

# Deep Appearance Models: A Deep Boltzmann Machine Approach for Face Modeling

Chi Nhan Duong · Khoa Luu · Kha Gia Quach · Tien D. Bui

Received: date / Accepted: date

**Abstract** The “interpretation through synthesis” approach to analyze face images, particularly Active Appearance Models (AAMs) method, has become one of the most successful face modeling approaches over the last two decades. AAM models have ability to represent face images through synthesis using a controllable parameterized Principal Component Analysis (PCA) model. However, the accuracy and robustness of the synthesized faces of AAM are highly depended on the training sets and inherently on the generalizability of PCA subspaces. This paper presents a novel Deep Appearance Models (DAMs) approach, an efficient replacement for AAMs, to accurately capture both shape and texture of face images under large variations. In this approach, three crucial components represented in hierarchical layers are modeled using the Deep Boltzmann Machines (DBM) to robustly capture the variations of facial shapes and appearances. DAMs are therefore superior to AAMs in inferencing a representation for new face images under various challenging conditions. The proposed approach is evaluated in various applications to demonstrate its robustness and capabilities, i.e. facial super-resolution reconstruction, facial off-angle reconstruction or face frontalization, facial occlusion removal and age estimation using challenging face databases, i.e. Labeled Face Parts in the Wild (LFPW), Helen and FG-NET. Comparing to AAMs and other deep learning based approaches, the proposed DAMs achieve competitive results in those applications, thus this showed their advantages in handling occlusions, facial representation, and reconstruction.

Chi Nhan Duong<sup>1,2</sup>, Khoa Luu<sup>2</sup>, Kha Gia Quach<sup>1,2</sup>, Tien D. Bui<sup>1</sup>  
<sup>1</sup> Concordia University, Computer Science and Software Engineering, Montréal, Québec, Canada.

<sup>2</sup> CyLab Biometrics Center and the Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA, USA.  
E-mail: {c.duong, k.q, bui}@encs.concordia.ca; kluu@andrew.cmu.edu



**Fig. 1** An illustration in facial interpretation using the AAMs and our DAMs approach in real world images, e.g. low resolution, blurred faces, occlusions, off-angle faces, etc. The first row: original images; The second row: shape free images; The third row: facial interpretation using PCA-based AAMs; The fourth row: facial interpretation using our proposed DAMs approach.

**Keywords** Deep Boltzmann Machines, Deep Appearance Models, Active Appearance Models, Principal Component Analysis, Facial Super-resolution Reconstruction, Facial Off-angle Reconstruction, Face Frontalization, Age Estimation.

## 1 Introduction

Modeling faces with large variations has been a challenging task in computer vision. These variations such as expressions, poses and occlusions are usually complex and non-linear. Moreover, new facial images also come with their own characteristic artifacts greatly diverse. Therefore, a good face modeling approach needs to be carefully designed for flexibly adapting to these challenging issues. Over the last two decades, the “interpretation through synthesis” approach has become one of the most successful and popular face

modeling approaches. This approach aims to “describe” a given face image by generating a new synthesized image similar to it as much as possible. This purpose can be achieved by an optimization process on the appearance parameters of the model based *a priori* on constrained solutions. The subspace model then plays a key role that decides the robustness of the whole system. Therefore, in order to be applicable, it must provide a basis for a broad range of variations that are usually unseen.

Active Appearance Models (AAMs), one of the most successful face interpretation methods, were first introduced by Cootes et al (1998). Since then, it has been widely applied in many applications such as face recognition (Edwards et al, 1998), facial expression recognition (Sung and Kim, 2008), face tracking (Zhu et al, 2006), expressive visual text-to-speech (Anderson et al, 2013) and many other tasks. Although the framework of AAMs is general and effective, their generalization ability is still limited especially when dealing with unseen variations. Gross et al (2005) showed that AAMs perform well in person-specific cases rather than generic ones. Cootes and Taylor (2006) pointed out the problem of the pre-computed Jacobian matrix computed during the training step. Since it is only an approximation for testing image, it may lead to poor convergence when the image is very different from training data. Lighting changes (Pizarro et al, 2008) also make AAMs difficult to synthesize new images.

To overcome these disadvantages, there have been numerous improvements and adaptations based on the original approach (Matthews and Baker, 2004; Donner et al, 2006; Amberg et al, 2009; Joan Alabort-i Medina, 2014). However, even when these adaptations are taken into account, the capabilities of facial generalization and reconstruction are still highly dependent on the characteristics of training databases. This is because at the heart of AAMs, Principal Component Analysis (PCA) is used to provide a subspace to model variations in training data. The limitation of PCA to generalize to illumination and poses, particularly for faces, is very well known. Therefore, it is not surprising that AAMs have difficulties in generalizing to new faces under these challenging conditions. On the other hand, the variations in data are not only large but also non-linear. For example, the variations in different facial expressions or poses are non-linear. It apparently violates the linear assumptions of PCA-based models. Thus, single PCA model is unable to interpret the facial variations well. Figure 1 presents some faces with various challenging factors, i.e. low-resolution, blurred faces, occlusions, poses. The AAM interpretations presented in the third row of the figure have a major negative impact from these wide range of variations.

Recently, Deep Boltzmann Machines (DBM) (Salakhutdinov and Hinton, 2009) have gained significant attention as one of the emerging research topics in both the higher-level

representation of data and the distribution of observations. In DBM, non-linear latent variables are organized in multiple connected layers in a way that variables in one layer can simultaneously contribute to the probabilities or states of variables in the next layers. Each layer learns a different factor to represent the variations in a given data. Thanks to the nonlinear structure of DBM and the strength of latent variables organized in hidden layers, it efficiently captures variations and structures in complex data that could be higher than second order.

Moreover, DBM is shown to be more robust with ambiguous input data (Salakhutdinov and Hinton, 2009). There are some recent works using DBM as shape prior model (Eslami et al, 2014; Wu et al, 2013; Taylor et al, 2010). Far apart from these methods, the higher-level relationships of both shape and texture are exploited in our proposed DAMs so that the reconstruction of one can benefit from the information on the other. This paper proposes a novel Deep Appearance Models (DAMs<sup>1</sup>) approach to find a set of parameters in both shape and texture to characterize the identity, facial poses, facial expressions, lighting conditions of a given face. In addition, our proposed approach also has ability to generate a compact set of parameters in a robust model that can later be used for classification. Specifically, the DBM-based shape and texture models are first independently constructed. Then the interactions between these shapes and textures are further modeled using a deeper hidden layer. By this way, after fitting the model to new images, these interactions can be used as a compact set of parameters that represent both shape and appearance of faces for further discriminative problems. Furthermore, unlike other deep learning based approaches such as CNNs, with a specific topology of a stochastic neural network and sampling based weight update process (i.e. Contrastive Divergence), the need for large-scale training data is also alleviated in our DAMs structure.

A preliminary version of our work can be found in (Duong et al, 2015). In that work, our proposed DAMs<sup>2</sup> are able to capture a wide range of face variations as well as efficiently interpret the connections between face shape and texture. In this paper, we further extend the fitting stage to make it more robust to occlusions and noise. As a result, the model represented in this work is more advanced in terms of both face representation and model fitting. The paper is organized as follows. Section 2 reviews some recent AAMs-based and deep learning based approaches in both face representation and modeling. Section 3 presents the structure of DAMs and

<sup>1</sup> Noted that the term DAM is also used for “Direct Appearance Models” in (Hou et al, 2001).

<sup>2</sup> The implementation of DAMs will be available at <https://github.com/dchan/DeepAppearanceModels> and our project page <http://www.contrib.andrew.cmu.edu/~kluu/faceaging.html>

their properties. The model fitting is given in Section 4 followed by the experimental results in Section 5. Finally, the conclusion together with future work are given in Section 6.

## 2 Related Work

This section briefly reviews recent advances of AAMs-based approaches for constructing and fitting deformable models, and deep learning based methods for modeling human faces.

### 2.1 AAMs Modeling and Fitting

The basic AAMs (Cootes et al, 2001) build a statistical unified appearance model describing both shape and texture variation. One of the major drawbacks of AAMs is that the models only capture small amounts of appearance variations which can lead to poor performance on unknown variations caused by changes in the real-world environment, e.g. poses, lighting and camera conditions. The second drawback is that *person specific* AAMs substantially outperform generic AAMs trained across numerous subjects. Addressing the first drawback, some improvements have been made by applying the ideas of mixture models (Van Der Maaten and Hendriks, 2010) and probabilistic PCAs (Joan Alabort-i Medina, 2014) to represent as much variations as possible especially in the appearance model. Descriptive feature-based approaches were employed instead of intensity-based to deal with the second drawback of AAMs. Ge et al (2013) proposed three Gabor-based texture representations for AAMs capturing the properties of both Gabor magnitude and phase. These Gabor-based texture representations are more compact and robust to various conditions, e.g. expression, illumination and pose changes. Antonakos et al (2014) proposed to use dense histogram of oriented gradients features with AAMs to enhance their robustness and performance on unseen faces. Haase et al (2014) proposed a transfer learning based approach which incorporates related knowledge obtained from another training set with unseen illumination conditions to the existing AAMs to improve their generalization ability. Fitting steps in AAMs are an iterative optimization process measuring the cost between a testing image and a model texture in the coordinate of a reference frame. Generally, previous fitting techniques can be divided into two categories, i.e. discriminative and generative approaches. In the first category, the optimizing process is updated using a train-ed parameter-updating model. The model can be trained in several ways, e.g. perturbing the parameters and recording the residuals (Cootes et al, 1998), directly using texture information to predict the shape (Hou et al, 2001), linear regression (Donner et al, 2006) and non-linear regression methods (Saragih and Goecke, 2007). These techniques usually require low computational costs but their qual-

ity is still limited since the mapping function is fixed and independent of current model parameters. In the second category, the fitting steps are formulated as an image alignment problem and iteratively solved via the Gaussian-Newton optimization technique. Matthews and Baker (2004) presented a project out inverse algorithm to work on the orthogonal complement of the texture subspace. Navarathna et al (2011) proposed a computationally efficient fitting algorithm based on a variant of the Lucas-Kanade (LK) algorithm, called Fourier LK or Fourier AAM, to provide invariance to both expression and illumination. Other methods find the shape and texture increments either simultaneously (Gross et al, 2005) or alternatively (Papandreou and Maragos, 2008). Amberg et al (2009) presented the compositional framework. Tzimiropoulos and Pantic (2013) presented a fitting algorithm that works effectively in both forward and inverse cases. Mollahosseini and Mahoor (2013) proposed bidirectional warping method based on image alignment for AAMs fitting.

