

Dual Attention MobDenseNet(DAMDNNet) for Robust 3D Face Alignment

Lei Jiang^{1,3}

Xiao-Jun Wu^{1,3,*}

Josef Kittler²

¹School of IoT Engineering, Jiangnan University 214122, Wuxi, China.

²Center for Vision, Speech and Signal Processing(CVSSP), University of Surrey, GU27XH, Guildford, UK.

³Jiangsu Provincial Engineering Laboratory of Pattern Recognition and Computational Intelligence, Jiangnan University, 214122, Wuxi, China.

ljiang-jnu@outlook.com xiaojun.wu-jnu@163.com j.kittler@surrey.ac.uk

Abstract

3D face alignment of monocular images is a crucial process in the recognition of faces with disguise. 3D face reconstruction facilitated by alignment can restore the face structure which is helpful in detecting disguise interference. This paper proposes a dual attention mechanism and an efficient end-to-end 3D face alignment framework. We build a stable network model through Depthwise Separable Convolution, Densely Connected Convolutional and Lightweight Channel Attention Mechanism. In order to enhance the ability of the network model to extract the spatial features of the face region, we adopt Spatial Group-wise Feature enhancement module to improve the representation ability of the network. Different loss functions are applied jointly to constrain the 3D parameters of a 3D Morphable Model (3DMM) and its 3D vertices. We use a variety of data enhancement methods and generate large virtual pose face data sets to solve the data imbalance problem. The experiments on the challenging AFLW, AFLW2000-3D datasets show that our algorithm significantly improves the accuracy of 3D face alignment. Our experiments using the field DFW dataset show that DAMDNNet exhibits excellent performance in the 3D alignment and reconstruction of challenging disguised faces. The model parameters and the complexity of the proposed method are also reduced significantly. The code is publicly available at <https://github.com/LeiJiangJNU/DAMDNNet>

1. Introduction

The aim of face alignment is to locate the feature points of the human face, such as the corners of the eyes, the corners of the mouth, tip of the nose. In general it involves fitting a face model to an image and extracting the semantic meaning of facial pixels. This is a fundamental step for many face analysis tasks, such as face recognition [6], face

expression analysis [3] and facial animation [10, 9]. In this paper we investigate face alignment in the context of face disguise detection. The problem of detecting a face disguise is concerned with determining whether a given pair of images belong to the same person even if one of them is subject to a disguise, or to different persons (one of them being an imposter). In view of the importance of this problem, face alignment has been widely studied since the Active Shape Model (ASM) of Cootes in the early 1990s [13]. Especially in recent years, face alignment has become a hot topic in computer vision.

The existing methods of face alignment can be divided into three categories: Constrained Local Model (CLM) methods (e.g., [13, 35]), Active Appearance Model (AAM) methods (e.g., [29, 30]) and regression methods (e.g., [10, 41]).

3D face shape reconstruction from 2D image is very challenging by its nature if no prior knowledge is provided. This is mainly because 2D data does not convey unambiguous depth information. A common method to solve the problem of monocular 2D face shape reconstruction is to use a set of 3D base shapes to capture the subspace, or a morphological model of face shape variations. Blanz and Vetter[5] proposed a comprehensive approach to minimizing the difference between the input 2D image and its 3D face rendering. Although this method has been able successfully to solve the problem of 3D face reconstruction, it is not friendly to changing lighting conditions, and its computational cost is high. To overcome this limitation, Blanz *et al.*[4] proposed to predict 3D parameters of a 3D face model from 2D facial feature points by linear regression. Although this method is efficient, it abandons the most useful information in the image and learns very simple regression functions. Recently, some innovative methods have been proposed, such as estimating 3D parameters through CNN and related cascaded regression operations, to achieve 3D face reconstruction. However, the network structure used by these methods is complex and the model

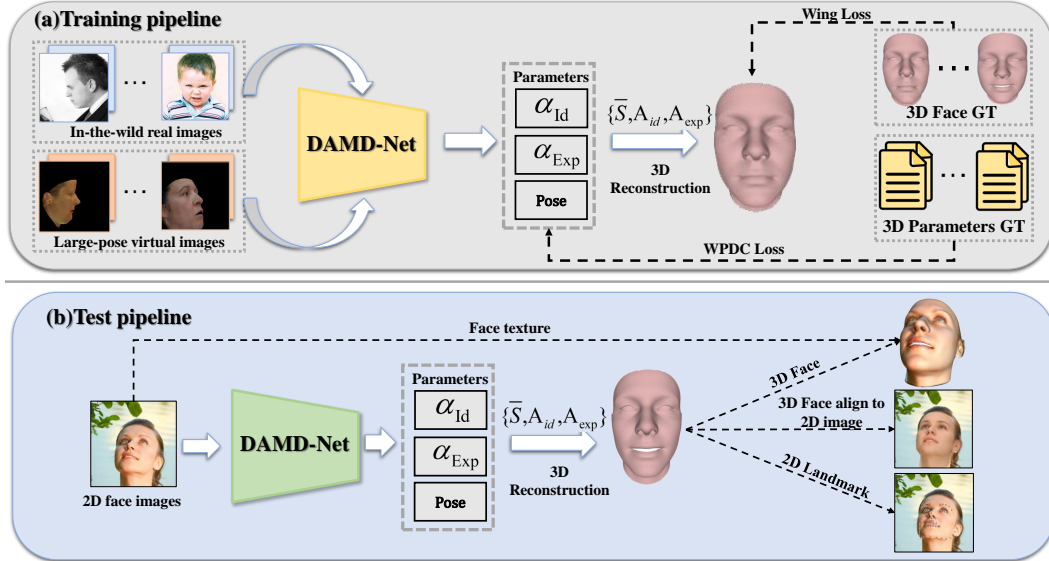


Figure 1. Overview of our method. As efficient dual attention convolutional neural network (DAMDNet). (a) Training pipeline for a single image 3D face Alignment, (b) Test pipeline. Figure 2 describes the details of DAMDNet.

parameter space is large, so the network is difficult to train to achieve convergence.

