

1 International Journal of Pattern Recognition and Artificial Intelligence  
2 © World Scientific Publishing Company

### 3 **Face Occlusion Detection Using Deep Convolutional Neural Networks**

4 Yizhang Xia\*, Bailing Zhang

5 *Department of Computer Science and Software Engineering, Xian Jiaotong-Liverpool*  
6 *University, SIP, Suzhou 215123, China*  
7 *yizhang.xia@xjtlu.edu.cn*

8 Frans Coenen

9 *Department of Computer Science, University of Liverpool, Liverpool L69 3BX, UK*  
10 *Coenen@liverpool.ac.uk*

11 With the rise of crimes associated with Automated Teller Machines (ATMs), security  
12 reinforcement by surveillance techniques has been a hot topic on the security agenda. As  
13 a result, cameras are frequently installed with ATMs, so as to capture the facial images  
14 of users. The main objective is to support follow-up criminal investigations in the event  
15 of an incident. However, in the case of miss-use, the user's face is often occluded. There-  
16 fore, face occlusion detection has become very important to prevent crimes connected  
17 with ATM usage. Traditional approaches to solving the problem typically comprise a  
18 succession of steps: localization, segmentation, feature extraction and recognition. This  
19 paper proposes an end-to-end facial occlusion detection framework, which is robust and  
20 effective by combining region proposal algorithm and Convolutional Neural Networks  
21 (CNN). The framework utilizes a coarse-to-fine strategy, which consists of two CNNs.  
22 The first CNN detects the head element within an upper body image while the second  
23 distinguishes which facial part is occluded from the head image. In comparison with  
24 previous approaches, the usage of CNN is optimal from a system point of view as the  
25 design is based on the end-to-end principle and the model operates directly on image  
26 pixels. For evaluation purposes, a face occlusion database consisting of over fifty thou-  
27 sand images, with annotated facial parts, was used. Experimental results revealed that  
28 the proposed framework is very effective. Using the bespoke face occlusion dataset, Alex  
29 and Robert (AR) face dataset and the Labeled Face in the Wild (LFW) database, we  
30 achieved over 85.61%, 97.58% and 100% accuracies for head detection when the Inter-  
31 section over Union-section (IoU) is larger than 0.5, and 94.55%, 98.58% and 95.41%  
32 accuracies for occlusion discrimination, respectively.

33 *Keywords:* Automated Teller Machine (ATM); Convolutional Neural Network (CNN);  
34 Face occlusion detection; Multi-Task Learning (MTL)

#### 35 **1. Introduction**

36 Automated Teller Machines (ATMs) have always been the targets of criminal ac-  
37 tivity since their widespread introduction in the 1970s. For example, fraudsters can  
38 obtain card details and PINs using a wide range of tactics. Among the possible

\*Corresponding author. Tel:+86 512 88161502.

39 techniques to defend against ATM crime, real time automatic alarm systems seem  
40 to be the most straightforward technical solution to maximize protection. This is  
41 because the surveillance cameras are installed in nearly all ATMs. However, cur-  
42 rent video surveillance for ATM requires constant staff monitoring, which has the  
43 obvious disadvantages of human error caused by fatigue or distraction.

44 Face occlusion detection has been studied for several years with a number of  
45 methods published [26, 31], many of which aim to reinforce ATM security. The  
46 published approaches can be roughly categorized into two categories: face or head  
47 detection approaches and occlusion classification approaches.

48 In the first category, the objective is robust face detection algorithms in the p-  
49 resence of partial occlusions, with two common practices, namely, facial component-  
50 based approaches and shape-based approaches. Facial component-based approach,  
51 such as that presented [20], detected facial components such as eyes, nose and  
52 mouth, and determines a face area based on the component detection result. For  
53 example, the method proposed in [20] combined seven AdaBoost-based classifier-  
54 s for whole face with individual face-part classifiers trained on non-occluded face  
55 sample sets, and a decision tree and Linear Discriminant Analysis (LDA) to classify  
56 non-occluded faces and various types of occluded faces. Inspired by discriminative-  
57 ly trained part-based models, Ahmed [7] proposed a Selective Part Model (SPM)  
58 to detect faces under certain types of partial occlusions. Gul [15] applied what is  
59 known as the Viola-Jones approach [38], with free rectangular features, to detect  
60 left half faces, right half faces and the holistic faces. AdaBoost-based face detection  
61 was also improved upon in [5] in order to detect partially occluded faces, which,  
62 however, only worked for frontal faces with sufficient resolutions.

63 Shape-based approaches [3, 4, 17, 21, 28] detect faces based on the prior knowl-  
64 edge of head, neck and shoulder shapes. In the scenario of ATM video surveillance,  
65 [28] proposed to compute the lower boundary of the head by moving object edge ex-  
66 traction and head tracking. Motion information was also exploited in [21] to detect  
67 the head and shoulder shape with the aid of B-spline active contouring. Color has  
68 also been applied as a major clue to detect the head or face, following appropriate  
69 template fitting strategies [3, 4, 17]. These approaches, however, are limited to con-  
70 strained poses and well-controlled illumination conditions. In [3, 4, 17] the head or  
71 shoulder were detected by using ellipse or what are known as “omega templates“,  
72 which, however, will fail when a face is severely occluded and/or the shapes are  
73 severely changed with different face poses.

74 Some of the previous published researches have tried to solve the face occlusion  
75 detection problem by straightforward classification [22]. For example, by separating  
76 a face area into upper and lower parts, Principal Component Analysis (PCA) and  
77 Support Vector Machine (SVM) were combined to distinguish between normal faces  
78 and partially occluded faces in [22].

79 Until recently, the most successful approaches to object detection utilized the  
80 well-known sliding window paradigm [10], in which a computationally efficient clas-  
81 sifier tests for object presence in every candidate image window. The steady increase

82 in complexity of the core classifiers has led to improved detection quality, but at the  
83 cost of significantly increased computation time per window [6, 11, 16, 35, 40]. One  
84 approach for overcoming the tension between computational tractability and high  
85 detection quality is through the use of "detection proposals" [8, 37]. If high object  
86 recall can be reached with considerably fewer windows than used by sliding win-  
87 dow detectors, significant performance improvement can be achieved. Current top  
88 performing object detectors, when applied to PASCAL benchmark image datasets  
89 [9] and ImageNet [30], all used detection proposals [6, 11, 13, 16, 35, 40]. According  
90 to [18], approaches for generating object proposals can be divided into four type-  
91 s: grouping methods, window scoring methods, alternative methods and baseline  
92 methods. Grouping methods attempt to generate multiple (possibly overlapping)  
93 segments that are likely to correspond to objects. Window scoring methods are used  
94 to score each candidate window according to how likely it is to contain an object.  
95 Inspired by the success of applying object proposal approaches in different object  
96 detections, this paper proposes a face occlusion detection system using the highly  
97 ranked object proposal technique, EdgeBoxes [47].

98 Over the last several years, there has been increasing interests in deep neural  
99 network models for solving various vision problems. One of the most successful  
100 deep learning frameworks is the CNN architecture [24], which is a bio-inspired hier-  
101 archical multilayered neural network that can learn visual representations directly  
102 from raw images. CNN possesses some key properties, namely translation invariance  
103 and spatially local connections (receptive fields). Pre-trained CNN models can be  
104 exploited as generic feature extractors for different vision tasks [24]. Among the var-  
105 ious advantages of deep neural networks over classical machine learning techniques,  
106 the most frequently cited examples include the conveniences for the implementation  
107 of knowledge transfer, Multi-Task Learning (MTL), attribute learning, multi label  
108 classification, and weakly supervised learning.