### 2.2 Deep Learning based Approaches

There has been significant recent interest in deep learning, e.g. deep convolutional networks and stacked auto-encoders, for face modeling or face representation, and face alignment. Zhu et al (2013) proposed to learn the face identity-preserving (FIP) features to represent faces while reducing significantly pose and illumination variances and preserving discriminative features for face recognition task. The authors designed a deep network that contains feature extraction layers, which produce FIP features, and a reconstruction layer, which reconstructs the canonical faces, i.e. frontal faces with neutral expression and normal illumination. Later, Zhu et al (2014) proposed a deep neural network so called multi-view perceptron (MVP) designed to separate the identity and view features. MVP can also generate multi-view images under unobserved viewpoints from a single 2D face image. Huang et al (2012) constructed local convolutional restricted Boltzmann machines that could exploit the global structure while achieving scalability and robustness to small misalignments. Sun et al (2014) proposed to learn a set of high-level features so called Deep hidden IDentity features for face verification. They trained 60 deep convolutional neural networks (CNNs) to extract complementary and over-complete representations from the last hidden layers of the networks. Taigman et al (2014) aim at improving the alignment and the representation step in face verification pipeline by using 3D model-based alignment and a nine-layer CNN. Both unsupervised similarity, i.e. the inner vector product, and supervised metric, i.e. the  $\chi^2$  similarity and the Siamese network are employed.

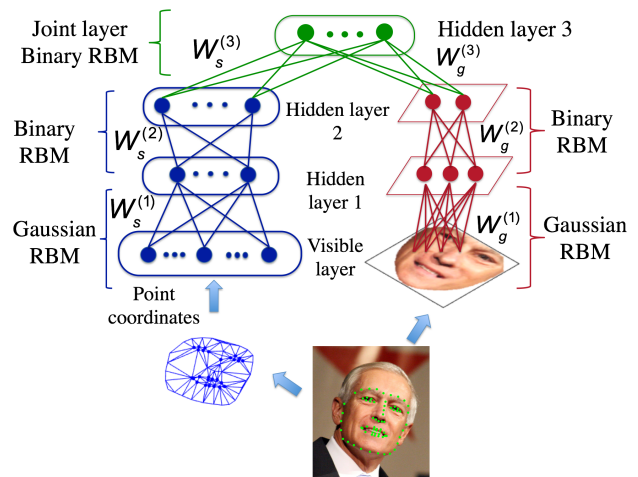
Kan et al (2014) proposed a deep network called stacked progressive auto-encoders to learn pose-robust features for

face recognition by modeling the complex non-linear transformation from the non-frontal face images to frontal ones. Similarly, Gao et al (2015) proposed to learn robust face representation features using a deep architecture based on supervised auto-encoder. The learning aims at transforming the faces with variants to the canonical view, and extracting similar features from the same subject. Ding and Tao (2015) proposed a deep learning framework to jointly learn face representation using multimodal information. The model consists of a set of CNNs to extract complementary facial features, and a three-layer stacked auto-encoder to compress the extracted features. Besides learning identity features for face recognition or verification tasks, age-related features could also be extracted (Zhai et al, 2015; Liu et al, 2016). Locating facial key points is also an essential step to represent facial shape. Sun et al (2013) proposed three-level cascaded CNNs for coarse-to-fine facial point detection (only detecting left eye, right eye, nose, and two mouth corners).

In addition, dealing with noise and occlusions, Robust Boltzmann machines (RoBMs) (Tang et al, 2012b) were proposed to extend RBM’s ability of estimating noise and distinguishing corrupted and uncorrupted pixels to find the optimal latent representations. The structure of RoBM consists of three components: a Gaussian RBM to model the “clean” data, a binary RBM for noise modeling, and a multiplicative gating mechanism to separate the clean data from noise/occlusion. Similar structure was used in (Tang et al, 2012a) to combine RBM with Lambertian reflectance model including the albedo and surface normals modeling instead of the occlusion/noise modeling. Li et al (2014) presented a single-face image decomposition method for image editing operations like relighting and re-texturing. They improved decomposition for faces by using human face priors including skin reflectance model and facial geometry. Yildirim et al (2015) proposed to combine a generative model based on 3D computer graphics and a discriminative model based on a CNN. This model can reconstruct the approximate shape and texture of a novel face from a single view.

### 3 Deep Appearance Models (DAMs)

The structure of DAMs consists of three main parts, i.e. two prior models for shape and texture and an additional higher-level hidden layer for appearance modeling. The shape model is used to learn the facial shape structure while texture model is used for texture variations. Both of them are mathematically modeled using the Deep Boltzmann Machines that are capable to model high-order correlations among input data. Their undirected connections provide both bottom-up and top-down passes to efficiently send updates between the texture model and the shape model. These modeling shape and texture parameters are then embedded in a higher-level layer



**Fig. 2** Deep Appearance Models that consists of shape model (left), texture model (right) and the joint representation of shape and texture.

that can be learned by clamping both shapes and textures as observations for the model. In this section, we present three main steps to construct the model, i.e. shape, texture and appearance modeling. Then in the next section, a fitting algorithm will be presented in order to synthesize any given new face image.

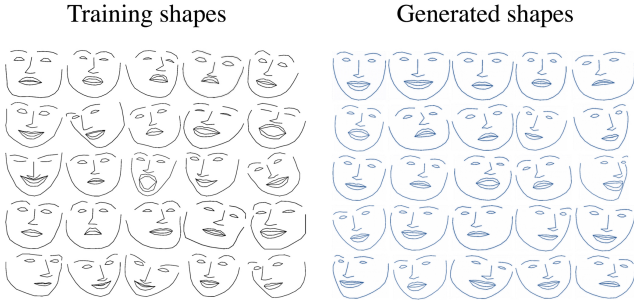
#### 3.1 Shape Modeling

In order to generalize possible patterns of facial shapes, we employ a two-layer DBM to learn the distributions of their landmark points. As illustrated in Figure 2, the shape model, i.e. left part of DAMs, consists of a set of visible units encoding the coordinates of landmark points and two sets of hidden units that are latent variables. The connections are symmetric and only those connecting units in adjacent layers are employed.

Let a shape  $\mathbf{s} = [x_1, y_1, \dots, x_N, y_N]^T$  with  $N$  landmark points  $\{x_i, y_i\}$ ,  $x_i \in \mathbb{R}, y_i \in \mathbb{R}$  be the visible units; and  $\mathbf{h}_s^{(1)} \in \{0, 1\}^{F_s^1}, \mathbf{h}_s^{(2)} \in \{0, 1\}^{F_s^2}$  be the binary variables of the first and second hidden layers respectively.  $F_s^1$  and  $F_s^2$  stand for the number of units in these hidden layers. Since  $\{x_i, y_i\}$  are real values while  $\mathbf{h}_s^{(1)}$  and  $\mathbf{h}_s^{(2)}$  are binary, we employ the Gaussian-Bernoulli Restricted Boltzmann Machines (GRBM) for the first layer and binary-binary RBM for the subsequent one. The energy of the joint configuration  $\{\mathbf{s}, \mathbf{h}_s^{(1)}, \mathbf{h}_s^{(2)}\}$  in facial shape modeling is formulated as follows:

$$E(\mathbf{s}, \mathbf{h}_s^{(1)}, \mathbf{h}_s^{(2)}; \theta_s) = \sum_i \frac{(s_i - b_{s_i})^2}{2\sigma_{s_i}^2} - \sum_{i,j} \frac{s_i}{\sigma_{s_i}} W_{sij}^{(1)} h_{sj}^{(1)} - \sum_{j,l} h_{sj}^{(1)} W_{sjl}^{(2)} h_{sl}^{(2)} \quad (1)$$

where  $\theta_s = \{\mathbf{W}_s^{(1)}, \mathbf{W}_s^{(2)}, \sigma_s^2, \mathbf{b}_s\}$  are the model parameters representing connecting weights of visible-to-hidden



**Fig. 3** A subset of training shapes and generated shapes from shape model with 10-step Gibbs sampling.

and hidden-to-hidden interactions, the variance, and the bias of visible units.