Inspired by the efficiency of MobileNet[20], achieved by the use of the Depthwise Separable Convolution in the network structure, and DenseNet[22] strengthened by the transmission of features, we propose a network structure that extends both the idea of Depthwise Separable Convolution, and the feature reuse of Densely Connected networks. As dense connection convolution may lead to channel information redundancy, this paper adds a lightweight Channel Attention Mechanism in the network structure, which improves the representation ability of the network without increasing the number of network parameters.

In convolutional neural networks, in addition to channel feature re-calibration, another important dimension that should be considered is the spatial dimension. For a specific semantic group, it is desirable and beneficial to identify the semantic features in the correct spatial location of the original image. Based on the channel attention mechanism, we enhance spatial features by grouping. Spatial Group-wise Enhancement[28] is feature re-calibration in spatial dimension. By combining channel and spatial attention mechanisms, the proposed network structure is a dual attention convolutional neural network.

In order to solve the problem of paucity of training samples in the case of large poses, this paper also presents a side-face data augmentation to enhance the robustness of the network model for arbitrary pose. Extensive experiments are conducted on AFLW dataset[26] with a wide range of poses, and the AFLW2000-3D dataset[45], in com-

parison with a number of methods. We also provide the means for subjective evaluation by visualizing the 2D/3D face alignment and face reconstruction on the DFW[27, 37] dataset.

An overview of our method is shown in Figure 1.

In summary, our contributions are as follows:

1) We propose a novel efficient network structure (DAMDNet). To the best of our knowledge, this is the first time that Depthwise Separable Convolution scheme, a Densely Connected network structure, a Channel Attention Mechanism and Spatial Group-wise Feature Enhancement are combined to create a DNN novel architecture.

2) Different loss functions are used to optimize the parameters of 3D Morphable Model and its 3D vertices. The resulting method can estimate 2D/3D landmarks of faces with an arbitrary pose.

3) The training data set is augmented by integrating various data enhancement techniques. The face profile technique and virtual sample technique are used to increase its number of the large pose face training data set.

4) We experimentally demonstrate that our algorithm has significantly improved the 3D alignment performance, compared to the state of the art methods. The proposed face alignment method can deal with arbitrary poses and it is more efficient.

2. Related Work

In this section, we review the prior work in generic face alignment and 3D face alignment.

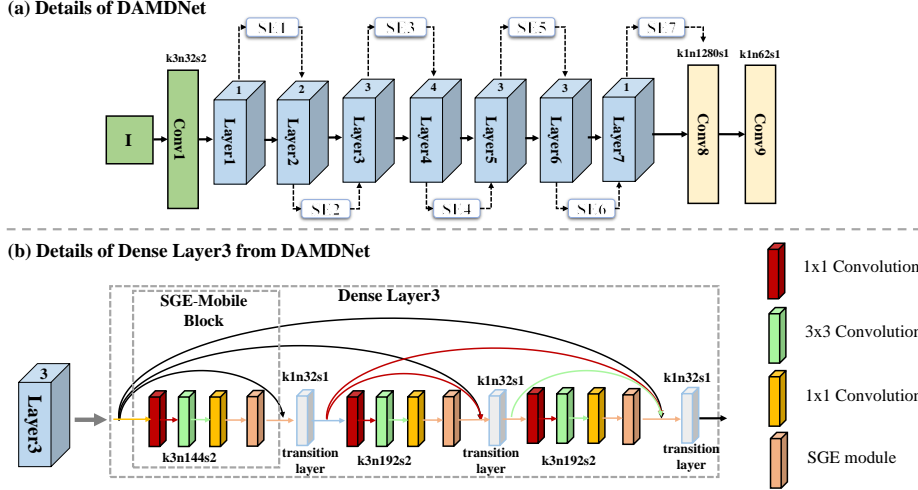


Figure 2. (a)Details of DAMDNet. $k3n64s1$ corresponds to the kernel size(k),number of feature maps(n) and stride(s) of conv1. (b)The details of one of the DenseBlock layers, namely *Layer3*. The convolution layer of a set of $1 \times 1, 3 \times 3, 1 \times 1$ filters and a SGE[28] module in DAMDNet as a basic unit called SGE-MobileBlock. The transition layer is the number of channels to match the input and output feature maps.

2.1. Generic Face Alignment

Face alignment research can boast many achievements, included the active appearance model(AAM)[12, 34] and the active shape model(ASM)[11].These methods consider face alignment as an optimization problem to find the best shape and appearance parameters, which allow the appearance model to achieve the best possible fit to the input face.The basic idea of the Constrained Local Model (CLM) method [14, 1, 36] in the Discriminative approaches category is to learn a set of local appearance models, one for each landmark.The output of the local models is combined with the help of a global shape model. Cascaded regression gradually refines initial predictions through a series of regressions. Each regression unit relies on the output of the previous regression unit to perform simple image operations. The entire system automatically learns from the training samples[15]. The ESR[10] (Explicit Shape Regression) proposed by Sun *et al.* includes three methods, namely two-level boosted regression, shape-indexed features and a correlation-based feature selection method.