109 In this paper, a novel CNN based approach to face occlusion detection is pro-  
110 posed. A CNN cascade paradigm is adopted, which tackles the occlusion detection  
111 problem in a coarse-to-fine manner. The first CNN implements head/shoulder de-  
112 tection by taking a person's upper body image as input. The second CNN takes the  
113 output of the previous CNN as input and locates and classifies different facial parts.  
114 To facilitate the study of various face occlusion problems, a database directed at  
115 different kind of facial occlusions was created. The database consists of over fifty  
116 thousand images which are demarcated with four facial parts: two eyes, nose and  
117 mouth.

118 To the best of our knowledge, this is the first work directed at analyzing how  
119 face occlusion can be detected using multi-task CNN. The approach was verified  
120 on our face occlusion dataset, AR dataset [29] and LFW dataset [19], obtaining  
121 94.55%, 95.58% and 95.41% accuracies, respectively.

122 The rest of the paper is organized as follows. Section 2 provides the problem  
123 descriptions with the introduction of our face occlusion dataset. Section 3 overviews  
124 the proposed method and elaborates on the details of the coarse-to-fine framework.

4 *Yizhang Xia, Frans Coenen, and Bailing Zhang*

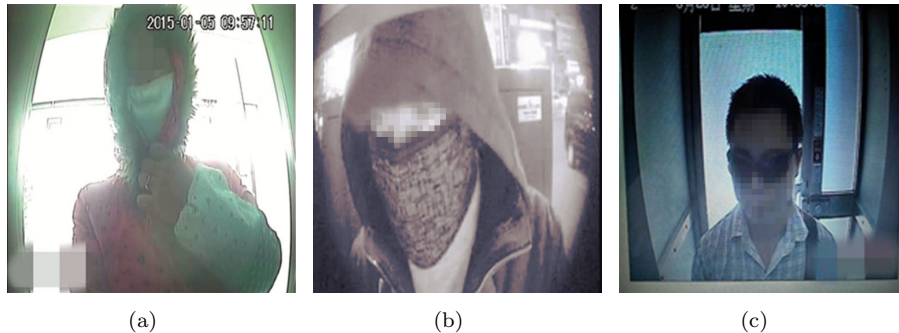


Fig. 1. Face occlusions examples for the withdrawing cash from an ATM scenario.

125 Section 4 reports the experiment results, followed by conclusion in Section 5.

## 126 2. Face Occlusion Dataset

127 A normal face image consists of two eyes, a nose and a mouth. The geometrical  
 128 and textual information from these components is critical to the recognition or  
 129 identification of a person. If any of these components is blocked or covered, the face  
 130 image is considered to be occluded. Commonly found face occlusions include the  
 131 faces being partially covered by a hat, sunglasses, mask or muffler. Some examples  
 132 are given in Fig. 1.

133 To facilitate research on face occlusion detection, a database of face images  
 134 which has different facial regions intentionally covered or occluded was created.  
 135 The images were taken using the Microsoft webcam studio camera and AMCap  
 136 9.20 edition. Some example images are illustrated in Figs. 2 and Fig. 3. The video  
 137 was filmed with a white background. There were 220 people involved, including 140  
 138 males and 80 females.

139 During the photography, a subject was asked to stand in front of the camera  
 140 with a set of different specified poses, including looking right ahead, up and down  
 141 roughly 45 degree, and right and left roughly 45 degree. In addition to taking face  
 142 images without any occlusions, a subject was also asked to wear sunglasses, hat (in  
 143 yellow and black), white mask and black helmet, again in five different poses. Each  
 144 subject has 6 video clips recorded, with 30 seconds for each clip and 25 frames per  
 145 second. The illumination was normal office lighting condition. Other conditions,  
 146 for example, clothing, make-up, hair style and expression, have not been strictly  
 147 controlled.

148 Although the created dataset is comprehensive, it cannot become public currently  
 149 due to legal considerations. To alleviate the problem, we also utilized the AR  
 150 face database [29], which is one of the earliest and most popular benchmark face  
 151 databases. AR faces have been often used for robust face recognition. It includes a  
 152 number of different types occlusions: faces with sunglasses and face partially cov-



Fig. 2. Upper body of our created dataset



Fig. 3. Head of our created dataset

153 ered by a scarf. The AR faces dataset consists of over 3200 frontal face images taken  
 154 from 126 subjects, with some examples given in Fig. 4.

155 The LFW dataset [19] was also employed to further evaluate our approach.  
 156 There are 13000 images and 1680 people in the LFW dataset collected from the  
 157 Internet. The faces were detected by the OpenCV implementation of the Viola-  
 158 Jones [38] face detector. The cropped region returned by the detector was then  
 159 automatically enlarged by a factor of 2.2 in each dimension to capture more of the  
 160 head and then scaled to a uniform size, 250\*250. For the evaluation presented in this  
 161 paper, 1000 images were selected and the heads manually cropped. Occlusions were  
 162 created using black rectangles and facial landmark localization [45]. Some examples  
 163 are given in Fig. 5.

6 *Yizhang Xia, Frans Coenen, and Bailing Zhang*

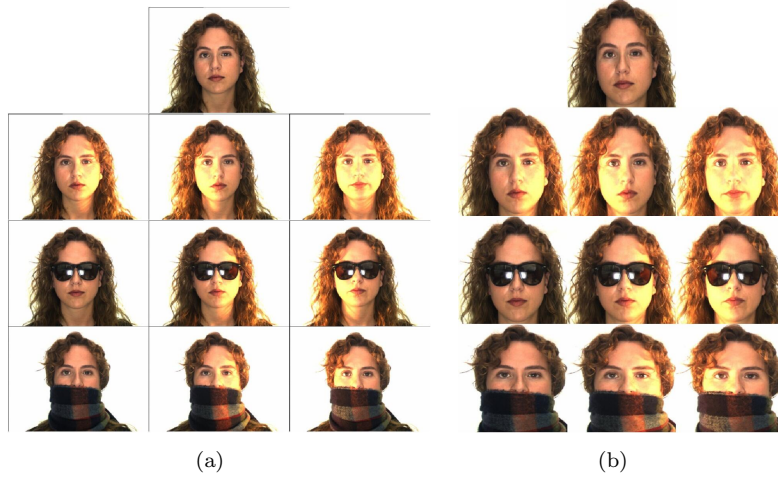


Fig. 4. Examples of face from the AR face Database. (a) head and shoulder, (b) head.



Fig. 5. Examples of images from the LFW database with occlusions superimposed

164 **3. System Overview**

165 Though deep neural networks have achieved remarkable performance, it is still  
 166 difficult to solve many real-world problems by a single CNN model. With the current  
 167 technology, the resolution of an input to CNN must be relatively small. This will  
 168 cause some details of an image lost. To get better performance, a common practice  
 169 of coarse-to-fine paradigm has been applied with CNN design [34, 48]. Following  
 170 the same line of thought, we proposed a two-stage convolutional neural network for  
 171 face occlusion detection, as illustrated in Fig. 6. The first CNN detects the head  
 172 from a person's upper body image while the second CNN distinguishes which facial  
 173 part is occluded from the head image.