Notice that in Eqn. (1), the bias terms of hidden units are ignored to simplify the equation. Its corresponding probability is then given by the Boltzmann distribution:

$$\begin{aligned} P(\mathbf{s}; \theta_s) &= \sum_{\mathbf{h}_s^{(1)}, \mathbf{h}_s^{(2)}} P(\mathbf{s}, \mathbf{h}_s^{(1)}, \mathbf{h}_s^{(2)}; \theta_s) \\ &= \frac{1}{Z(\theta_s)} \sum_{\mathbf{h}_s^{(1)}, \mathbf{h}_s^{(2)}} e^{-E(\mathbf{s}, \mathbf{h}_s^{(1)}, \mathbf{h}_s^{(2)}; \theta_s)} \end{aligned} \quad (2)$$

where  $Z(\theta_s)$  is the partition function.

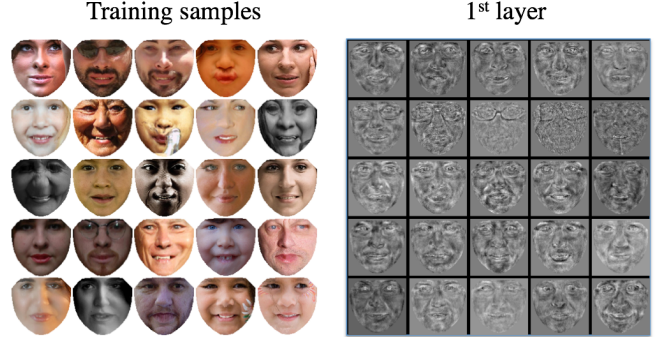
The conditional distributions over  $\mathbf{s}$ ,  $\mathbf{h}_s^{(1)}$ , and  $\mathbf{h}_s^{(2)}$  are then given as in Eqn. (3).

$$\begin{aligned} p(h_{s_j}^{(1)} | \mathbf{s}, \mathbf{h}_s^{(2)}) &= \delta \left( \sum_i W_{s_{ij}}^{(1)} \frac{s_i}{\sigma_{s_i}} + \sum_l W_{s_{jl}}^{(2)} h_{s_l}^{(2)} \right) \\ p(h_{s_l}^{(2)} | \mathbf{h}_s^{(1)}) &= \delta \left( \sum_j W_{s_{jl}}^{(2)} h_{s_j}^{(1)} \right) \\ s_i | \mathbf{h}_s^{(1)} &\sim \mathcal{N} \left( \sigma_{s_i} \sum_j W_{s_{ij}}^{(1)} h_{s_j}^{(1)} + b_{s_i}, \sigma_{s_i}^2 \right) \end{aligned} \quad (3)$$

where  $\delta(x) = 1/(1 + \exp(-x))$  is the logistic function. Figure 3 illustrates a subset of training shapes together with samples generated from shape model after 10-step Gibbs sampling. From this, one can see that the shape model is able to capture the overall shape structure as well as a wide range of head poses and expressions.

### 3.2 Texture Modeling

As opposed to facial shapes, the appearance of human face usually varies drastically due to numerous factors such as identities, lighting conditions, facial occlusions, expressions, image resolutions, etc. These factors can significantly change pixel values presented in these textures and result in much higher non-linear variations. Therefore, the process of texture modeling is more complicated and requires the texture model to be sophisticated enough to represent the variations.



**Fig. 4** A subset of training faces and learned features of the first layer texture model.

The structure of texture model is represented in the right part of DAMs in Figure 2. Different from the shape model which directly works with landmark coordinates in image domain  $\mathcal{I} \subset \mathbb{R}^2$ , the given facial image is first warped from  $\mathcal{I}$  to texture domain  $\mathcal{D} \subset \mathbb{R}^2$  using a reference candidate obtained from the training data. Then the obtained shape-free image is vectorized and used as the visible units for texture model. The purpose of warping step is to remove the effect of shape factors from the texture model and, therefore, making it more robust to shape changes during modeling process. Specifically, given an image  $I$ , the texture  $\mathbf{g}$  is computed as

$$\mathbf{g} = \text{vec}(I(W(r_{\mathcal{D}}, \mathbf{s}))) \quad (4)$$

where  $\text{vec}(\cdot)$  is the vectorization operator;  $W(r_{\mathcal{D}}, \mathbf{s}) = r_{\mathcal{I}}$  is the warping operator;  $r_{\mathcal{I}} = (x_I, y_I)$  and  $r_{\mathcal{D}} = (x_i, y_i)$  are the 2D locations in image domain  $\mathcal{I}$  and texture domain  $\mathcal{D}$ , respectively.

A two-layer DBM is then employed to model the distributions of texture feature represented in  $\mathbf{g}$ . Similar to shape model, the GRBM is used in the bottom layer while interactions between hidden units in higher layers are formulated by a binary-binary RBM. The energy of a state  $\{\mathbf{g}, \mathbf{h}_g^{(1)}, \mathbf{h}_g^{(2)}\}$  in texture modeling is given as in Eqn. (5) where  $\{\mathbf{h}_g^{(1)}, \mathbf{h}_g^{(2)}\}$  denote the set of hidden units and  $\theta_g = \{\mathbf{W}_g^{(1)}, \mathbf{W}_g^{(2)}, \sigma_g^2, \mathbf{b}_g\}$  are the model parameters.

$$\begin{aligned} E(\mathbf{g}, \mathbf{h}_g^{(1)}, \mathbf{h}_g^{(2)}; \theta_g) &= \sum_k \frac{(g_k - b_{g_k})^2}{2\sigma_{g_k}^2} - \sum_{k,t} \frac{g_k}{\sigma_{g_k}} W_{g_{kt}}^{(1)} h_{g_t}^{(1)} \\ &\quad - \sum_{t,v} h_{g_t}^{(1)} W_{g_{tv}}^{(2)} h_{g_v}^{(2)} \end{aligned} \quad (5)$$

The probability of  $\mathbf{g}$  assigned by the model is computed as:

$$\begin{aligned} P(\mathbf{g}; \theta_g) &= \sum_{\mathbf{h}_g^{(1)}, \mathbf{h}_g^{(2)}} P(\mathbf{g}, \mathbf{h}_g^{(1)}, \mathbf{h}_g^{(2)}; \theta_g) \\ &= \frac{1}{Z(\theta_g)} \sum_{\mathbf{h}_g^{(1)}, \mathbf{h}_g^{(2)}} e^{-E(\mathbf{g}, \mathbf{h}_g^{(1)}, \mathbf{h}_g^{(2)}; \theta_g)} \end{aligned} \quad (6)$$

The conditional distributions over  $\mathbf{g}$ ,  $\mathbf{h}_g^{(1)}$ , and  $\mathbf{h}_g^{(2)}$  are derived similar to those of shape model as in Eqn. (7). Figure 4 illustrates a subset of training texture and the learned feature obtained using the first layer of the presented texture model.

$$\begin{aligned} p(h_{gt}^{(1)} | \mathbf{g}, \mathbf{h}_g^{(2)}) &= \delta \left( \sum_k W_{gkt}^{(1)} \frac{g_k}{\sigma_{gk}} + \sum_v W_{gtv}^{(2)} h_{gv}^{(2)} \right) \\ p(h_{gv}^{(2)} | \mathbf{h}_g^{(1)}) &= \delta \left( \sum_t W_{gtv}^{(2)} h_{gt}^{(1)} \right) \\ g_k | \mathbf{h}_g^{(1)} &\sim \mathcal{N} \left( \sigma_{gk} \sum_t W_{gkt}^{(1)} h_{gt}^{(1)} + b_{gk}, \sigma_{gk}^2 \right) \end{aligned} \quad (7)$$

### 3.3 Appearance Modeling

A straightforward way to extract model parameters for both shape and texture is to employ a weighted concatenation and apply a dimensional reduction method such as PCA. However, this is not an optimal solution since these parameters are presented in different domains, i.e. shape parameters  $\alpha_s$  determine the coordinates of landmark points while texture parameters  $\alpha_g$  present facial appearance in the texture domain  $\mathcal{D}$ . Therefore, the gaps between them still exist in the final model parameters although weight values are employed to balance the combined features.

Meanwhile, our Deep Appearance Models also aim to produce a robust facial shape and texture representation. It, however, can be considered as the problem of data learning from multiple sources. In this problem, the information learned from multiple input channels can complement each other and boost the overall performance of the whole model. Particularly, captions and tags can be used to improve the classification accuracy (Huiskes et al, 2010; Ngiam et al, 2011; Srivastava and Salakhutdinov, 2012).

In order to generate a robust feature in DAMs, one should notice that the hidden units are powerful in terms of increasing the flexibility of deep model. Besides the ability of capturing different factors from the observations, the higher layer these hidden units are in, the more independent of the specific correlations of an input source (Srivastava and Salakhutdinov, 2012). Therefore, we can use them as a source-free representation. From that reason, we construct one more high-level layer to interpret the connections between face shape and its texture. Since  $\mathbf{h}_s^{(2)}$  and  $\mathbf{h}_g^{(2)}$  are independent of the spaces where the coordinates and appearance are in, the new layer can encode the shape and texture information more naturally and effectively.