Besides the traditional models, deep convolutional neural networks have recently been used for feature point localization of faces. Sun *et al.*[38] were first to use CNN to regress the raw face image landmark locations,accurately positioning 5 key points of the face from coarse to fine. The work of [19] uses the human body pose estimation, and the boundary information for the key point regression. In recent years, most of the landmark detection methods have been adopted some form of "coarse to fine" strategy. On the other hand, Feng *et al.*[17] have taken a different approach,

using the idea of cascaded convolutional neural networks. A [17] compared the commonly used loss functions for face landmark detection, and based on this, the concept of wing loss was proposed.

2.2. 3D Face Alignment

Although traditional methods provide a guide to successful face alignment, they are affected by non-frontal pose, illumination and occlusion in real-life applications. The most common approach to deal with pose variation is the multi-view framework [39], which uses different landmark configurations for different views. For example, TSPM [47] and CDM [44] use the DPM-like [18] method to align faces of different shape models, and finally select the most probable model as the final result. However, since each view requires testing, the computational cost of the multiview approach is always high.

Apart from multi-view solutions, 3D face alignment is also a popular approach. 3D face alignment [19, 23] aims to fit a 3D morphable model (3DMM) [6] to a 2D image. The 3D Morphable Model is a typical statistical 3D face model. It has a clear understanding of 3D faces based on a statistical analysis. Zhu *et al.*[45] proposed a localization method based on 3D face shape, which solves the problem of some feature points being invisible in extreme poses (such as side faces), as well as the face appearance in different poses varying greatly, making it difficult to locate landmarks. Liu *et al.*[24] used a cascade of 6 convolutional neural networks to solve the problem of locating facial feature points in images of faces with extreme poses by means of 3D face mod-

elling. This method not only predicts the 3D face shape and projection matrix, but also calculates whether each feature point is visible or not. If a feature point is invisible, the feature block about the invisible point is not used as input, which is difficult to achieve for common 2D face alignment methods. Paper [16] designed a UV position map to represent 3D shape features of a complete human face in a 2D. The purpose of 3D face alignment is to reconstruct the 3D face from a 2D image, and then align the 3D face to the 2D image, so that 2D/3D face feature points can be located. Our approach is also based on convolutional neural networks, but we have redesigned the network structure to make it efficient and robust. At the same time, we use different loss functions for 3D parameters and 3D vertices to constrain the semantic information being recovered.

3. Proposed Method

In this section we introduce our proposed robust 3D face alignment method, which fits a 3D morphable model using DAMDNet.

3.1. 3D Morphable Model

The 3D Morphable model is one of the most successful methods for describing 3D face space. Blanz *et al.* [6] proposed a 3D morphable model (3DMM) of 3D face based on Principal Component Analysis (PCA). It is expressed as follows:

$$S = \bar{S} + A_{id}\alpha_{id} + A_{exp}\alpha_{exp} \quad (1)$$

where S is a specific 3D face, \bar{S} is the mean face, A_{id} are the principle axes trained on the 3D face scans with neutral expression and α_{id} is the shape parameter vector, A_{exp} are the principle axes trained on the offsets between expression scans and neutral scans and α_{exp} is the expression parameter vector. The coefficients $\{\alpha_{id}, \alpha_{exp}\}$ define a unique 3D face. In this work A_{id} adopted come from the Basel Face Model (BFM)[31] and A_{exp} comes from the FaceWarehouse model[8].

In the process of 3DMM fitting, we use the Weak Perspective Projection to project 3DMM onto the 2D face plane. This process can be expressed as follows:

$$S_{2d} = f * Pr * R * \{S + t_{3d}\} \quad (2)$$

where S_{2d} is the 2D coordinate matrix of the 3D face after Weak Perspective Projection, rotation and translation. f is the scaling factor. Pr is the projection matrix $\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}$. R is a rotation matrix constructed according to three rotation angles of pitch, yaw and roll respectively. t_{3d} is the 3D translation vector. For the modeling of a specific face, we only need to find the 3D parameters $P = [f, pitch, yaw, roll, t_{3d}, \alpha_{id}, \alpha_{exp}]$

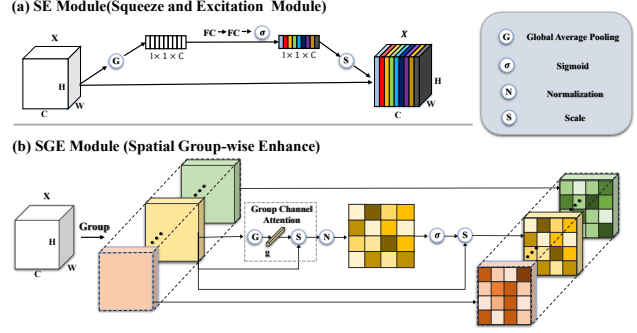


Figure 3. (a)SE Module(Squeeze and Excitation Module),(b)SGE Module(Spatial Group-wise Enhance).

3.2. Dual Attention Mechanism

Extracting the main facial features for the 3D face alignment task is a critical step. A 2D convolutional neural network typically performs feature extraction in the Channel and Spatial dimensions. This paper enhances the feature representation of convolutional neural networks by adding lightweight attention mechanisms in both, spatial and channel dimensions.