174 **3.1. Head Detection**

175 To identify the locations of the head in an image, advantages were taken of Region  
 176 with Convolutional Neural Networks (R-CNN) [12], which is the state-of-the-art

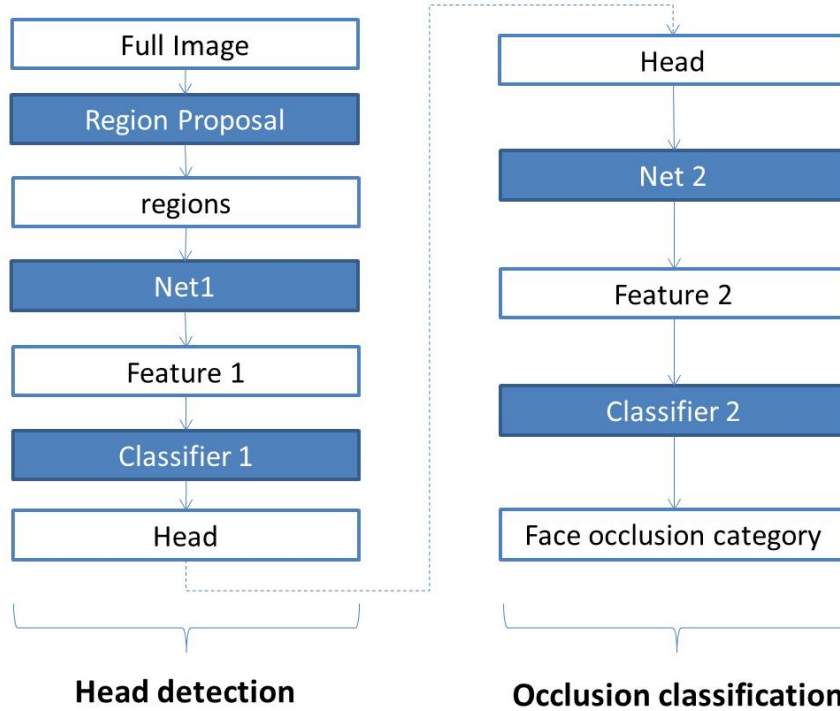


Fig. 6. The flow diagram of whole system

177 object detector that classifies candidate object hypothesis generated by appropriate  
 178 region proposal algorithms. The R-CNN leverages several advantages of computer  
 179 vision development and CNN, including the superb feature expression capability  
 180 from a pre-trained CNN, fine-tuning flexibility for specific objects to be detected  
 181 and the ever-increasing efficiency of object proposal generation schemes.

182 Among the off-the-shelf object proposal generation algorithms, the EdgeBoxes  
 183 technique [47] was chose which has attracted much interest in recent years. Edge-  
 184 Boxes is built on the Structural Edge Map to locate object boundaries and find  
 185 object proposals. The number of enclosed edges inside a bounding box is used to  
 186 rank the likelihood of the box containing an object.

187 The overall convolutional net architecture is shown in Fig. 7. The network  
 188 consists of three convolution stages followed by three fully connected layers. A  
 189 convolution stage includes a convolution layer, a non-linear activation layer, a lo-  
 190 cal response normalization layer and a max pooling layer. The non-linear activa-  
 191 tion layer and local response normalization layers are not included in Fig. 7 and  
 192 Fig. 8 as data size was not changed. Using shorthand notation, the full architec-  
 193 ture is  $C(8,5,1)-\tilde{A}-N-P-C(16,5,1)-\tilde{A}-N-P-C(32,5,1)-\tilde{A}-N-P-FC(1024)-\tilde{A}-FC(128)-\tilde{A}-$   
 194  $FC(4)-\tilde{A}$ , where  $C(d,f,s)$  indicates a convolutional layer with  $d$  filters of spatial size

8 *Yizhang Xia, Frans Coenen, and Bailing Zhang*

195  $f \times f$ , applied to the input with stride  $s$ .  $A$  is the non-linear activation function,  
 196 which uses the ReLU activation function[14].  $FC(n)$  is a fully connected layer with  
 197  $n$  output nodes. All pooling layers  $P$  use max-pooling in non-overlapping  $2 \times 2$  re-  
 198 gions and all normalization layers  $N$  are defined as described in Krizhevsky et al.  
 199 [24]. The final layer is connected to a soft-max layer with dense connections. The  
 200 structure of the networks and the hyper-parameters were empirically initialized  
 201 based on previous works using CNNs.

202 The well-known overfitting problem has been taken into account in our design  
 203 with the following considerations. Firstly, we empirically compared a set of different  
 204 CNN architectures with varying number of kernels and selected the one which is  
 205 deemed as the most effective with regard to the trade-off between network complex-  
 206 ity and performance. Secondly, the overfitting has been avoided to a large extent as  
 207 the CNN size is very moderate, which compares sharply with some published CNN  
 208 models for large-scale image classifications [24].

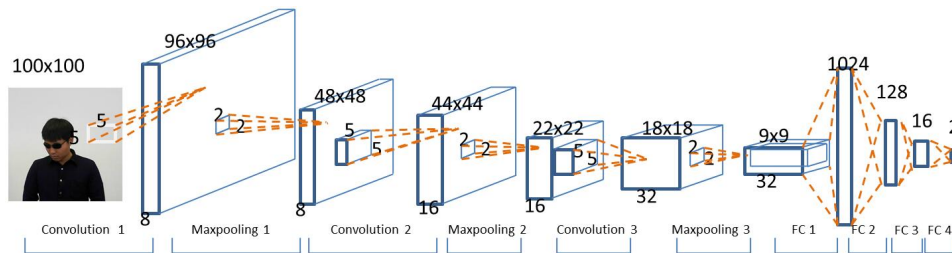


Fig. 7. Architecture of head detector CNN

209 The adopted CNN used the shared weight neural network architecture [25], in  
 210 which the local receptive field (kernel or filter) is replicated across the entire visual  
 211 field to form a feature map, which is known as convolution operation. The sharing  
 212 of weights reduces the number of free variables, and increases the generalization  
 213 performance of the network. Weights (kernels or filters) are initialized at random  
 214 and will learn to be edge, color or specific pattern detectors.

215 In deep CNN, the classical sigmoidal function has been replaced by a Rectifier  
 216 Linear Unit(ReLU) to accelerate training speed. Recent CNN-based approaches  
 217 [23, 24, 33, 36, 41, 43] applied the ReLU as the nonlinear activation function for  
 218 both the convolution layer and the full connection layer, often with faster training  
 219 speed as reported in [24].

220 Typical pooling functions include average-pooling and max-pooling layers. Av-  
 221 erage pooling takes the arithmetic mean of the elements in each pooling region  
 222 while max-pooling selects the largest element from the input.

223 The four layers of convolution, nonlinear activation, pooling and normalization,  
 224 are combined hierarchically to form a convolution stage (block). Generally, an input  
 225 image will be passed through several convolution stages for extracting complex



226 descriptive features. In the output of the topmost convolution stage, all small-sized  
 227 feature maps are concatenated into a long vector. Such a vector plays the same  
 228 role as hand-coded features and it is fed to a full connection layer. A standard full  
 229 connection operation can be either the conventional Multi-Layer Perceptron (MLP)  
 230 or SVM.

231 As the fully connected layers receive feature vector from the topmost convolution  
 232 stage, the output layer can generate a probability distribution over the output  
 233 classes. Toward this purpose, the output of the last fully-connected layer is fed to  
 234 a K-way softmax (where K is the number of classes) layer, which is the same as a  
 235 multi-class logistic regression.

### 236 3.2. Face Occlusion Classification

237 The second stage CNN takes the output from head detector as input and implements  
 238 the face occlusion classification. The CNN trains the classifier with the implicit,  
 239 highly discriminative features to differentiate facial parts and distinguish whether  
 240 a facial part is occluded or not at the same time. This is aided by a multi-task  
 241 learning paradigm described in more detail below.

242 The intuition of multi-task learning is to jointly learn multiple tasks by exploit-  
 243 ing a shared structural representation and improving the generalization by using  
 244 related tasks. Deep neural networks such as CNN have been proven advantageous  
 245 in multi-task learning due to their powerful representation learning capability and  
 246 the knowledge transferability across similar tasks [32, 42, 44].