Let  $\mathbf{h}^{(3)}$  be the connection layer and  $\theta = \{\theta_s, \theta_g\}$ , the energy of the joint configuration  $\{\mathbf{s}, \mathbf{g}, \mathbf{h}_s^{(1)}, \mathbf{h}_s^{(2)}, \mathbf{h}_g^{(1)}, \mathbf{h}_g^{(2)}, \mathbf{h}^{(3)}\}$  in DAMs is defined as the summation of three energy

functions of shape model, texture model and the joint layer.

$$\begin{aligned} E(\mathbf{s}, \mathbf{g}, \mathbf{h}_s, \mathbf{h}_g; \theta) &= \sum_i \frac{(s_i - b_{s_i})^2}{2\sigma_{s_i}^2} - \sum_{i,j} \frac{s_i}{\sigma_{s_i}} W_{sij}^{(1)} h_{sj}^{(1)} - \sum_{j,l} h_{sj}^{(1)} W_{sjl}^{(2)} h_{sl}^{(2)} \\ &+ \sum_k \frac{(g_k - b_{g_k})^2}{2\sigma_{g_k}^2} - \sum_{k,t} \frac{g_k}{\sigma_{g_k}} W_{gkt}^{(1)} h_{gt}^{(1)} - \sum_{t,v} h_{gt}^{(1)} W_{gtv}^{(2)} h_{gv}^{(2)} \\ &- \sum_{l,n} h_{sl}^{(2)} W_{sln}^{(3)} h_n^{(3)} - \sum_{v,n} h_{gv}^{(2)} W_{gvn}^{(3)} h_n^{(3)} \end{aligned} \quad (8)$$

where  $\mathbf{h}_s = \{\mathbf{h}_s^{(1)}, \mathbf{h}_s^{(2)}\}$  and  $\mathbf{h}_g = \{\mathbf{h}_g^{(1)}, \mathbf{h}_g^{(2)}\}$ . The joint distribution over the multimodal input can be written as:

$$\begin{aligned} P(\mathbf{s}, \mathbf{g}; \theta) &= \sum_{\mathbf{h}_s^{(2)}, \mathbf{h}_g^{(2)}, \mathbf{h}^{(3)}} P(\mathbf{h}_s^{(2)}, \mathbf{h}_g^{(2)}, \mathbf{h}^{(3)}) \\ &\left( \sum_{\mathbf{h}_s^{(1)}} P(\mathbf{s}, \mathbf{h}_s^{(1)}, \mathbf{h}_s^{(2)}) \right) \left( \sum_{\mathbf{h}_g^{(1)}} P(\mathbf{g}, \mathbf{h}_g^{(1)}, \mathbf{h}_g^{(2)}) \right) \end{aligned} \quad (9)$$

and the conditional distributions over  $\mathbf{h}_s^{(2)}$ ,  $\mathbf{h}_g^{(2)}$ , and  $\mathbf{h}^{(3)}$  are derived as

$$\begin{aligned} p(h_{sl}^{(2)} | \mathbf{h}_s^{(1)}, \mathbf{h}^{(3)}) &= \delta \left( \sum_j W_{sjl}^{(2)} h_{sj}^{(1)} + \sum_n W_{sln}^{(3)} h_n^{(3)} \right) \\ p(h_{gv}^{(2)} | \mathbf{h}_g^{(1)}, \mathbf{h}^{(3)}) &= \delta \left( \sum_t W_{gtv}^{(2)} h_{gt}^{(1)} + \sum_n W_{gvn}^{(3)} h_n^{(3)} \right) \\ p(h_n^{(3)} | \mathbf{h}_s^{(2)}, \mathbf{h}_g^{(2)}) &= \delta \left( \sum_l W_{sln}^{(3)} h_{sl}^{(2)} + \sum_v W_{gvn}^{(3)} h_{gv}^{(2)} \right) \end{aligned} \quad (10)$$

Other conditional distributions over  $\mathbf{s}$ ,  $\mathbf{g}$ ,  $\mathbf{h}_s^{(1)}$  and  $\mathbf{h}_g^{(1)}$  are the same as in Eqns. (3) and (7).

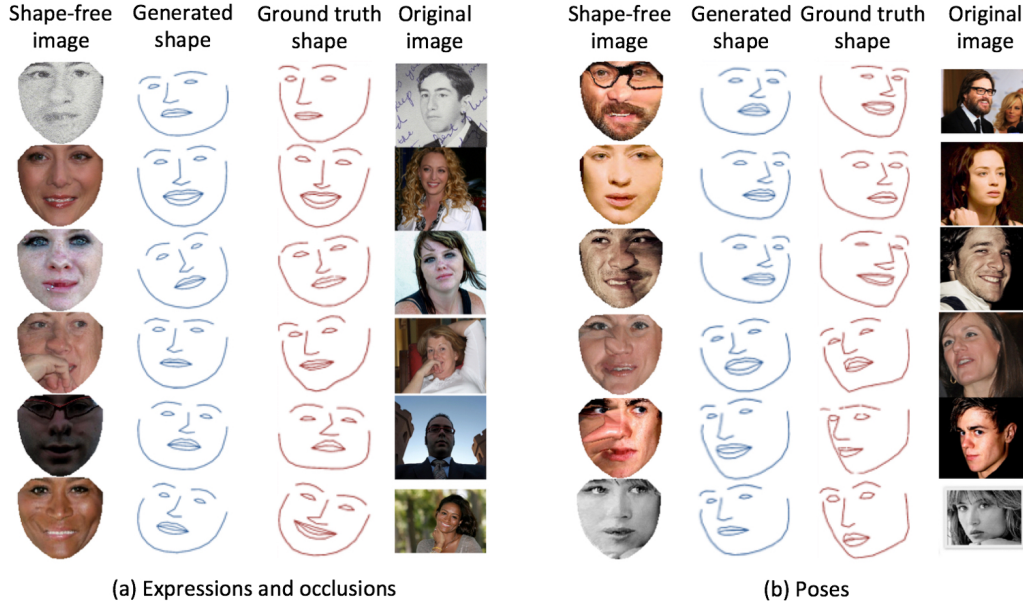
### 3.4 Model Learning

The parameters in the model are optimized in order to maximize the log likelihood  $\theta^* = \arg \max_{\theta} \log P(\mathbf{s}, \mathbf{g}; \theta)$ . Then the optimal parameter values can be obtained in a gradient descent fashion given by

$$\frac{\partial}{\partial \theta} \mathbb{E} [\log P(\mathbf{s}, \mathbf{g}; \theta)] = \mathbb{E}_{\text{data}} \left[ \frac{\partial E}{\partial \theta} \right] - \mathbb{E}_{\text{model}} \left[ \frac{\partial E}{\partial \theta} \right] \quad (11)$$

where  $\mathbb{E}_{\text{data}} [\cdot]$  and  $\mathbb{E}_{\text{model}} [\cdot]$  are the expectations with respect to data distribution, i.e. *data-dependent expectation*, and distribution estimated by DAM, i.e. *model's expectation*. The former term can be approximated by mean-field inference while the latter term can be estimated using Markov-chain Monte-Carlo (MCMC) based stochastic approximation.

**Computing Data-dependent Expectation:** Mean-field approximation can be used to compute the first term of Eqn. (11) (Salakhutdinov and Hinton, 2009). The main idea of this technique comes from the variational approach where the lower bound of the log-likelihood is maximized with respect to the variational parameters  $\mu$ . In the mean-field approximation, for each training face with its shape and texture  $\mathbf{s}, \mathbf{g}$ , all visible units corresponding to  $\mathbf{s}$  and  $\mathbf{g}$  are fixed and



**Fig. 5** Facial shape generation using texture information with (a) expressions and occlusions; and (b) poses. In both cases, given the shape-free image (first column), DAMs are able to generate the facial shape (second column) by sampling from  $P(\mathbf{s}|\mathbf{g}, \theta)$ . The ground truth shapes and original images are also given in the third and fourth columns, respectively.

the states of hidden units in the models are set to  $\mu$  which are iteratively updated through layers using mean-field fixed-point equations:

$$\begin{aligned}
 \mu_{sj}^{(1)} &\leftarrow \delta \left( \sum_i W_{sij}^{(1)} \frac{s_i}{\sigma_{s_i}} + \sum_l W_{sjl}^{(2)} \mu_{sl}^{(2)} \right) \\
 \mu_{sl}^{(2)} &\leftarrow \delta \left( \sum_j W_{sji}^{(2)} \mu_{sj}^{(1)} + \sum_n W_{sln}^{(3)} \mu_n^{(3)} \right) \\
 \mu_{gt}^{(1)} &\leftarrow \delta \left( \sum_k W_{gkt}^{(1)} \frac{g_k}{\sigma_{g_k}} + \sum_v W_{gtv}^{(2)} \mu_{gv}^{(2)} \right) \\
 \mu_{gv}^{(2)} &\leftarrow \delta \left( \sum_t W_{gtv}^{(2)} \mu_{gt}^{(1)} + \sum_n W_{gvn}^{(3)} \mu_n^{(3)} \right) \\
 \mu_n^{(3)} &\leftarrow \delta \left( \sum_l W_{sln}^{(3)} \mu_{sl}^{(2)} + \sum_v W_{gvn}^{(3)} \mu_{gv}^{(2)} \right)
 \end{aligned} \tag{12}$$

Using these variational parameters, the data-dependent statistics are then computed by averaging over training cases.