In the channel dimension, we opt for a SE[21] attention mechanism module. The SE module uses a new feature recalibration strategy. Specifically, it learns the importance of each feature channel automatically, and enhances the useful features according to the learnt importance measure, and suppresses the non informative features. Figure 3(a) describes the basic operation of the SE module. It first uses a global average pooling layer as a Squeeze operation. Then the two fully connected layers form a Bottleneck structure to model the correlation between the channels and output the same number of weights as the number of input features. A normalized weight between 0 and 1 is obtained by a Sigmoid function to weight each channel.

The spatial attention mechanism induces the model to pay more attention to the contribution of the key feature areas of the human face and reduces the influence of other unrelated features. We introduce the SGE(Spatial Group-wise Enhancement)[28] mechanism to strengthen the the spatial distribution of facial semantic features. A comprehensive face feature is composed of many sub features, and these sub features are distributed in groups in each feature layer. By generating an attention factor for each group, SGE module can gauge the importance of each sub feature, and help to suppress noise in a targeted way. This attention factor is determined by the similarity of global and local features within each group, so SGE is very lightweight. Figure 3(b) describes the specific computational operations of SGE. First, the features are grouped, and each set of features is spatially compared with the global pooling feature

(similarity) to get the initial attention mask. This part we call Group Channel Attention. After, normalizing the attention mask, we obtain the final attention mask through a sigmoid operation, and scale the features of each position to the original feature group.

3.3. DAMDNet(Dual Attention MobDenseNet) Structure

The DAMDNet proposed in this paper applies the depth separable convolution, dense connection, channel attention and spatial attention mechanism to the 3D face alignment task for the first time. The architecture of DAMDNet is illustrated in Figure 2(a). Conv1 is a convolution layer with kernel size(k) of 3, stride(s) of 2 and the number of feature maps(n) totalling 32 to extract rough features. *Layer1* to *Layer7* are 7 dense blocks for extracting deep features. An SE[21] module is added between each DenseBlock to explicitly model the interdependencies between feature channels. Figure 2(b) shows the details of one of the DenseBlock, *Layer3*. The convolution layer of a set of $1 \times 1, 3 \times 3, 1 \times 1$ filters and the SGE[28] module in DAMDNet form the basic unit called SGE-MobileBlock. DenseLayer3 contains three sets of SGE-MobileBlock(each SGE-MobileBlock output is cascaded as the input of the next SGE-MobileBlock). As shown in Figure 2(b), Layer3 contains three sets of SGE-MobileBlock. In order to match the number of channels connected to the Dense connection, we add a transition layer after each SGE-MobileBlock (the convolution layer filter is 1×1), the purpose is to adjust the number of channels in the preview SGE-MobileBlock output feature map.

3.4. Loss Function

We use two different Loss Functions to jointly train DAMDNet. For predicting 3D parameters we make use of the Weighted Parameter Distance Cost (WPDC) of Zhu *et al.* [45] to calculate the difference between the ground truth of 3D parameters and the predicted 3D parameters. The basic idea is explicitly to model the importance of each parameter:

$$L_{wpdc} = (P_{gt} - \bar{P})^T W (P_{gt} - \bar{P}) \quad (3)$$

where \bar{P} is an estimate and P_{gt} is the ground truth. The diagonal matrix W contains the weights. For each element of the shape parameter p , its weight is the inverse of the standard deviation that was obtained from the data used in 3DMM training. Our ultimate goal is to accurately obtain 68 landmarks of the human face, For 3D face vertices reconstructed with the estimated 3D parameters, we use Wing Loss[17] which is defined as:

$$L_{wing}(\Delta V(P)) = \begin{cases} \omega \ln(1 + |\Delta V(P)| / \epsilon) & \text{if } |\Delta V(P)| < \omega \\ |\Delta V(P)| - C & \text{otherwise} \end{cases} \quad (4)$$

where $\Delta V(P) = V(P_{gt}) - V(\bar{P})$, $V(P_{gt})$ and $V(\bar{P})$ are the ground truth of the 3D facial vertices and the 3D facial vertices reconstructed using the 3D parameters predicted by the network, respectively. ω and ϵ are the log function parameters. $C = \omega - \omega \ln(1 + \omega / \epsilon)$ is a constant that smoothly links the piecewise-defined linear and nonlinear parts.

Overall, the framework is optimized by the following loss function:

$$L_{loss} = \lambda_1 L_{wpdc} + \lambda_2 L_{wing} \quad (5)$$

where λ_1 and λ_2 are parameters, which balance the contribution of L_{wpdc} and L_{wing} . The selection of these parameters will be discussed in the next section.

3.5. Data Augmentation and Training

The input to DAMDNet is a 2D image with the facial ROI localized by a face detector. In this paper, we use the Dlib¹ SDK for face detection. We first enlarge the detected face bounding box by a factor of 0.25 of its original size and crop a square image patch of the face ROI, which is scaled to 120×120 . DAMDNet outputs a 62-dimensional 3D parameter vector, including 40-dimensional identity parameter vector, 10-dimensional expression parameter vector and 12-dimensional pose vector. We use both real face images and generated face images to train our DAMDNet. We use the same method as [32] to generate a virtual face sample with full parameters. The generated face samples contain a large number of large poses.