247 Inspired by the successes of CNN multi-task learning, we configured the second  
 248 stage CNN to enable its shared representation learned in the feature layers for two  
 249 independent MLP classifiers in the final layer. Specifically, we jointly trained the  
 250 facial parts (left eye, right eye, nose and mouth) classification and occlusion/non-  
 251 occlusion decision simultaneously.

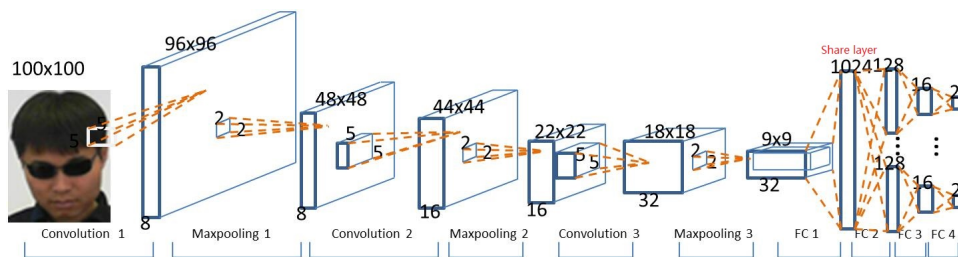


Fig. 8. Architecture of face occlusion classifier CNN

252 In our experiments, we adopt a multi-task CNN as shown in Fig. 8. The ar-  
 253 chitecture is same as for the CNN for head detector. When multi-task learning is  
 254 performed, we minimize the linear combination of individual task loss [42] as:

10 *Yizhang Xia, Frans Coenen, and Bailing Zhang*

$$L_{joint} = \sum_{i=1}^N \alpha_i L_i \quad (1)$$

255 where  $N$  is the total number of tasks,  $\alpha_i$  is weighting factor for the  $i$ -th task and  $L_i$   
 256 is the  $i$ -th task loss. When one of  $\alpha_i$ s takes 1, it will degenerate to classical single  
 257 task learning.

### 258 **3.3. Bounding Box Regression**

259 A bounding box regression module is employed to improve the detection accuracy.  
 260 Many bounding boxes generated from the EdgeBoxes algorithm are not close to  
 261 the object ground truth, but might be judged as positive samples. We can regard  
 262 detection as a regression problem to find the location of an object. This formulation  
 263 can work well for localizing a single object. Based on the error analysis, we imple-  
 264 mented a simple method to reduce localization errors. Inspired by the bounding-box  
 265 regression employed in the Deformable Parts Model (DPM) [10], we trained a linear  
 266 regression model to predict a new detection window given by the last pooled  
 267 features for the region proposals produced by the EdgeBoxes. This simple approach  
 268 can fix a large number of mislocalized detections, thus substantially boosting the  
 269 accuracy.

### 270 **3.4. Pre-train**

271 By supervised learning with sufficient annotated training data, CNNs that con-  
 272 tain millions of parameters have demonstrated competitive performance for visual  
 273 recognition tasks [23, 24, 33, 36, 41, 43] when starting from a random initializa-  
 274 tion. However, CNN architecture has a property that is strongly dependent on large  
 275 amounts of training data for good generalization. When the amount of labeled data  
 276 is limited, directly training a high capacitor CNN may become problematic. Re-  
 277 searches [39] have shown an alternative solution to compensate the problem by  
 278 choosing an optimised starting point which can be pre-trained by transferring pa-  
 279 rameters from either supervised learning or unsupervised learning, as opposite to a  
 280 random initialized start. We first trained the CNN model in the supervised mode  
 281 using the ImageNet data and then fine-tuned it on the domain-specific labeled  
 282 images as the head detector. To be specific, the pre-trained model is designed to  
 283 recognize objects in natural images. The leveraged knowledge from the source task  
 284 could reflect some common characteristics shared in these two types of images such  
 285 as corners or edges.

## 286 **4. Experiment Results**

287 The proposed approach was analyzed using our face occlusion dataset, the AR face  
 288 dataset and the LFW dataset with pre-training and fine-tuning. For head detector,

289 the pre-training stage employed 10% of images from the ILSVRC2012 [30] for the  
290 classification task. For the occlusion classification, the pre-training was implemented  
291 by the face recognition.

#### 292 **4.1. Implementation Details**

293 The experiments were conducted on a computer Dell Tower 5810 with Intel Xeon  
294 E5-1650 v3 and 64G of memory. In order to speed up CNN training, a GPU,  
295 NVIDIA GeForce GTX TITAN, is plugged on the board. The program operated  
296 under 64-bit windows 7 Ultimate with Matlab 2013b, Microsoft visual studio 2012  
297 and CUDA7.0[46].

298 We trained our models using stochastic gradient decent with a batch size of  
299 128 examples and momentum of 0.01. The learning rate was initialized as 0.01 and  
300 adapted during training. More specifically, we monitored the overall loss function.  
301 If the loss was not reduced for 5 epochs continually, the learning rate was dropped  
302 by 50%. We deem the network converged if the loss is stabilized.

303 In the training procedure, firstly, 10% of images from the ILSVRC2012 are  
304 used to train Net1 from random initialization. Then, the fine tuning on Net1 is  
305 implemented by head detection dataset in the AR face dataset, our dataset and the  
306 LFW dataset. Thirdly, Net2 is initialized from the Net1 after fine turning. Finally,  
307 the fine tuning on Net2 is trained by face occlusion classification dataset in the AR  
308 face dataset, our dataset and the LFW dataset. And the test procedure is described  
309 in Fig. 6.

#### 310 **4.2. Head Detector**

311 As the state-of-the-art object proposal method with regard to the trade-off between  
312 the speed and recall, EdgeBoxes method needs to tune the parameters at each  
313 desired overlap threshold [47]. We estimated three pairs of  $\alpha$  and  $\beta$  parameters to  
314 evaluate the object proposal corresponding to the head hypothesis, following the  
315 standard object proposal evaluation framework. The detection recall is calculated  
316 as the ratio of ground truth bounding boxes that have been predicted among the  
317 EdgeBoxes proposals with an Intersection Over Union (IoU) larger than a given  
318 threshold.

319 Three useful pairs of  $\alpha$  and  $\beta$  values in EdgeBoxes [47] were estimated to provide  
320 the head candidates from our dataset, the AR face dataset and the LFW dataset,  
321 as shown in Fig. 9. The parameters  $\alpha$  and  $\beta$  control the step size of the sliding  
322 window search and the NMS threshold, respectively. The two figures in Fig. 9  
323 illustrate the algorithm's behavior with varying  $\alpha$  and  $\beta$  when the same number  
324 of object proposals are generated. More specifically, three pairs of  $\alpha$  and  $\beta$  are  
325 0.65/0.55, 0.65/0.75, and 0.85/0.95 respectively. They were tested when the max  
326 region proposal was fixed at 500 with respect to our database and the AR database  
327 and 700 on the LFW database. It is obvious that the biggest value obtain the best  
328 recall, such that we were able to achieve 97.53% on our face occlusion database,

12 *Yizhang Xia, Frans Coenen, and Bailing Zhang*

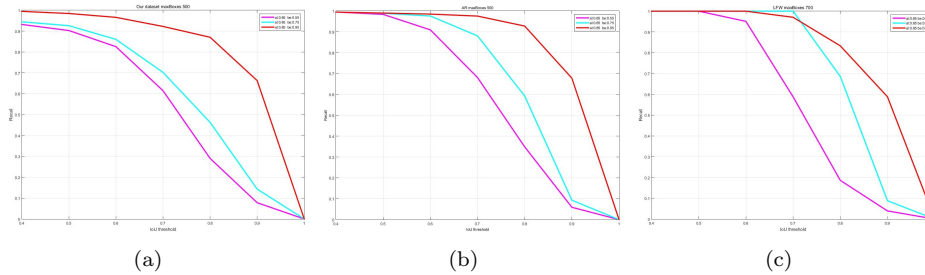


Fig. 9. The three useful variants of EdgeBoxes when the max number of region proposal is same.