**Computing Expectation of the Model:** For the second term of Eqn. (11), the MCMC sampling can be applied (Salakhutdinov, 2009). Specifically, given the current state of visible and hidden units, their new states are obtained by employing a few steps of persistent Gibbs sampling using Eqns. (3), (7) and (10). Then the  $\mathbb{E}_{\text{model}}[\cdot]$  is approximated by the expectations with respect to the new states of the model units.

### 3.5 Properties of Deep Appearance Models

Deep Appearance Models provide the capability of generating facial shapes using texture information and vice versa. For example, one can predict a facial shape from the appearance using DAMs as follows: (1) clamping the texture information  $\mathbf{g}$  as observations for the texture model and initializing hidden units with random states; (2) performing standard Gibbs sampling as a posterior inference step; and (3) obtaining the reconstructed shape from  $P(\mathbf{s}|\mathbf{g}; \theta)$ . To generate the appearance from a given shape, one can apply the same way with reversed pathways after clamping the shape information to the shape model. Figure 5 represents the generated shapes given textures in three cases of expressions, occlusions and poses. In all these cases, the DAMs model is able to predict the shape correctly.

In addition, it is more natural to interpret both shapes and textures using higher hidden layers. In order to obtain this representation, one can clamp both observed shape  $\mathbf{s}$  and texture  $\mathbf{g}$  together before applying the Gibbs sampling procedure to estimate  $P(\mathbf{h}^{(3)}|\mathbf{s}, \mathbf{g}; \theta)$ . Eventually, probabilities of these hidden layers can be used as features. Notice that, besides the advantage of better features for discriminative tasks, one can easily see that even when one of two inputs is missing (i.e. shape),  $P(\mathbf{h}^{(3)}|\mathbf{g}; \theta)$  is still able to approximate. Hence, DAMs can be considered as a more generative model compared to other appearance models.

In terms of the number of training samples, with a specific topology of a stochastic neural network and sampling based weight update process (i.e. Contrastive Divergence),



**Fig. 6** Facial image super-resolution reconstruction at different scales of down-sampling. The 1st row: original image, the 2nd row to the 5th row: down-scaled images with factors of 4, 6, 8, 12 (left) and reconstructed facial images using DAMs (right).

the need for large-scale training data is alleviated in our DAMs structure. Therefore, this is an advantage of our deep structure since it does not require a large-scale training dataset as other CNN based approaches.

The proposed method can also deal with facial reconstruction in various challenging conditions, such as: facial occlusions, facial expressions, facial off-angles, etc. These advantages of this method will be shown in Section 5.

## 4 Fitting in Deep Appearance Models

### 4.1 Forward Composition Based Fitting

Given a testing face  $I$ , the fitting process in DAMs can be formulated as finding an optimal shape  $\mathbf{s}$  that maximizes the probability of the shape-free image as in Eqn. (13).

$$\mathbf{s}^* = \arg \max_{\mathbf{s}} P(I(W(r_{\mathcal{D}}, \mathbf{s})) | \mathbf{s}; \theta) \quad (13)$$

Since the connections between textures and hidden units  $\mathbf{h}_g^{(1)}$  are modeled by a GRBM, the probability of texture  $\mathbf{g}$  given hidden units  $\mathbf{h}_g^{(1)}$  is computed as:

$$P(\mathbf{g} | \mathbf{h}_g^{(1)}; \mathbf{s}, \theta) = \mathcal{N}(\sigma_g \mathbf{W}_g^{(1)} \mathbf{h}_g^{(1)} + \mathbf{b}_g, \sigma_g^2 \mathbf{A}) \quad (14)$$

where  $\mathbf{A}$  is the identity matrix;  $\{\sigma_g, \mathbf{b}_g\}$  are the standard-deviation and bias of visible units in the texture model; and  $\mathbf{W}_g^{(1)}$  are learned weights of the visible-hidden texture.

During the fitting steps, the states of hidden units  $\mathbf{h}_g^{(1)}$  are estimated by clamping both the current shape  $\mathbf{s}$  and the texture  $\mathbf{g}$  to the model. The Gibbs sampling method is then applied to find the optimal estimated texture of the testing face given a current shape  $\mathbf{s}$ . By this way, the hidden units in DAMs can take into account both shape and texture information in order to reconstruct a better texture for further refinement.

Let  $\mathbf{m} = \sigma_g \mathbf{W}_g^{(1)} \mathbf{h}_g^{(1)} + \mathbf{b}_g$  be the mean of the Gaussian distribution, we have the following approximation:

$$P(I(W(r_{\mathcal{D}}, \mathbf{s})) | \mathbf{h}_g^{(1)}; \theta) = \mathcal{N}(\mathbf{m}, \sigma_g^2 \mathbf{A}) \quad (15)$$

The maximum likelihood can be then estimated as follows:

$$\begin{aligned} \mathbf{s}^* &= \arg \max_{\mathbf{s}} (P(I(W(r_{\mathcal{D}}, \mathbf{s})) | \mathbf{s}; \theta)) \\ &= \arg \max_{\mathbf{s}} \mathcal{N}(I(W(r_{\mathcal{D}}, \mathbf{s})) | \mathbf{m}, \sigma_g^2 \mathbf{A}) \\ &= \arg \min_{\mathbf{s}} \frac{1}{\sigma_g^2} \sum (I(W(r_{\mathcal{D}}, \mathbf{s})) - \mathbf{m})^2 \end{aligned} \quad (16)$$

Then the forward compositional algorithm can be used to solve the problem (16) by finding the updating parameter  $\Delta \mathbf{s}$  that increases the likelihood:

$$\Delta \mathbf{s} = \arg \min_{\Delta \mathbf{s}} \|I(W(W(r_{\mathcal{D}}, \Delta \mathbf{s}), \mathbf{s})) - \mathbf{m}\|^2 \quad (17)$$

The linearization is taken place of the test image coordinate using first order Taylor expansion  $I(W(W(r_{\mathcal{D}}, \Delta \mathbf{s}), \mathbf{s})) = I(W(r_{\mathcal{D}}, \mathbf{s})) + \mathbf{J}_I \Delta \mathbf{s}$  and the update parameter is given as:

$$\Delta \mathbf{s} = -(\mathbf{J}_I^T \mathbf{J}_I)^{-1} \mathbf{J}_I^T [I(W(r_{\mathcal{D}}, \mathbf{s})) - \mathbf{m}] \quad (18)$$

where  $\mathbf{J}_I = \nabla I \frac{\partial W}{\partial \mathbf{s}}$  is the Jacobian.

### 4.2 Dictionary Learning Based Fitting

In this section, we further improve the fitting process so that it can deal with occlusions and other variations. From Eqn. (17), we can see that the shape update  $\Delta \mathbf{s}$  mostly relies on the difference between the shape-free image and its DAMs reconstruction. However, this metric is easily affected by the presence of occlusions. In particular, when part of the face is occluded, the occlusion still remained in the shape-free image as the result of warping operator but will be removed in DAMs reconstruction. This will lead to the cases where the  $\ell_2$  distance between the two images is still very large even the optimization approaches the correct shape. As a result, it will mislead the optimization process and the final shape cannot be optimized. Therefore, the  $\ell_2$ -norm of their difference is not robust enough to guide the fitting process to the true shape when occlusions occur.



To address this problem more effectively, instead of working directly in texture space, we define a mapping function  $f : \mathcal{I} \mapsto \mathcal{C}$  to map  $I(W(r_{\mathcal{D}}, \mathbf{s}))$  and  $\mathbf{m}$  from texture space  $\mathcal{I}$  to a parameter space  $\mathcal{C}$  such that the relationship between  $f(I(W(r_{\mathcal{D}}, \mathbf{s})))$  and  $f(\mathbf{m})$  is more robust to occlusions. Then this relationship can be used for fitting process. The mapping function  $f$  can be defined as

$$\begin{aligned} f : \mathcal{I} &\mapsto \mathcal{C} \\ \mathbf{c}_1 &= f(I(W(r_{\mathcal{D}}, \mathbf{s}))) \\ \mathbf{c}_2 &= f(\mathbf{m}) \end{aligned} \quad (19)$$

Then it can be parameterized by dictionaries and representation coefficients as follows.

$$\begin{aligned} f(I(W(r_{\mathcal{D}}, \mathbf{s}))) &= \arg \min_{\mathbf{c}_1} \| I(W(r_{\mathcal{D}}, \mathbf{s})) - \hat{\mathbf{D}}_I \mathbf{c}_1 \|_2^2 \\ &\quad + \lambda_1 \| \mathbf{c}_1 \|_1 \\ f(\mathbf{m}) &= \arg \min_{\mathbf{c}_2} \| \mathbf{m} - \hat{\mathbf{D}}_m \mathbf{c}_2 \|_2^2 + \lambda_2 \| \mathbf{c}_2 \|_1 \end{aligned} \quad (20)$$

where  $\{\hat{\mathbf{D}}_I, \mathbf{c}_1\}$  and  $\{\hat{\mathbf{D}}_m, \mathbf{c}_2\}$  are the dictionaries and representation coefficients of the shape free image and its DAMs reconstruction, respectively. The DAMs fitting can be decomposed into two steps, i.e. training and testing.

