint training of cascaded CNN for face detection, in: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, 2016.
- [36] R. Girshick, Fast R-CNN, in: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, 2015, pp. 1440–1448.
- [37] K. Zhang, Z. Zhang, Z. Li, Y. Qiao, Joint face detection and alignment using multitask cascaded convolutional networks, *IEEE Signal Process. Lett.* 23 (10) (2016) 1499–1503.
- [38] Z. Hao, Y. Liu, H. Qin, J. Yan, X. Li, X. Hu, Scale-aware face detection, in: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, 2017.
- [39] N. Pinto, D.D. Cox, J.J. Dicarlo, Why is real-world visual object recognition hard, *PLoS Comput. Biol.* 4 (1) (2008) 415–434.
- [40] K. Jarrett, K. Kavukcuoglu, M. Ranzato, Y. LeCun, What is the best multi-stage architecture for object recognition? in: *Proceedings of the IEEE International Conference Computer Vision*, 2009, pp. 2146–2153.
- [41] D. Scherer, A. Müller, S. Behnke, Evaluation of pooling operations in convolutional architectures for object recognition, in: *Proceedings of the International Conference on Artificial Neural Networks: Part III*, 2010, pp. 92–101.
- [42] E. Shelhamer, J. Long, T. Darrell, Fully convolutional networks for semantic segmentation, *IEEE Trans. Pattern Anal. Mach. Intell.* 39 (4) (2017) 640–651.
- [43] B. Ghazi, H. Hassanieh, P. Indyk, D. Katabi, E. Price, L. Shi, Sample-optimal average-case sparse fourier transform in two dimensions, *Proceedings of the Annual Allerton Conference on Communication, Control, and Computing* (2013) 1258–1265.
- [44] S.V. Kozyrev, Classification by ensembles of neural networks, *P-Adic Num. Ultram. Anal. Appl.* 4 (1) (2012) 27–33.
- [45] L.G. Valiant, A theory of the learnable, *Commun. ACM* 27 (11) (1984) 1134–1142.

- [46] R. Herbrich, T. Graepel, A Pac-Bayesian margin bound for linear classifiers: Why SVMs work, in: *Proceedings of the Advances in Neural Information Processing Systems*, 2001, pp. 224–230.
- [47] R.B. Girshick, J. Donahue, T. Darrell, J. Malik, Region-based convolutional networks for accurate object detection and segmentation, *IEEE Trans. Pattern Anal. Mach. Intell.* 38 (1) (2016) 142–158.
- [48] H. Nam, B. Han, Learning multi-domain convolutional neural networks for visual tracking, in: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, 2016, pp. 4293–4302.
- [49] Y. Sun, X. Wang, X. Tang, Deep convolutional network cascade for facial point detection, in: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, 2013, pp. 3476–3483.
- [50] W. Gao, B. Cao, S. Shan, X. Chen, D. Zhou, X. Zhang, D. Zhao, The CAS-peal large-scale chinese face database and baseline evaluations, *IEEE Trans. Sys. Man Cyber. Part A* 38 (1) (2008) 149–161.
- [51] A. Vedaldi, K. Lenc, Matconvnet-convolutional neural networks for matlab, in: *Proceeding of the ACM International Conference on Multimedia*, 2015, pp. 689–692.
- [52] V. Jain, E. Learned-Miller, Fddb: A Benchmark for Face Detection in Unconstrained Settings, Technical Report, University of Massachusetts, Amherst, 2010.
- [53] D. Ramanan, Face detection, pose estimation, and landmark localization in the wild, in: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, 2012, pp. 2879–2886.
- [54] S. Zhan, Q.-Q. Tao, X.-H. Li, Face detection using representation learning, *Neurocomputing* 187 (2016) 19–26.
- [55] B. Yang, J. Yan, Z. Lei, S.Z. Li, Aggregate channel features for multi-view face detection, in: *Proceedings of the International Joint Conference on Biometrics*, 2014, pp. 1–8.
- [56] M. Mathias, M.P. R. Benenson, L.V. Gool, Face detection without bells and whistles, in: *Proceedings of the European Conference on Computer Vision*, 2014, pp. 720–735.
- [57] J. Yan, Z. Lei, L. Wen, S.Z. Li, The fastest deformable part model for object detection, in: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, 2014, pp. 2497–2504.
- [58] H. Li, Z. Lin, J. Brandt, X. Shen, G. Hua, Efficient boosted exemplar-based face detection, in: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, 2014, pp. 1843–1850.
- [59] J. Li, Y. Zhang, Learning surf cascade for fast and accurate object detection, in: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, 2013, pp. 3468–3475.
- [60] M. Inc., Face++ research toolkit, 2013, www.faceplusplus.com.
- [61] H. Li, G. Hua, Z. Lin, J. Brandt, J. Yang, Probabilistic elastic part model for unsupervised face detector adaptation, in: *Proceedings of the IEEE International Conference Computer Vision*, 2013, pp. 793–800.
- [62] P. Viola, M. Jones, Robust real-time object detection, *Int. J. Comput. Vis.* 57 (2) (2001) 137–154.
- [63] H. Li, Z. Lin, X. Shen, J. Brandt, G. Hua, A convolutional neural network cascade for face detection, in: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, 2015, pp. 5325–5334.
- [64] Z. Kalal, K. Mikolajczyk, J. Matas, Face-TLD: Tracking-learning-detection applied to faces, in: *Proceedings of the IEEE International Conference on Image Processing*, 2010, pp. 3789–3792.
- [65] X. Shen, Z. Lin, J. Brandt, Y. Wu, Detecting and aligning faces by image retrieval, in: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, 2013, pp. 3460–3467.