329 99.35% on the AR face database and 100% on the LFW database when IoU is  
 330 larger than 0.5. In conclusion, when the density of the sampling rate increases, we  
 331 will get higher recall, but the run time becomes longer.

332 The second experiment compared the different number of bounding boxes pro-  
 333 duced by EdgeBoxes according to ranking score. A suite of maximal number of  
 334 boxes, 100, 300, 500, 700, 900, was evaluated when the density of the sampling rate  
 335 and NMS threshold are fixed to 0.85 and 0.95 respectively, as shown in Fig. 10. In  
 336 order to make sure the face occlusion classification is valid, the IoU between region  
 337 and ground truth is larger than 0.5. The recall is approximately 100% on our face  
 338 occlusion database, the AR face database and the LFW database (Fig. 10). The re-  
 339 call does not increase when the maximum number of region proposal is larger than  
 340 500 on our dataset and the AR dataset and 700 on the LFW dataset. In consider-  
 341 ation of the tolerance of Net1, the  $\alpha$ ,  $\beta$  and max number of region proposal were  
 342 set to 0.85, 0.95 and 500 on our dataset and AR dataset respectively. By taking  
 343 account of the complex background of LFW dataset, the maximal number of boxes  
 344 was selected as 700 and the others parameters were the same for our dataset and  
 345 AR dataset. The same parameters were selected for our face occlusion database and  
 346 the AR face database. However, the recall obtained with respect to the AR dataset  
 347 is lower due to the greater variability of our database. The recall curves (Fig. 9 and  
 348 Fig. 10) demonstrate that almost all of the head has been selected as proposals by  
 349 the EdgeBoxes.

350 After obtaining the candidate regions from EdgeBoxes, the regions are classified  
 351 into head and non-head by the trained CNN module. Before applying the CNN for  
 352 the classification, a hypothesis object proposal will be discarded if the proposal can  
 353 be simply judged as head or useless patch based on the following reasoning: a region  
 354 is head if the overlap with the whole image is less than 5% on AR face database,  
 355 a region is head or useless patch if the overlapping with the full image is between  
 356 2% and 30% on our face occlusion database. Then, the features are extracted via  
 357 Net1 from normalized regions,  $100 \times 100$  gray images. Next, the subsequent MLP  
 358 predicts whether the region is a head region or not. After MLP, several positive  
 359 regions are merged into a box by the non-max suppression with 0.3 overlap threshold

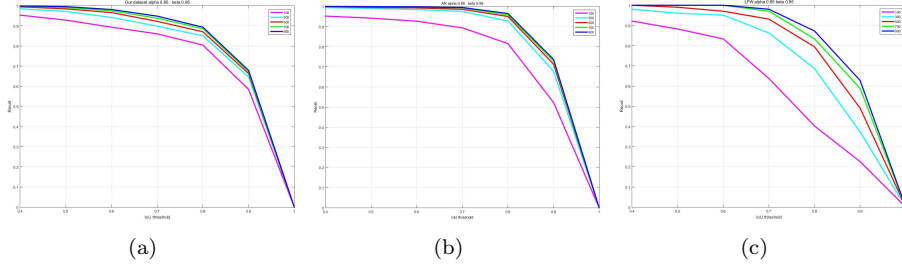


Fig. 10. The difference maximal number of region proposal when the  $\alpha$  is 0.85 and  $\beta$  is 0.95.

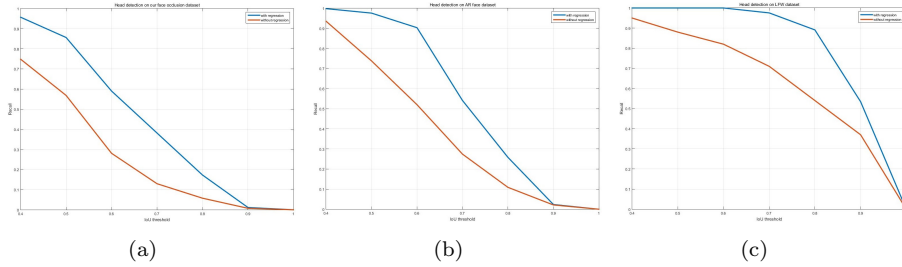


Fig. 11. The recall of head detection, the red line is the performance of head detector without regression and the blue line indicate the performance of head detector with regression.

360 sorted by the reliability of the MLP, which indicated the accuracy of detection.  
 361 The performance of the approach was tested on the AR face database and our face  
 362 occlusion database. The accuracies using our face occlusion database, the AR face  
 363 database and the LFW database were 56.83%, 73.79% and 88.52% with 0.5 IoU,  
 364 respectively (Fig. 10). The performance can be further improved by bounding box  
 365 regression, as expounded in the following.

366 We use a simple bounding-box regression stage to improve the localization per-  
 367 formance. After scoring each object proposal by MLP, we predict a new bounding  
 368 box for detection using a class-specific bounding-box regressor. This is similar in  
 369 spirit to the approach used in deformable part models [12]. The better location of  
 370 head is regressed from features computed by the CNN (Fig. 11). After regressing  
 371 the positive features from CNN, the performance is improved, with increases of  
 372 28.78%, 23.79% and 11.48% on our face occlusion database, AR face database and  
 373 LFW database respectively with IoU 0.5. The reason why the regressor improves  
 374 the performance is that the head is centered in the images. However, the regres-  
 375 sor only improves when the classification result from CNN is correct. Examples of  
 376 the head detection from our face occlusion database, AR face database and LFW  
 377 database are shown in the Fig. 12, Fig. 13 and Fig. 14, respectively.

378 We created a head detector based on HoG feature extraction [1] and SVM  
 379 classification [2] as a baseline approach. The parameters for the HOG was set as

14 *Yizhang Xia, Frans Coenen, and Bailing Zhang*

380 orientations in the range  $[0, 360]$ , 40 orientations bins and 3 spatial levels, which  
 381 means that the inputs for SVM have a dimension of 840,  $(1 + 4 + 16) \times 40$ . For  
 382 the SVM, we employed C-SVC and radial basis kernel function [2]. The comparison  
 383 performance is illustrated in Table 1, which demonstrates that our approach is  
 384 robust and effective in complex scenes. The computational complexity is showed  
 385 Table 2. The major part of the computation of our approach is from the region  
 386 proposal algorithm, EdgeBoxes.

Table 1. Accuracy of head detection when IoU is larger than 0.5

	LFW	AR	Our dataset
HOG+SVM	91.87%	97.44%	72.41%
Our method	100%	97.58%	85.61%

Table 2. Computational complexity on head detection(fps)

	LFW	AR	Our created dataset
HOG+SVM	32.48	24.06	419.68
Our method	0.94	2.02	1.76

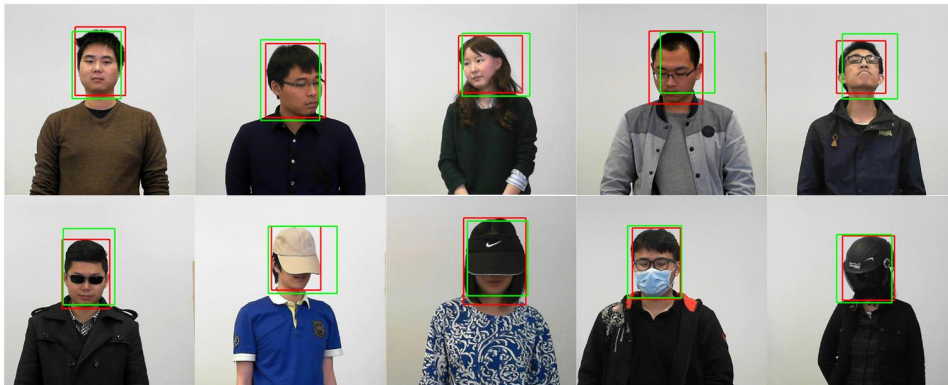


Fig. 12. Examples of head detection results from our face occlusion database, red rectangle is ground truth and the green rectangle is detection location.