**Training step:** Given a training dataset with  $N$  images and their shapes  $\{(I^i, \mathbf{s}^i)\}_{i=1}^N$ , the dictionaries are learned by minimizing the objective function

$$\begin{aligned} \{\hat{\mathbf{D}}_I, \hat{\mathbf{D}}_m\} &= \arg \min_{\mathbf{D}_I, \mathbf{D}_m \in \mathbb{R}^{k \times l}} \frac{1}{N} \sum_{i=1}^N \{ \min_{\mathbf{c}^i \in \mathbb{R}^l} \| \mathbf{I}_W^i - \mathbf{D}_I \mathbf{c}^i \|_2^2 \\ &\quad + \| \mathbf{m}^i - \mathbf{D}_m \mathbf{c}^i \|_2^2 \\ &\quad + \lambda \| \mathbf{c}^i \|_1 \} \end{aligned} \quad (21)$$

where  $\mathbf{I}_W^i = I^i(W(r_{\mathcal{D}}, \mathbf{s}))$ ,  $k$  is the length of texture vector and  $l$  is the size of dictionaries. With this objective function, the two dictionaries  $\hat{\mathbf{D}}_I$  and  $\hat{\mathbf{D}}_m$  are learned in a way that regardless of the present of occlusions in the input face, the extracted representation coefficients of  $I(W(r_{\mathcal{D}}, \mathbf{s}))$  and  $\mathbf{m}$  are forced to share the same (i.e.  $\mathbf{c} = \mathbf{c}_1 = \mathbf{c}_2$ ) when the true shape is approaching. To solve (21), we apply the four-step iterative procedure as in (Xing et al, 2014). The main steps of this procedure are summarized in Algorithm 1. There are two main advantages of learning the dictionaries as in Eqn. (21). Firstly, since both shape-free image and its DAMs reconstruction are forced to share the same representation  $\mathbf{c}^i$ , their underlying relationships are naturally embedded in these coefficients. Secondly, when the vector  $\mathbf{c}^i$  is sparse, the optimization will result in the most related features between the shape-free image and its reconstruction. Therefore, it will be more robust to occlusions and other variations.

After obtaining the dictionaries, instead of following the Gaussian-Newton optimization as in Eqn. (17), we learn a linear regressor to directly infer the shape update  $\Delta \mathbf{s}$  from the difference between  $f(I_W)$  and  $f(\mathbf{m})$ . Specifically, given training images with their initial estimate  $\{\bar{\mathbf{s}}^i\}_{i=1}^N$  of their

---

### Algorithm 1 Dictionary Learning for fitting

---

**Require:** Training data  $\{(I^i, \mathbf{s}^i)\}_{i=1}^N$ , regularization parameter  $\lambda$

**Ensure:** Learned dictionaries  $\{\hat{\mathbf{D}}_I, \hat{\mathbf{D}}_m\}$

- 1: Construct the matrix  $\mathbf{Y} \in \mathbb{R}^{k \times N}$  whose  $i$ -th column is shape-free image  $\mathbf{I}_W^i$ .
  - 2: Construct the matrix  $\mathbf{M} \in \mathbb{R}^{k \times N}$  whose  $i$ -th column is  $\mathbf{m}^i$ .
  - 3: Initialize  $\mathbf{D}_I \in \mathbb{R}^{k \times l}$  and  $\mathbf{D}_m \in \mathbb{R}^{k \times l}$  with random samples from a normal distribution with zero mean and unit variance.
  - 4: **while** not converged **do**
  - 5:   (1) Fix  $\mathbf{D}_m$ , learn  $\mathbf{D}_I$  and coefficient matrix  $\mathbf{C} \in \mathbb{R}^{l \times N}$ 

$$\{\mathbf{D}_I, \mathbf{C}\} = \arg \min_{\mathbf{D}_I, \mathbf{C}} \| \mathbf{Y} - \mathbf{D}_I \mathbf{C} \|_2^2 + \lambda \| \mathbf{C} \|_1$$
  - 6:   (2) Update  $\mathbf{D}_m$  as  $\mathbf{D}_m = \mathbf{M} / \mathbf{C}$ . Notice that this result is used as initial  $\mathbf{D}_m$  for step (3).
  - 7:   (3) Fix  $\mathbf{D}_I$ , learn  $\mathbf{D}_m$  and new coefficient matrix  $\mathbf{C}$ 

$$\{\mathbf{D}_m, \mathbf{C}\} = \arg \min_{\mathbf{D}_m, \mathbf{C}} \| \mathbf{M} - \mathbf{D}_m \mathbf{C} \|_2^2 + \lambda \| \mathbf{C} \|_1$$
  - 8:   (4) Update  $\mathbf{D}_I$  as  $\mathbf{D}_I = \mathbf{Y} / \mathbf{C}$ .
  - 9: **end while**
  - 10: Set  $\hat{\mathbf{D}}_I = \mathbf{D}_I$  and  $\hat{\mathbf{D}}_m = \mathbf{D}_m$ .
- 

ground truth shape, the linear regressor is learned by minimizing

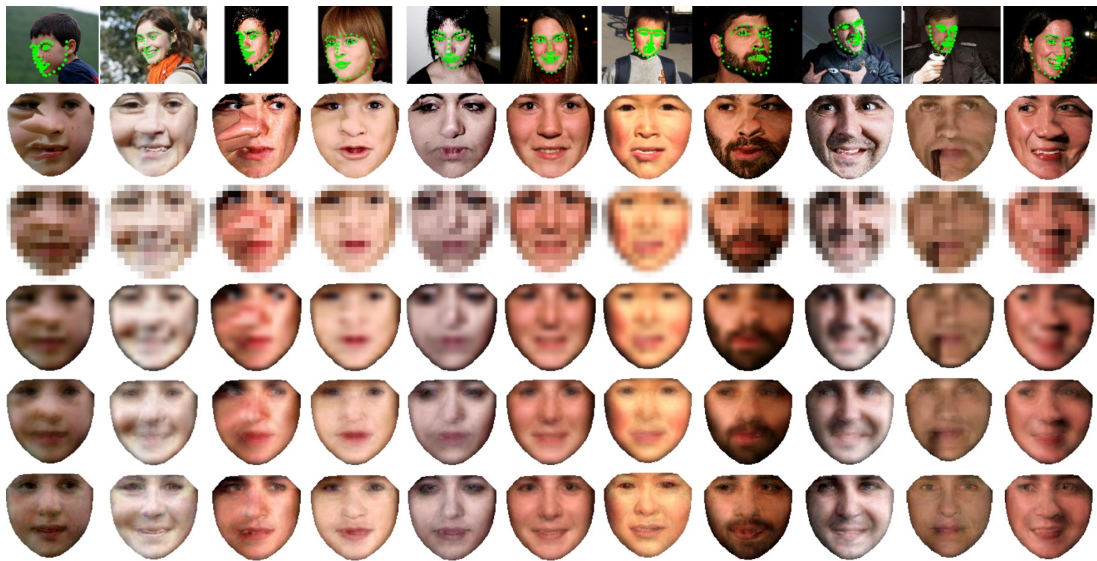
$$\arg \min_{\mathbf{H}, \mathbf{b}} \sum_{i=1}^N \| \Delta \mathbf{s}^i - \mathbf{H} [f(I_W^i) - f(\mathbf{m}^i)] - \mathbf{b} \|^2 \quad (22)$$

where  $\Delta \mathbf{s}^i = \mathbf{s}^i - \bar{\mathbf{s}}^i$ ;  $\{\mathbf{H}, \mathbf{b}\}$  are the regressor's parameters.

**Testing step:** In the fitting process, given an input face with its initial shape, the shape-free image  $I_W$  and DAMs reconstruction  $\mathbf{m}$  are first computed. Their representation coefficients  $\mathbf{c}_1^*$  and  $\mathbf{c}_2^*$  are also estimated using Eqn. (20). Then the shape is updated using the difference  $\mathbf{c}_1^* - \mathbf{c}_2^*$  together with the learned regressor. After that the image pair  $(I_W, m)$  is recomputed for the next iteration.

## 5 Experimental Results

In section 5.1, we briefly introduce the main features of the three databases used in our evaluations. They consist of two ‘‘face in the wild’’ and one aging databases. By using these databases with numerous challenging factors, we aim to show the robustness and efficiency of our proposed model. Then, in the next three sections, we validate the generative capabilities of our DAMs in both facial representation and reconstruction via four applications, i.e. facial super-resolution, facial off-angle reconstruction, facial occlusion removal, and facial age estimation. The experiments are also made to be more challenging by including numerous variations in poses, occlusions and impulsive noise. Comparing to AAMs, bicubic interpolation and other deep learning based approaches, our DAMs achieve better reconstructions without blurring



**Fig. 7** Facial image super-resolution. The original images (first row) are warped to shape-free images in texture domain (second row); then they are down-sampled by a factor of 8 from  $117 \times 120$  to  $15 \times 15$  (third row) The next three rows are the high-resolution reconstructed using Bicubic method (the fourth row), PCA-based AAMs (the fifth row) and Deep Appearance Models (the sixth row).

effects or spreading out the errors caused by occlusions or noise. We also represent in section 5.5 an experiment to evaluate our proposed DAMs method in the ability of synthesizing new face images. Its performance is compared with AAMs and other face alignment methods such as RCPR (Burgos-Artizzu et al, 2013), and CNN based approach (Sun et al, 2013).