Most of the current real training data sets contain images of small and medium poses and unoccluded faces. In order to improve the robustness of the algorithm for arbitrary poses, we conduct a facial profile processing of real face images using the methods proposed by Zhu *et al.* [46]. The face standardization process divides the face image into three regions: the face region, the area around the face, and the background region. Similar to face normalization, the basic idea of face profile is to predict the depth of the face image and generate a contour view that can be rotated in three dimensions. When the depth information is estimated, the face image can be rotated in three dimensions to produce the appearance of larger poses. In this paper, we rotate each real sample by 10 to 90 degrees on the z-axis to generate a face image of new poses. Figure 4, (a) and (b) show the effect of a 2D face image rotated by $0^\circ, 15^\circ, 30^\circ$ and 60° respectively, and Figure 4(c) its 3D mesh

4. Experiments

In this section, we evaluate the performance of our method on three common face alignment tasks, face alignment in small and medium poses, face alignment in large

¹<http://dlib.net/>

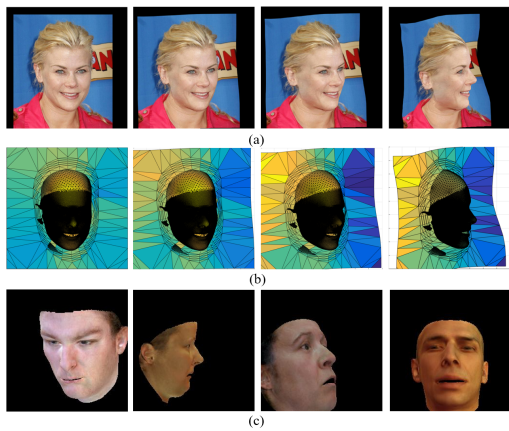


Figure 4. (a) and (b) from left to right are face images and 3D mesh diagrams rotated by 0° , 15° , 30° , and 60° around the z-axis; (c) virtual face samples.

poses, and face reconstruction in extreme poses ($\pm 90^\circ$ yaw angles), respectively.

4.1. Implementation details

We use the Pytorch² deep learning framework to train the DAMDNet models. The loss weights of our method are empirically set to $\lambda_1 = 0.5$ and $\lambda_2 = 1$. In our experiments, we set the parameters of the Wing loss as $\omega = 10$ and $\epsilon = 2$. The Adam solver[25] is employed with the mini-batch size and the initial learning rate set to 128 and 0.01, respectively. There are 680,000 face images in our training set, including 430,000 real face images and 250,000 synthetic face images. Real face images come from AFLW[47] and LFPW[2] data sets, and various data enhancement algorithms are adopted to expand the datasets. We run the training for a total of 40 epochs. After 15, 25 and 30 epochs, we reduced the learning rate to 0.002, 0.0004 and 0.00008 respectively.

4.2. Evaluation databases

We evaluate the performance of our method on three publicly available face data sets AFLW [26], AFLW2000-3D[45] and DFW[27, 37]. These AFLW and AFLW2000-3D data sets contain small and medium poses, large poses and extreme poses ($\pm 90^\circ$ yaw angles). We divide the dataset AFLW and AFLW2000-3D into three sections of $[0^\circ, 30^\circ]$, $[30^\circ, 60^\circ]$, and $[60^\circ, 90^\circ]$ according to the face absolute yaw angle.

AFLW AFLW face database is a large-scale face database including multi-poses and multi-views, and each face is annotated with 21 feature points. This database contains very diverse images, including pictures of various poses, expressions, lighting, and ethnicity. The AFLW face database consists of approximately 250 million hand-labeled face images, of which 59% are women and 41% are

men. Most of the images are color images, only a few are gray images. We only use the part of extreme pose face images of the AFLW database for qualitative analysis.

AFLW2000-3D AFLW2000-3D is constructed [45] to evaluate 3D face alignment on challenging unconstrained images. This database contains the first 2000 images from AFLW and expands its annotations with fitted 3DMM parameters and 68 3D landmarks. We use this database to evaluate the performance of our method for the face alignment task.

DFW Disguised Faces in the Wild (DFW)[27, 37] dataset containing 11,157 images pertaining to 1,000 identities with variations in terms of different disguise accessories. For a given subject there are four types of images: normal, validation, disguised, and impersonator. We visualized DAMDNet’s 3D face alignment effect on the DFW data set, proving that our algorithm also has excellent performance for disguised face.

4.3. The evaluation metric

We are given the ground truth 2D landmarks U_i , their visibility v_i , and estimated landmarks \hat{U}_i for N_t test images. Normalized Mean Error (NME) is the average of the normalized estimation error of visible landmarks, defined as,

$$NME = \frac{1}{N_t} \sum_i \left(\frac{1}{d_i |v_i|_1} \sum_j v_i(j) \|\hat{U}_i(:, j) - U_i(:, j)\| \right) \quad (6)$$

where d_i is the square root of the face bounding box size. Note that normally d_i is the distance of the two centers of the eyes in most prior face alignment work dealing with near-frontal face images.

4.4. Comparative evaluation

4.4.1 Comparison on AFLW

In the AFLW dataset, 21,080 images were selected as test samples, with 21 landmarks available for each sample. During testing, we divide the test set into 3 subsets according to their absolute yaw angles: $[0^\circ, 30^\circ]$, $[30^\circ, 60^\circ]$, and $[60^\circ, 90^\circ]$ with 11,596, 5,457 and 4,027 samples respectively. Since a few experiments have been conducted on AFLW, we choose baseline methods for which the code is available, including CDM [44], RCPR [7], ESR [10], SDM [43], 3DDFA[45] and nonlinear 3DMM[40]. Table 1 presents the results, given in terms of NME(%) of face alignment on AFLW with the best results highlighted. The results of the provided alignment models are identified with their references. Figure 5 shows the corresponding CED curves. Our CED curve is only compared to the best method in Table 1. Since the best nonlinear 3DMM method currently only provides data for the

²<https://pytorch.org/>

Table 1. The NME(%) of face alignment results on AFLW and AFLW2000-3D.