### 387 **4.3. Face Occlusion Classification**

388 After obtaining the head position by the head detector, the face occlusion classifier  
 389 was used to classify the type of face occlusion.



Fig. 13. Examples of head detection results from the AR face database, the red rectangle is the ground truth and the green rectangle is the detection location.

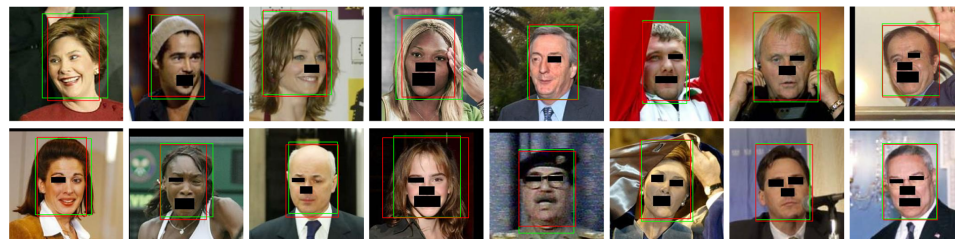


Fig. 14. Examples of head detection results from the LFW database, the red rectangle is the ground truth and the green rectangle is the detection location.

390 The AR face dataset [29] is a popular dataset for face recognition. It contains  
 391 over 4,000 color images from 126 people's faces (70 males and 56 females). Images  
 392 cover the frontal view faces with occlusions (sun glasses and scarf), different facial  
 393 expressions and illumination conditions (Fig. 1). For the face occlusion classifier,  
 394 the dataset is categorized into two occlusion conditions: face with eyes occluded by  
 395 sunglasses and face with mouth occluded by scarf. Half of the un-occluded faces were  
 396 use to pre-train the CNN model, following a general face recognition methodology.  
 397 More specifically, the pre-train dataset consisted of 31 male and 25 female faces  
 398 (frontal view face with different expressions and illumination).

399 After pre-training, the fully connected layers were replaced by a new MLP,  
 400 initialized from random connection values for the fine-tuning, by applying the face

16 *Yizhang Xia, Frans Coenen, and Bailing Zhang*

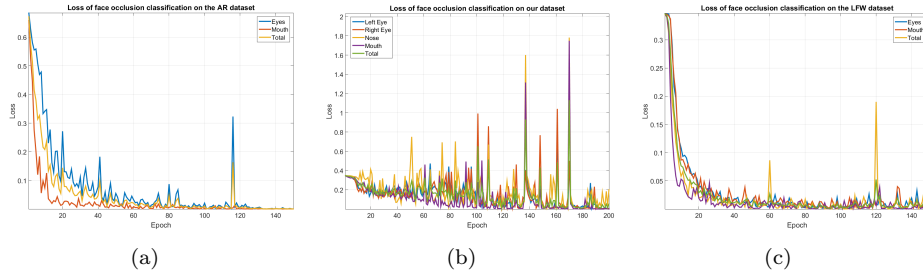


Fig. 15. Comparison of loss function values during training

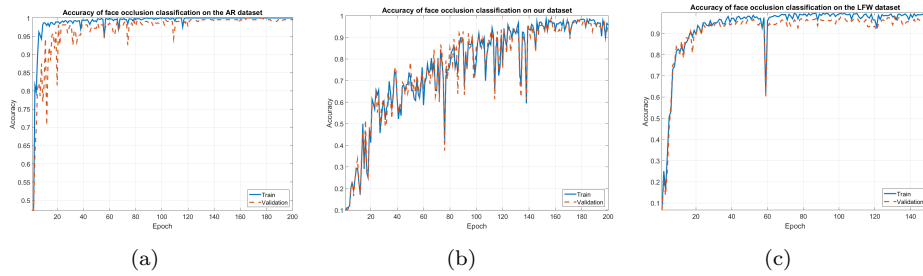


Fig. 16. Comparison of accuracy values during training and validation

401 occlusion dataset. The CNN with multi-tasks learning was designed to predict the  
 402 presence or absence of different facial parts, thus indicating occlusion or not. The  
 403 structure of CNN is illustrated in Fig. 8.

404 The loss function values during training is illustrated in Fig. 15. After the loss  
 405 becomes stabilized, the converged model is applied to test a testing sample with  
 406 accuracy as shown in Table 3.

407 The accuracy values during training and validation are given in Fig. 16. Cross  
 408 validation and early stopping are employed to avoid overfitting. Fig. 16 illustrates  
 409 that the accuracies during training and validation continually rise. Early stopping is  
 410 used to select the appropriate trained model during training and avoid overfitting.

Table 3. Face occlusion classification on AR dataset

	eyes	Mouth	Total
Accuracy	98.58%	100%	98.58%

411 The system performance for face occlusion classification was also evaluated using  
 412 our face occlusion dataset and the LFW dataset. The experiment procedures are  
 413 similar to the AR face dataset. The difference is that there is a variety of occlusions  
 414 from different facial parts, namely, left eye, right eye, nose and mouth. Accordingly,



415 there are four MLPs following the shared layer to verify whether each facial part is  
 416 occluded or not, as explained in Fig. 8.

417 The training loss is shown in Fig. 15, and the occlusion accuracies on our dataset  
 418 and LFW are 94.55% and 95.41% respectively, with further details provided in  
 419 Table 4. Our model is a multi-task framework that shares the front layers. The  
 420 summed loss changes with the updating of CNN parameters. And it can be seen  
 421 that the loss fluctuates more obviously on our face dataset because it is much more  
 422 complex compared with AR dataset and LFW dataset.

423 The face detector combined with Haar feature [27] extractor and Viola-Jones  
 424 [38] classifier is used as the base line of face occlusion classification. The experiments  
 425 result is summarized in Table 5, which indicates that our method outperforms Haar-  
 426 based face detection. And the corresponding computational complexity is reported  
 427 in Table 6 showing that our method is faster than the classical approach.

Table 4. The Accuracy on LFW and our face occlusion dataset

The Accuracy on our face occlusion dataset					
	Left eye	Right eye	Nose	Mouth	Total
Our dataset	98.15%	99.07%	98.15%	99.07%	94.55%
LFW	97.79%	98.91%	99.63%	99.02%	95.41%

Table 5. Accuracy on face occlusion classification

	LFW	AR	Our dataset
Haar+VJ	45.22%	47.98%	21.20%
Our method	95.41%	98.58%	94.55%

Table 6. Computational complexity on face occlusion classification(ms)

	LFW	AR	Our created dataset
Haar+VJ	16.29	58.23	22.35
Our method	13.10	20.67	17.31

#### 428 4.4. Error Analysis

429 There are two factors that may impact the effectiveness of the proposed approach.  
 430 Firstly, the region proposal algorithm should be robust with regard to the gener-  
 431 ated candidate bounding boxes. However, as the EdgeBoxes produce the candidate