## 5.1 Databases

We aim to build a model that can represent face texture in-the-wild. Therefore, in the first three applications, we evaluate DAMs on two face databases in-the-wild, i.e. Labeled Face Parts in the Wild (LFPW) (Belhumeur et al, 2011) and Helen (Le et al, 2012). These databases contain unconstrained facial images collected from various multimedia resources. These facial images have considerable resolutions and contain numerous variations such as poses, occlusions and expressions. For the age estimation application, FG-NET face aging database<sup>3</sup> is used to evaluate the method.

The LFPW database contains 1400 images in total with 1100 training and 300 testing images. However, a part of it is no longer accessible. Therefore, in our experiments, we only use 811 training and 224 testing images, the available remaining. Each facial image is annotated with 68 landmark points provided by 300-W competition (Sagonas et al, 2013).

The Helen database provides a high-resolution dataset with 2000 images used for training and 330 images for testing. The variations consist of pose changing from  $-30^\circ$  to  $30^\circ$ ; several types of expression such as neutral, surprise, smile, scream; and occlusions. Similar to LFPW, all faces in Helen are also annotated with 68 landmark points.

FG-NET is a popular face aging database. There are 1002 face images of 82 subjects with age ranges from 0 to 69 years. The annotations in FG-NET are also 68 landmarks in the same format as LFPW and Helen databases.

## 5.2 Facial Super-resolution Reconstruction

The proposed DAMs method is evaluated in its capability to recover high-resolution face images given their very low-resolution versions. Moreover, since LFPW and Helen databases also include numerous variations in poses, expressions and occlusions, the experiment becomes more challenging. Our proposed method is very potential in dealing with the problem of super-resolution in various conditions of facial poses and occlusions.

In order to train the DAMs model, we combine 811 training images from LFPW and 2000 images from Helen database into one training set. The coordinates of facial landmarks are normalized to zero mean before setting as observations to train the shape model. In the texture modeling, shape-free images are first extracted by warping faces into the texture domain  $\mathcal{D}$ . The size of the shape-free image is set to  $117 \times 120$  pixels based on the mean shape of the training data. Then texture model is trained to learn the facial variations represented in these shape-free images.

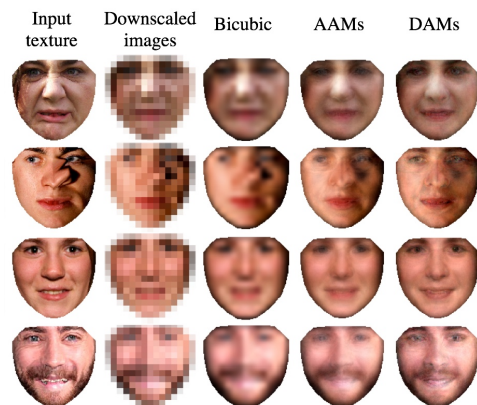
<sup>3</sup> The FG-NET Aging Database, <http://www.fgnet.rsunit.com/>.



**Fig. 8** Comparisons of different facial image super-resolution methods. The 1st row: ground truth faces. The 2nd row: down-scaled images with factors of 6 (left) and 8 (right). From the 3rd row to the 8th row: reconstructed faces using bicubic, PCA-based AAMs, ScSR (Yang et al, 2010), SFH (Yang et al, 2013), SRCNN (Dong et al, 2014) and DAMs, respectively.

During the testing phase, since the number of visible units in the texture model is fixed, the testing low-scale facial shape-free image is first resized to  $117 \times 120$  using bicubic interpolation method. Then both the shape and the shape-free image are clamped to DAMs. After 50 epochs in the alternating Gibbs updates, the face texture is reconstructed based on the current states of hidden unit  $h_g^{(1)}$ . Different magnification factors  $\alpha$  are used for evaluating the quality of DAMs reconstructions. Testing images are down-sampled in different magnification levels ranging from 4 to 12. They are then used as inputs to the reconstruction module of our approach. Figure 6 shows the reconstructions using the DAMs approach. Remarkable results are achieved using DAMs with very low-resolution input images, i.e.  $10 \times 10$  pixels with the magnification factor  $\alpha = 12$ .

**Comparisons against Baseline Methods:** Our proposed approach is also compared with two base-line methods, i.e. bicubic interpolation method and PCA-based AAMs (Tzimiropoulos and Pantic, 2013). Root Mean Square Error (RMSE) is used as a performance measurement. RMSE is a common metric that is usually used for evaluating image recovery task. Although this metric is not always reliable for rating image quality visually (Wang and Bovik, 2009), it could provide a qualitative view for comparing DAMs and other methods.



**Fig. 9** Results of average RMSEs over 4 images: Bicubic interpolation (RMSE = 19.68); PCA-based AAMs reconstruction (RMSE = 19.96); (d) Deep Appearance Models reconstruction (RMSE = 20.44).

From the results shown in Figure 9, our method gives better reconstruction results in visualization than the others. However, the RMSE results are not much better as shown in Table 1. This is because RMSE cannot fully evaluate the quality of reconstructed images in the task of image super-resolution (Yang et al, 2010). Especially, we don't have the ground-truth for RMSE evaluation in these databases.

**Table 1** The average RMSEs of reconstructed images using different methods against LFPW and Helen databases with  $\alpha = 16$ .

| Methods     | LFPW         | Helen        |
|-------------|--------------|--------------|
| Bicubic     | 19.53        | 22.13        |
| AAMs        | 19.74        | 22.3         |
| DAMs (Ours) | <b>19.24</b> | <b>21.24</b> |

**Table 2** The average PSNRs (dB) of different methods on LFPW and Helen.

| Methods | $\alpha = 4$ |              | $\alpha = 6$ |              | $\alpha = 8$ |              |
|---------|--------------|--------------|--------------|--------------|--------------|--------------|
|         | LFPW         | Helen        | LFPW         | Helen        | LFPW         | Helen        |
| Bicubic | 19.47        | 19.22        | 18.82        | 18.59        | 18.29        | 18.07        |
| AAMs    | 26.43        | <b>26.56</b> | 25.72        | 25.83        | 24.82        | 24.97        |
| SRCNN   | <b>26.73</b> | 26.46        | 24.61        | 24.30        | 23.22        | 22.85        |
| DAMs    | 26.46        | 26.54        | <b>25.91</b> | <b>25.97</b> | <b>25.06</b> | <b>25.16</b> |



**Fig. 10** Facial off-angle reconstruction: the 1st row: original image, the 2nd row: shape-free image, the 3rd row: PCA-based AAMs reconstruction, and the 4th-row: DAMs reconstruction

For example, in the cases of occlusions and poses in those databases, although the reconstructed images obtained using PCA-based AAMs and bicubic methods are very blurry, their RMSEs are still low. This is because the reconstructed images still contain occlusion components or pose features which are quite similar to the original ones. Figure 7 illustrates further reconstruction results obtained using bicubic method, PCA-based AAMs method and our DAMs approach. The PCA-based AAMs method is trained using the same dataset as DAMs and the length of texture parameter vector is 200, the highest level used in (Tzimiropoulos and Pantic, 2013)).

**Comparisons against Other Super-resolution Methods:** For further evaluations, we compare DAMs with three other super-resolution methods in Figure 8. The other three super-resolution methods are: sparse representation based image super-resolution (ScSR) (Yang et al, 2010), Structured Face Hallucination (SFH) (Yang et al, 2013), and super-resolution CNN (SRCNN) based approach (Dong et al, 2014). The main difference between the first two is that the former is designed for images in general while the latter is more specific for facial images. The third approach is a deep learning based method. For each subject, the low-resolution (LR) faces (i.e.  $LR_6$  and  $LR_8$ ) are obtained by

down-sampling the ground truth face with factors of 6 and 8. Their high-resolution (HR) reconstructed faces (i.e.  $HR_6$  and  $HR_8$ ) of different methods are shown in the left and right columns underneath each ground truth face in the first row, respectively. The results of Bicubic and AAMs are also presented in this figure. The resolution of  $LR_6$  is  $20 \times 20$  and that of  $LR_8$  is  $15 \times 15$ .

It is clear in the figure that SFH performs better than ScSR in terms of reconstruction details. This is because SFH was already trained with the face structure and contour's statistical priors. However, some noisy and blocky effects are still remained in the reconstructed faces of SFH. Especially, when parts of face images are blurred due to the effects of warping operator, artifacts may appear in the final results. The SRCNN also shows some advantages with small magnification factor  $\alpha$  but the blurry effects are still presented when  $\alpha$  increases. Meanwhile, remarkable results can be achieved by DAMs in terms of keeping fine details without noisy effects. In addition, these results also show the advantages of DAMs when dealing with higher magnification factor  $\alpha$ . Whereas all five methods fail to produce high quality reconstructions when  $\alpha$  increases from 6 to 8, DAMs still perform well and generate faces with consistent quality.

Table 2 presents another qualitative comparison in terms of average PSNR values between our DAMs against one deep learning based (i.e. SRCNN) and two baseline methods. These results again show the effectiveness of DAMs in this task. With small value of  $\alpha$ , DAMs produce comparable results to AAMs and SRCNN. When  $\alpha$  increases (i.e. from 4 to 8), DAMs achieve the highest PSNRs as compared to other methods.