Method	AFLW DataSet(21 pts)					AFLW2000-3D DataSet(68 pts)				
	[0°, 30°]	[30°, 60°]	[60°, 90°]	Mean	Std	[0°, 30°]	[30°, 60°]	[60°, 90°]	Mean	Std
CDM[44]	8.150	13.020	16.170	12.440	4.040	-	-	-	-	-
RCPR[7]	5.430	6.580	11.530	7.850	3.240	4.260	5.960	13.180	7.800	4.740
ESR[10]	5.660	7.120	11.940	8.240	3.290	4.600	6.700	12.670	7.990	4.190
SDM[43]	4.750	5.550	9.340	6.550	2.450	3.670	4.940	9.760	6.120	3.210
3DDFA(CVPR16)[45]	5.000	5.060	6.740	5.600	0.990	3.780	4.540	7.930	5.420	2.210
Nonlinear 3DMM(CVPR18)[40]	-	-	-	-	-	-	-	-	4.700	-
Ours	4.359	5.209	6.028	5.199	0.682	2.907	3.830	4.953	3.897	0.837

Table 2. The NME(%) of face alignment results on AFLW and AFLW2000-3D with the different network structures.

Method	GFLOPs	Params(M)	AFLW DataSet(21 pts)					AFLW2000-3D DataSet(68 pts)				
			[0°, 30°]	[30°, 60°]	[60°, 90°]	Mean	Std	[0°, 30°]	[30°, 60°]	[60°, 90°]	Mean	Std
RestNeXt50	1.319	23.11	4.599	5.516	6.297	5.471	0.694	3.122	4.065	5.351	4.179	0.913
MobileNetV2	0.109	2.38	4.643	5.581	6.397	5.540	0.716	3.236	4.080	5.181	4.165	0.796
DenseNet121	0.800	7.02	4.442	5.249	6.168	5.286	0.705	3.051	3.912	5.297	4.087	0.925
MDNet	0.127	2.74	4.549	5.427	6.204	5.393	0.676	3.149	4.010	5.270	4.143	0.871
AMDNet	0.128	2.75	4.367	5.317	6.131	5.271	0.726	2.879	3.906	4.982	3.922	0.858
DAMNet	0.125	2.76	4.359	5.209	6.028	5.199	0.682	2.907	3.830	4.953	3.897	0.837

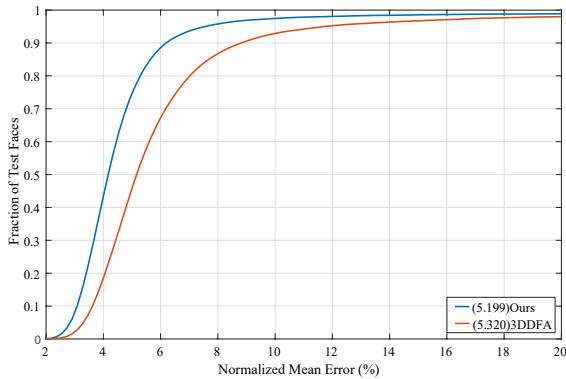


Figure 5. Cumulative errors distribution (CED) curves on AFLW.

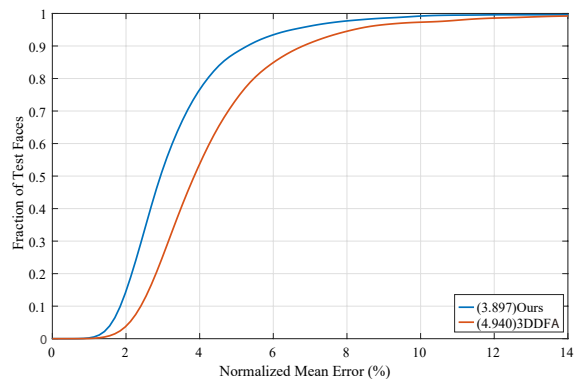


Figure 6. Comparisons of cumulative errors distribution (CED) curves on AFLW2000-3D.

AFLW2000-3D dataset, there is no CED baseline for it. The results show that our algorithm significantly improves the face alignment accuracy in a full range of poses. The minimum standard deviation of our method also proves its robustness to poses changes.

4.4.2 Comparison on AFLW2000-3D

In the case of the AFLW2000-3D dataset, 2000 images were selected as test samples. Considering the visible and invisible evaluation, the 3D face alignment evaluation can be transformed to a full landmark evaluation. We divide the test set into 3 subsets according to their absolute yaw angles: [0°, 30°], [30°, 60°], and [60°, 90°] with 1,312, 383 and 305 samples respectively. Table 1 presents the results (NME(%)) with the best results highlighted. The results

achieved by existing methods are identified by their references. Figure 6 shows the corresponding CED curves. Table 1 and Figure 6 demonstrate that our algorithm also achieves a significant improvement in the prediction of invisible regions, showing good robustness for face alignment in arbitrary poses.

4.4.3 A visualization experiment performed on DFW

In the DFW database, we select some face images for qualitative testing. DFW is by far the most complete data set of disguised faces in the wild. Figure 7 visualizes the results of our method on DFW, showing (a) 2D landmarks, (b) the fitted 3D face model with image face texture, (c) the reconstructed 3D face and (d) the result of 3D reconstruction



Figure 7. (a)the landmarks of 2D, (b) the 3D face model with image face texture, (c) the reconstructed 3D face, and (d) the result of 3d reconstruction of the mean texture of the model using z-buffer projection on the input image.

of the mean texture of the model using Z-buffer projection on the input image. Accurate 3D face alignment plays an important role in the next step of disguised face recognition. The results show that our algorithm is robust to disguise. Our algorithm can accurately locate the key points of a face and provide 3D face structure and depth information. This effectively improves the recognition accuracy of the disguised face.