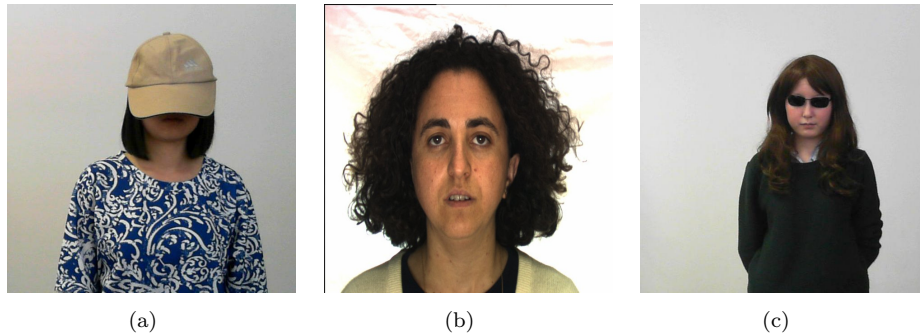
18 *Yizhang Xia, Frans Coenen, and Bailing Zhang*

Fig. 17. Examples of wrong region proposal

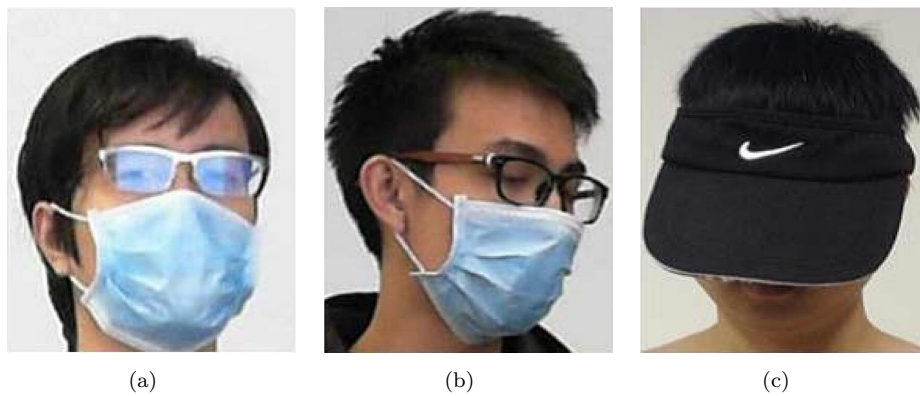


Fig. 18. Error prediction image sample

432 regions by the edge, it has several potential problems: (a) negative data will be  
 433 generated when a person wears clothing with complex texture (Fig. 17(a)); (b) a  
 434 head will appear larger with certain hairstyles (Fig. 17(b)); (c) the head will not be  
 435 segmented when a person has long hair and wears a black dust coat (Fig. 17(c)).  
 436 Secondly, there are three factors that may influence the performance of face occlusion  
 437 classifier, including the illumination variations, occlusions and the partial  
 438 occlusions. Fig. 18 further explains the difficult situations.

## 439 5. Conclusion

440 This paper proposed an approach for face occlusion detection to enhance the surveil-  
 441 lance security for ATM. The coarse-to-fine approach consists of a head detector and  
 442 a face occlusion classifier. The head detector is implemented with EdgeBoxes — re-  
 443 gion proposal, CNN and MLP. The method achieved detection accuracies of 97.58%,  
 444 85.61% and 100% on the AR face database, our face occlusion database and LFW  
 445 dataset. For face occlusion classification, CNN is applied with a pre-training strat-

446 egy via usual face recognition task, followed by fine-tuning with the face occlusion  
447 classification based on MTL to verify whether a facial part is occluded or not.  
448 Our approach is evaluated on the AR face dataset, our dataset and LFW dataset,  
449 achieving 98.58%, 94.55% and 95.41% accuracies, respectively. Further work is be-  
450 ing made toward the improvement of the model by more robust and accurate region  
451 proposal, which will render it more realistic for real-world applications.

## 452 6. Acknowledgement

453 The first author thanks Chao Yan and Rongqiang Qian for their valuable helps.  
454 The constructive comments and suggestions from the anonymous reviewers are  
455 much appreciated.

## 456 References

- 457 1. A. Bosch, A. Zisserman and X. Munoz, Representing shape with a spatial  
458 pyramid kernel, *the 6th ACM International Conference on Image and Video  
459 Retrieval*, pp. 401–408, New York, 2007.
- 460 2. C.-C. Chang and C.-J. Lin, Libsvm: A library for support vector machines,  
461 *ACM Trans. Intell. Syst. Technol.*, Vol. 2, pp. 27:1–27:27, 2011.
- 462 3. T. Charoenpong, C. Nuthong and U. Watchareeruetai, A new method for oc-  
463 cluded face detection from single viewpoint of head, *2014 11th International  
464 Conference on Electrical Engineering/Electronics, Computer, Telecommunica-  
465 tions and Information Technology (ECTI-CON)*, pp. 1-5, 2014.
- 466 4. T. Charoenpong, Face occlusion detection by using ellipse fitting and skin col-  
467 or ratio, *Burapha University International Conference (BUU)*, pp. 1145-1151,  
468 2013.
- 469 5. J. Chen, S. Shan, S. Yang, X. Chen and W. Gao, Modification of the adaboost-  
470 based detector for partially occluded faces, *18th International Conference on  
471 Pattern Recognition*, Vol. 2, pp. 516-519, 2006.
- 472 6. R. Cinbis, J. Verbeek and C. Schmid, Segmentation driven object detection  
473 with fisher vectors, *2013 IEEE International Conference on Computer Vision  
474 (ICCV)*, pp. 2968-2975, 2013.
- 475 7. A. El-Barkouky, A. Shalaby, A. Mahmoud and A. Farag, Selective part models  
476 for detecting partially occluded faces in the wild, *2014 IEEE International  
477 Conference on Image Processing (ICIP)*, pp. 268-272, 2014.
- 478 8. I. Endres and D. Hoiem, Category-independent object proposals with diverse  
479 ranking, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol.  
480 36, pp. 222–234, 2014.
- 481 9. M. Everingham, S. Eslami, L. Van Gool, C. Williams, J. Winn and A. Zisser-  
482 man, The pascal visual object classes challenge: A retrospective, *International  
483 Journal of Computer Vision*, Vol. 111, pp. 98-136, 2015.
- 484 10. P. Felzenszwalb, R. Girshick, D. McAllester and D. Ramanan, Object detection