Guanjun Guo is currently pursuing the Ph.D. degree with the School of information science technology, Xiamen University, Xiamen, China. His current research interests include computer vision and machine learning.



Hanzi Wang is currently a Distinguished Professor of “Minjiang Scholars” in Fujian province and a Founding Director of the Center for Pattern Analysis and Machine Intelligence (CPAMI) at XMU. He was an Adjunct Professor (2010–2012) and a Senior Research Fellow (2008–2010) at the University of Adelaide, Australia; an Assistant Research Scientist (2007–2008) and a Postdoctoral Fellow (2006–2007) at the Johns Hopkins University; and a Research Fellow at Monash University, Australia (2004–2006). He received his Ph.D. degree in Computer Vision from Monash University where he was awarded the Douglas Lampard Electrical Engineering Research Prize and Medal for the best Ph.D. thesis in the Department. His research interests are concentrated on computer vision and pattern recognition including visual tracking, robust statistics, object detection, video segmentation, model fitting, optical flow calculation, 3D structure from motion, image segmentation and related fields. He has published more than 100 academic papers. He is a senior member of the IEEE. He was an Associate Editor for IEEE Transactions on Circuits and Systems for Video Technology (from 2010 to 2015). He was the General Chair for IEEE FMT2018 and ACM ICIMCS2014, Program Chair for CVRS2012, and Area Chair for ICPR2018, ACCV2016, DICTA2010. He also serves on the program committee (PC) of ICCV, ECCV, CVPR, ACCV, PAKDD, ICG, ADMA, CISP, etc, and he serves on the reviewer panel for more than 40 journals and conferences.



Yan Yan is currently an associate professor at the School of Information Science and Technology at Xiamen University, China. He received the Ph.D. degree in Information and Communication Engineering from Tsinghua University, China, in 2009. He worked at Nokia Japan R&D center as a research engineer (2009–2010) and Panasonic Singapore Lab as a project leader (2011). He has published more than 30 papers in international journals and conferences including IEEE T-IP, IEEE T-ITS, PR, KBS, ICCV, ECCV, ACM MM, etc. His research interests include computer vision and pattern recognition.



Jin Zheng received the B.S. degree in applied mathematics and informatics from the College of Science in 2001, and the M.S. degree from the School of Computer Science, Liaoning Technical University, Fuxing, China, in 2004, and the Ph.D. degree from the School of Computer Science and Engineering, Beihang University, Beijing, China, in 2009. She is currently a Teacher with Beihang University. Her current research interests include moving object detection and tracking, object recognition, image enhancement, video stabilization, and video mosaic.



Bo Li received the B.S. degree in computer science from Chongqing University in 1986, the M.S. degree in computer science from Xi'an Jiaotong University in 1989, and the Ph.D. degree in computer science from Beihang University in 1993. Now he is a professor of Computer Science and Engineering at Beihang University, the Director of Beijing Key Laboratory of Digital Media, and has published over 100 conference and journal papers in diversified research fields including digital video and image compression, video analysis and understanding, remote sensing image fusion and embedded digital image processor.