4.4.4 A comparison of different network structures

In order to verify the effectiveness of our network structure, we compare our method and the current mainstream neural network structures on the task of face alignment. The experimental network structures include ResNeXt[42], MobileNetV2[33], DenseNet121[22], and our proposed DAMDNet. To the best of our knowledge, these three popular and efficient network structures are the first applied to the task 3D face alignment. Table 2 shows that our DAMDNet achieves a 5% and 6.7% reduction in error on the AFLW and AFLW2000 datasets compared to ResNeXt50. In terms of operational efficiency, the GFLOPs complexity is reduced 10.5 times and the number of model parameters is reduced 8.37 times. Compared with DenseNet121, DAMDNet has reduced the error on the AFLW and AFLW2000 data sets by 2.2% and 4.6%, the GFLOPs complexity 6.4 times, and the number of model parameters 2.5 times. Compared with MobileNetV2, DAMDNet is higher both in terms of GFLOPs and the number of network parameters due to the addition of Densely Connected Convolutional and Dual Attention Mechanism in our network structure. However, our model has obvious

advantages in terms of accuracy.

Similarly, in order to verify the validity of each module of our proposed network structure, we compare MDNet, AMDNet and DAMDNet respectively. Among them, MDNet only combines Depthwise Separable Convolution and Densely Connected structure, AMDNet adds an SE module of Channel Attention Mechanism, and DAMDNet includes Depthwise Separable Convolution, Densely Connected structure and the Dual Attention Mechanism. AMDNet adds a channel attention mechanism based on MDNet, which reduces face alignment error by 2.2% in the case of the AFLW dataset and 5.1% for the AFLW2000-3D dataset. However, the number of GFLOPs and network parameters is increased by 0.78% and 0.36% respectively. This result shows that SE module can improve the precision of the network significantly, without adding network parameters and GFLOPs. DAMDNet adds a Spatial Group-wise Feature Enhancement based on AMDNet, which reduces face alignment error by 1.4% in the AFLW dataset and 0.64% in the AFLW2000-3D dataset. The GFLOPs is reduced by 2.3% and the number of parameters is increased by 0.36%. DAMDNet, which add the SGE module and the Channel Attention Mechanism, significantly improves model efficiency and face alignment accuracy.

5. Conclusions

In this paper, we proposed a DAMDNet which solves the problem of 2D/3D face alignment for face images exhibiting a full range of poses. In order to improve the feature expression ability of the model, we incorporated a lightweight attention mechanism for channel and spatial dimensions respectively. We also proposed two novel loss functions to jointly optimize 3D reconstruction parameters and 3D vertices. We use a variety of data augmentation methods and generate a large virtual pose face data set to solve the problem of the large pose training sample imbalance. Our method achieved the best accuracy on both AFLW, AFLW2000-3D datasets, compared to existing methods. A qualitative evaluation of the proposed method on the DFW dataset showed its promise in dealing with disguised faces. In comparison to several popular networks, our algorithm achieves a good trade-off between accuracy and efficiency. In the future, will explore the texture and illumination features of face images.

Acknowledgments

The paper is supported by the National Natural Science Foundation of China (Grant No.61672265,U1836218),the 111 Project of Ministry of Education of China(Grant No.B12018) and UK EPSRC Grant EP/N007743/1, Muri/EPSRC/Dstl Grant EP/R018456/1.