## 20 REFERENCES

- 485 with discriminatively trained part-based models, *IEEE Transactions on Pattern*  
486 *Analysis and Machine Intelligence*, Vol. 32, pp. 1627-1645, 2010.
- 487 11. R. Girshick, J. Donahue, T. Darrell and J. Malik, Rich feature hierarchies for  
488 accurate object detection and semantic segmentation, *2014 IEEE Conference*  
489 *on Computer Vision and Pattern Recognition (CVPR)*, pp. 580-587, 2014.
- 490 12. R. Girshick, J. Donahue, T. Darrell and J. Malik, Region-based convolutional  
491 networks for accurate object detection and segmentation, *IEEE Transactions*  
492 *on Pattern Analysis and Machine Intelligence*, Vol. 38, PP. 142-158, 2016.
- 493 13. R. Girshick, Fast r-cnn, *The IEEE International Conference on Computer Vi-*  
494 *sion (ICCV)*, pp. 1440-1448, 2015.
- 495 14. X. Glorot, A. Bordes and Y. Bengio, Deep sparse rectifier neural networks,  
496 *Journal of Machine Learning Research*, Vol. 15, pp. 315-323, 2011.
- 497 15. S. Gul and H. Farooq, A machine learning approach to detect occluded faces in  
498 unconstrained crowd scene, *2015 IEEE 14th International Conference on Cog-*  
499 *nitive Informatics Cognitive Computing (ICCI\*CC)*, pp. 149-155, 2015.
- 500 16. K. He, X. Zhang, S. Ren and J. Sun, Spatial pyramid pooling in deep convolu-  
501 tional networks for visual recognition, *IEEE Transactions on Pattern Analysis*  
502 *and Machine Intelligence*, Vol. 37, pp. 1904-1916, 2015.
- 503 17. H. Sun, J. Wang, P. Sun and X. Zou, Facial area forecast and occluded face de-  
504 tection based on the ycbcr elliptical model, *International Conference on Mecha-*  
505 *tronic Sciences, Electric Engineering and Computer (MEC)*, pp. 1199-1202,  
506 2013.
- 507 18. J. Hosang, R. Benenson, P. Dollár and B. Schiele, What makes for effective  
508 detection proposals?, *IEEE transactions on pattern analysis and machine in-*  
509 *telligence* Vol. 38, pp. 814-830, 2016.
- 510 19. G. B. Huang, M. Ramesh, T. Berg and E. Learned-Miller, Labeled faces in the  
511 wild: A database for studying face recognition in unconstrained environments,  
512 *Technical Report*, University of Massachusetts, Amherst, pp. 07-49, 2007.
- 513 20. K. Ichikawa, T. Mita, O. Hori and T. Kobayashi, Component-based face detec-  
514 tion method for various types of occluded faces, *3rd International Symposium*  
515 *on Communications, Control and Signal Processing*, pp. 538-543, 2008.
- 516 21. G. Kim, J. K. Suhr, H. G. Jung and J. Kim, Face occlusion detection by using b-  
517 spline active contour and skin color information, *11th International Conference*  
518 *on Control Automation Robotics Vision (ICARCV)*, pp. 627-632, 2010.
- 519 22. J. Kim, Y. Sung, S. Yoon and B. Park, A new video surveillance system em-  
520 ploying occluded face detection, *Innovations in Applied Artificial Intelligence*,  
521 pp. 65-68, 2005.
- 522 23. J. Krause, T. Gebru, J. Deng, L.-J. Li and L. Fei-Fei, Learning features and  
523 parts for fine-grained recognition, *2014 22nd International Conference on Pat-*  
524 *tern Recognition (ICPR)*, pp. 26-33, 2014.
- 525 24. A. Krizhevsky, I. Sutskever and G. E. Hinton, Imagenet classification with  
526 deep convolutional neural networks, *Advances in Neural Information Processing*  
527 *Systems*, pp. 1097-1105, 2012.

- 528 25. Y. LeCun, Y. Bengio and G. Hinton, Deep learning, *Nature*, pp. 436-444, 2015.
- 529 26. S. Liao, A. K. Jain and S. Z. Li, Partial face recognition: Alignment-free ap-  
530 proach, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol.  
531 35, pp. 1193-1205, 2013.
- 532 27. R. Lienhart and J. Maydt, An extended set of haar-like features for rapid object  
533 detection, *International Conference on Image Processing*, Vol. 1, pp. 900-903,  
534 2002.
- 535 28. D.-T. Lin and M.-J. Liu, Face occlusion detection for automated teller machine  
536 surveillance, *Advances in Image and Video Technology*, pp. 641-651, 2006.
- 537 29. A. Martinez and R. Benavente, The ar face database, *CVC Technical Report*,  
538 1998.
- 539 30. O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang,  
540 A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg and L. Fei-Fei, ImageNet  
541 Large Scale Visual Recognition Challenge, *International Journal of Computer  
542 Vision (IJCV)* Vol. 115, pp. 211-252, 2015.
- 543 31. J. Shermine and V. Vasudevan, Recognition of the face images with occlusion  
544 and expression, *International Journal of Pattern Recognition and Artificial In-  
545 telligence* Vol. 26, 2012.
- 546 32. L. Sijin, L. Zhi-Qiang and A. Chan, Heterogeneous multi-task learning for hu-  
547 man pose estimation with deep convolutional neural network, *IEEE Conference  
548 on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 488-  
549 495, 2014.
- 550 33. K. Simonyan and A. Zisserman, Two-stream convolutional networks for action  
551 recognition in videos, *Advances in Neural Information Processing Systems*, pp.  
552 568-576, 2014.
- 553 34. Y. Sun, X. Wang and X. Tang, Deep convolutional network cascade for facial  
554 point detection, *IEEE Conference on Computer Vision and Pattern Recognition  
555 (CVPR)*, Oregon, pp. 3476-3483, 2013.
- 556 35. C. Szegedy, S. Reed, D. Erhan and D. Anguelov, Scalable, high-quality object  
557 detection, *arXiv preprint*, 2014.
- 558 36. Y. Taigman, M. Yang, M. Ranzato and L. Wolf, Deepface: Closing the gap to  
559 human-level performance in face verification, *IEEE Conference on Computer  
560 Vision and Pattern Recognition (CVPR)*, pp. 1701-1708, 2014.
- 561 37. K. van de Sande, J. Uijlings, T. Gevers and A. Smeulders, Segmentation as  
562 selective search for object recognition, *IEEE International Conference on Com-  
563 puter Vision (ICCV)*, pp. 1879-1886, 2011.
- 564 38. P. Viola and M. Jones, Rapid object detection using a boosted cascade of simple  
565 features, *Computer Vision and Pattern Recognition*, Vol. 1, pp. 511-518, 2001.
- 566 39. R. Wagner, M. Thom, R. Schweiger, G. Palm and A. Rothemel, Learning  
567 convolutional neural networks from few samples, *The 2013 International Joint  
568 Conference on Neural Networks (IJCNN)*, pp. 1-7, 2013.
- 569 40. X. Wang, M. Yang, S. Zhu and Y. Lin, Regionlets for generic object detection,  
570 *IEEE International Conference on Computer Vision (ICCV)*, pp. 17-24, 2013.

## 22 REFERENCES

- 571 41. S. Yi, W. Xiaogang and T. Xiaoou, Deep learning face representation from  
572 predicting 10,000 classes, *IEEE Conference on Computer Vision and Pattern  
573 Recognition (CVPR)*, pp. 1891-1898, 2014.
- 574 42. C. Zhang and Z. Zhang, Improving multiview face detection with multi-task  
575 deep convolutional neural networks, *IEEE Winter Conference on Applications  
576 of Computer Vision (WACV)*, pp. 1036-1041, 2014.
- 577 43. N. Zhang, M. Paluri, M. Ranzato, T. Darrell and L. Bourdev, Panda: Pose  
578 aligned networks for deep attribute modeling, *IEEE Conference on Computer  
579 Vision and Pattern Recognition (CVPR)*, pp. 1637-1644, 2014.
- 580 44. Z. Zhang, P. Luo, C. Loy and X. Tang, Facial landmark detection by deep  
581 multi-task learning, *European Conference on Computer Vision (ECCV)*, pp.  
582 94-108, 2014.
- 583 45. E. Zhou, H. Fan, Z. Cao, Y. Jiang and Q. Yin, Extensive facial landmark local-  
584 ization with coarse-to-fine convolutional network cascade, *IEEE International  
585 Conference on Computer Vision Workshops (ICCVW)*, pp. 386-391, 2013.
- 586 46. X. Zhu, K. Li, A. Salah, L. Shi and K. Li, Parallel implementation of mafft on  
587 cuda-enabled graphics hardware, *IEEE/ACM Trans. Comput. Biol. Bioinforma-  
588 tics* Vol. 12, pp. 205-218, 2015.
- 589 47. C. Zitnick and P. Dollar, Edge boxes: Locating object proposals from edges,  
590 *European Conference on Computer Vision (ECCV)*, pp. 391-405, 2014.
- 591 48. Z. Liu, P. Luo, X. Wang and X. Tang, Deep learning face attributes in the  
592 wild, *The IEEE International Conference on Computer Vision (ICCV)*, pp.  
593 3730-3738, 2015.