References

- [1] A. Asthana, S. Zafeiriou, S. Cheng, and M. Pantic. Robust discriminative response map fitting with constrained local models. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3444–3451, 2013.
- [2] P. N. Belhumeur, D. W. Jacobs, D. J. Kriegman, and N. Kumar. Localizing parts of faces using a consensus of exemplars. *IEEE transactions on pattern analysis and machine intelligence*, 35(12):2930–2940, 2013.
- [3] V. Bettadapura. Face expression recognition and analysis: the state of the art. *arXiv preprint arXiv:1203.6722*, 2012.
- [4] V. Blanz, A. Mehler, T. Vetter, and H.-P. Seidel. A statistical method for robust 3d surface reconstruction from sparse data. In *3D Data Processing, Visualization and Transmission, 2004. 3DPVT 2004. Proceedings. 2nd International Symposium on*, pages 293–300. IEEE, 2004.
- [5] V. Blanz and T. Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, pages 187–194. ACM Press/Addison-Wesley Publishing Co., 1999.
- [6] V. Blanz and T. Vetter. Face recognition based on fitting a 3d morphable model. *IEEE Transactions on pattern analysis and machine intelligence*, 25(9):1063–1074, 2003.
- [7] X. P. Burgos-Artizzu, P. Perona, and P. Dollár. Robust face landmark estimation under occlusion. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1513–1520, 2013.
- [8] C. Cao, Y. Weng, S. Zhou, Y. Tong, and K. Zhou. Facewarehouse: A 3d facial expression database for visual computing. *IEEE Transactions on Visualization and Computer Graphics*, 20(3):413–425, 2014.
- [9] C. Cao, H. Wu, Y. Weng, T. Shao, and K. Zhou. Real-time facial animation with image-based dynamic avatars. *ACM Transactions on Graphics*, 35(4), 2016.
- [10] X. Cao, Y. Wei, F. Wen, and J. Sun. Face alignment by explicit shape regression. *International Journal of Computer Vision*, 107(2):177–190, 2014.
- [11] T. Cootes, E. Baldock, and J. Graham. An introduction to active shape models. *Image processing and analysis*, pages 223–248, 2000.
- [12] T. F. Cootes, G. J. Edwards, and C. J. Taylor. Active appearance models. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (6):681–685, 2001.
- [13] T. F. Cootes, C. J. Taylor, and A. Lanitis. Active shape models: Evaluation of a multi-resolution method for improving image search. In *BMVC*, volume 1, pages 327–336. Citeseer, 1994.
- [14] D. Cristinacce and T. F. Cootes. Feature detection and tracking with constrained local models. In *Bmvc*, volume 1, page 3. Citeseer, 2006.
- [15] P. Dollár, P. Welinder, and P. Perona. Cascaded pose regression. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 1078–1085. IEEE, 2010.
- [16] Y. Feng, F. Wu, X. Shao, Y. Wang, and X. Zhou. Joint 3d face reconstruction and dense alignment with position map regression network. *arXiv preprint arXiv:1803.07835*, 2018.
- [17] Z.-H. Feng, J. Kittler, M. Awais, P. Huber, and X.-J. Wu. Wing loss for robust facial landmark localisation with convolutional neural networks. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2235–2245. IEEE, 2018.
- [18] D. Forsyth. Object detection with discriminatively trained part-based models. *Computer*, (2):6–7, 2014.
- [19] L. Gu and T. Kanade. 3d alignment of face in a single image. In *null*, pages 1305–1312. IEEE, 2006.
- [20] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [21] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018.
- [22] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *CVPR*, volume 1, page 3, 2017.
- [23] A. Jourabloo and X. Liu. Pose-invariant 3d face alignment. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3694–3702, 2015.
- [24] A. Jourabloo and X. Liu. Large-pose face alignment via cnn-based dense 3d model fitting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4188–4196, 2016.
- [25] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [26] M. Koestinger, P. Wohlhart, P. M. Roth, and H. Bischof. Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization. In *Computer Vision Workshops (ICCV Workshops), 2011 IEEE International Conference on*, pages 2144–2151. IEEE, 2011.
- [27] V. Kushwaha, M. Singh, R. Singh, M. Vatsa, N. Ratha, and R. Chellappa. Disguised faces in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 1–9, 2018.
- [28] X. Li, X. Hu, and J. Yang. Spatial group-wise enhance: Enhancing semantic feature learning in convolutional networks. 2019.
- [29] X. Liu. Discriminative face alignment. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(11):1941, 2009.
- [30] I. Matthews and S. Baker. Active appearance models revisited. *International journal of computer vision*, 60(2):135–164, 2004.
- [31] P. Paysan, R. Knothe, B. Amberg, S. Romdhani, and T. Vetter. A 3d face model for pose and illumination invariant face recognition. In *Advanced video and signal based surveillance, 2009. AVSS'09. Sixth IEEE International Conference on*, pages 296–301. Ieee, 2009.
- [32] E. Richardson, M. Sela, and R. Kimmel. 3d face reconstruction by learning from synthetic data. In *3D Vision (3DV), 2016 Fourth International Conference on*, pages 460–469. IEEE, 2016.
- [33] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen. Inverted residuals and linear bottlenecks: Mobile net-

- works for classification, detection and segmentation. *arXiv preprint arXiv:1801.04381*, 2018.
- [34] J. Saragih and R. Goecke. A nonlinear discriminative approach to aam fitting. In *Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on*, pages 1–8. IEEE, 2007.
 - [35] J. M. Saragih, S. Lucey, and J. F. Cohn. Face alignment through subspace constrained mean-shifts. In *Computer Vision, 2009 IEEE 12th International Conference on*, pages 1034–1041. Ieee, 2009.
 - [36] J. M. Saragih, S. Lucey, and J. F. Cohn. Deformable model fitting by regularized landmark mean-shift. *International Journal of Computer Vision*, 91(2):200–215, 2011.
 - [37] R. Singh, M. Vatsa, and A. Noore. Recognizing face images with disguise variations. In *Recent Advances in Face Recognition*. IntechOpen, 2008.
 - [38] Y. Sun, X. Wang, and X. Tang. Deep convolutional network cascade for facial point detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3476–3483, 2013.
 - [39] A. T. Tran, T. Hassner, I. Masi, and G. Medioni. Regressing robust and discriminative 3d morphable models with a very deep neural network. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 1493–1502. IEEE, 2017.
 - [40] L. Tran and X. Liu. Nonlinear 3d face morphable model. *arXiv preprint arXiv:1804.03786*, 2018.
 - [41] M. Valstar, B. Martinez, X. Binefa, and M. Pantic. Facial point detection using boosted regression and graph models. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 2729–2736. IEEE, 2010.
 - [42] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 5987–5995. IEEE, 2017.
 - [43] J. Yan, Z. Lei, D. Yi, and S. Li. Learn to combine multiple hypotheses for accurate face alignment. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 392–396, 2013.
 - [44] X. Yu, J. Huang, S. Zhang, and D. N. Metaxas. Face landmark fitting via optimized part mixtures and cascaded deformable model. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (11):2212–2226, 2016.
 - [45] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li. Face alignment across large poses: A 3d solution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 146–155, 2016.
 - [46] X. Zhu, Z. Lei, J. Yan, D. Yi, and S. Z. Li. High-fidelity pose and expression normalization for face recognition in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 787–796, 2015.
 - [47] X. Zhu and D. Ramanan. Face detection, pose estimation, and landmark localization in the wild. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 2879–2886. IEEE, 2012.