### 5.3 Facial off-angle Reconstruction and Occlusion Removal

This section illustrates the ability of DAMs to deal with facial poses and occlusions.

#### 5.3.1 Facial off-angle Reconstruction

Using the same trained model as in the previous experiment, facial images with different poses are represented in Figure 10. Comparing to AAMs, our DAMs achieve better reconstructions especially in the invisible regions of extreme poses. These regions in shape-free images are blurry and noisy due to the non-linear warping operator. Therefore, the errors are spread out in the reconstructions of PCA-based AAMs approaches. Meanwhile, the generative capability of our proposed DAMs method can solve those challenging cases. From the results, it is easy to see that the blurry effects are effectively removed in DAMs reconstructions.



Fig. 11 Face Frontalization: Top: input faces and Bottom: frontalized faces reconstructed using DAMs.

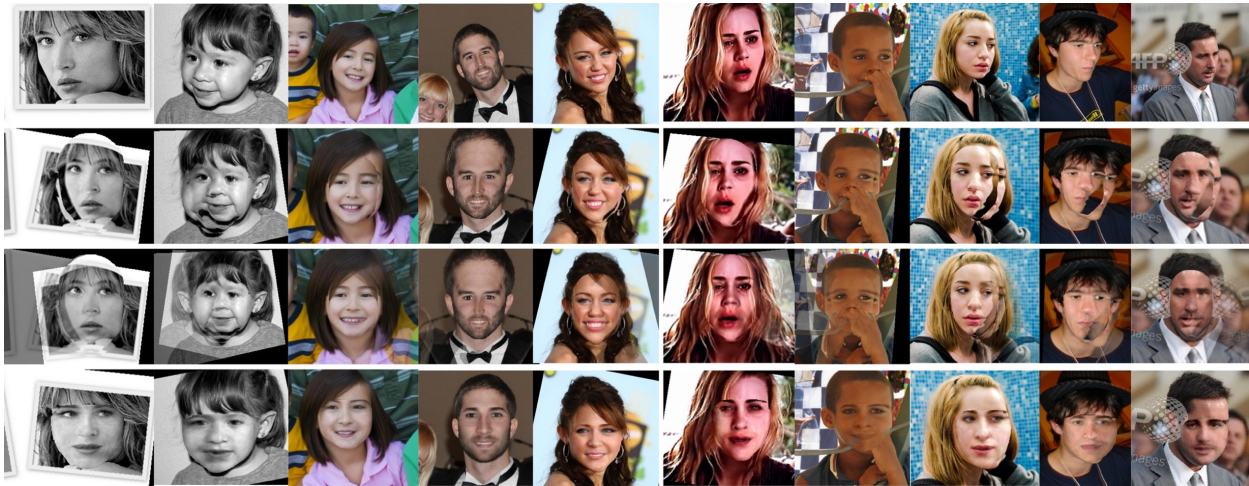


Fig. 12 Comparisons between DAMs and Face Frontalization approach (Hassner et al, 2015). The 1st row: input faces; the 2nd and 3rd rows: synthesized frontal view before and after applying soft symmetry (Hassner et al, 2015); the 4th row: frontalized faces produced by DAMs.

### 5.3.2 Face Frontalization

We next emphasize this ability of our DAMs approach on the face frontalization problem. Given an input face with pose, the process of “frontalization” is to synthesize the frontal view of that face. Notice that the facial photos are unconstrained and the subjects are not required to already be in the training data. Once again, in order to produce aesthetic frontal view, not only poses but other factors such as expressions and occlusions are needed to be taken into account. The frontalization can help to boost the performance of other subsequent processes such as face recognition, verification, gender estimation (Hassner et al, 2015), etc. There are several frontalization approaches proposed in literature (Sagonas et al, 2015; Hassner et al, 2015; Wang et al, 2015a; Jeni and Cohn, 2016; Ferrari et al, 2016). Sagonas et al (2015) employed a statistical model built on some frontal images to simultaneously reconstruct frontal view and align faces. Frontal and aligned faces are then obtained by solving a rank minimization problem. Inspired by the use of 3D model for frontal face reconstruction, Hassner et al (2015) proposed to approximate facial shapes using a single 3D surface while Wang et al (2015a) tried to fit the 3D model to the input 2D face by iteratively refining the 3D landmarks and the weighting coefficients of each landmark. Jeni and Cohn

Table 3 The face verification accuracies on LFW benchmark.

| Methods             | Original inputs | LFW3D  | DAMs          |
|---------------------|-----------------|--------|---------------|
| Pittpatt Classifier | 83.67%          | 86.63% | <b>87.72%</b> |

(2016) later introduced a cascade regression-based approach to estimate a canonical view of the eyes for 3D gaze estimation. The algorithm includes several steps: localizing a set of dense landmarks, fitting a part-based 3D model for 3D shape reconstruction and estimating head pose and gaze. Ferrari et al (2016) proposed to use a 3DMM to fit the input image and then map each image pixel to its 3D corresponding coordinate of the model to obtain a frontal view. By projecting the 3D model back to the frontalized image, image patches can be located and aligned for feature extraction over different images.

Figure 11 represents the frontalized views of input faces with different poses and expressions given in the top row. Our reconstruction results are also compared with the recent frontalization work (Hassner et al, 2015) against LFPW and Helen databases in Figure 12. From the second and third rows, one can see that the approach in (Hassner et al, 2015) achieves good reconstructions when the input poses are not so extreme (i.e. not greater than 30 degrees). However, in case of extreme poses (i.e. the first two and the last three



**Fig. 13** Occlusion removal: the 1st row: original image, the 2nd row: shape-free image, the 3rd row: PCA-based AAMs reconstruction still remains with occlusion and blurring effects, and the 4th-row: DAMs reconstruction can help to remove the occlusion

faces) or occlusions (the 7th face), even when the symmetry property is used, the full faces can not be reconstructed aesthetically. Meanwhile, the results in the last row show that DAMs can effectively synthesize the frontal views of these faces without further applying the soft symmetry property. Since the face priors are already learned, DAMs are able to produce more natural faces instead of duplicating the information from known side to the other side. To further compare the effectiveness of DAMs against this frontalization work (i.e. LFW3D), we also employ the face verification protocol on the Labeled Face in the Wild (LFW) dataset and achieve the results as in Table 3. Noticed that applying directly deep learning based classifiers may mask out the contributions of frontalization step, we employ the off-the-shelf commercial face recognition Pittpatt (developed by CMU and Google) in this evaluation. These results show that our DAMs can handle the face poses effectively and produce a higher accuracy boost comparing to LFW3D approach.

There are some other deep learning based frontalization approach such as FIP (Zhu et al, 2013) and MVP (Zhu et al, 2014). However, both approaches require several images of the same identity with different poses to train the deep models. Moreover, other variations such as expressions, illuminations, and occlusions are kept unchanged between the input and output images during the learning stage. These requirements have limited the use of these models with “in-the-wild” databases where each subject has only one image. Meanwhile, our DAMs approach provides a model structure that is able to handle many variations at the same time without requirements on the number of images per subject as well as the view labels during training stage.

### 5.3.3 Facial Occlusion Removal

Similarly, DAMs also show their capability in the problem of facial occlusion removal. In Figure 13, the occlusions, e.g. hands, glasses, hair, etc., can be removed successfully without blurring effects. More interestingly, the occlusions

**Table 4** The rank-1 recognition results (%) on the AR database.

| Methods    | Pixels | AAMs   | RoBM   | DAMs          |
|------------|--------|--------|--------|---------------|
| Sunglasses | 51.24% | 69.40% | 71.89% | <b>76.87%</b> |
| Scarf      | 55.97% | 61.44% | 64.18% | <b>70.41%</b> |

are removed from faces without losing facial features. For example, glasses are totally removed without making beard blurred as in AAMs reconstruction.

Using occluded faces as references and measuring the reconstruction quality by RMSE cannot illustrate the modeling capabilities of DAMs. To get a better evaluation protocol, we select a subset of 174 occluded faces of the first 29 subjects, i.e. 15 males and 14 females, from AR database (Martinez and Benavente, 1998). We employ DAMs to reconstruct these occluded faces and then use their corresponding neutral faces, i.e. frontal face without occlusions, as references to compute the RMSE. In this testing set, each subject includes two faces with scarf and four other faces with both illumination and scarf. The average RMSE of DAMs is 45.08 while that of PCA-based AAMs is 47.36. The trained models in DAMs and AAMs use LFPW and Helen databases as presented in Section 5.1. This experiment shows that DAMs achieve better reconstructions, i.e. closer to the neutral faces, compared to AAMs.

In order to further illustrate the effectiveness of DAMs in the problem of facial occlusion removal, we also compare our DAMs with Robust Boltzmann Machines (RoBM) (Tang et al, 2012b) in terms of recognition performance on AR. We set up a similar protocol as in (Tang et al, 2012b). However, in that protocol, occluded faces of all subjects are also included in the gallery set. This may easily cause a “match” between occluded faces of the same subject in the gallery and probe sets. Therefore, in this experiment, we use only the non-occluded faces (i.e. seven images per subject) to compose the gallery set and leave all occluded faces (i.e. three with sunglasses and three with scarf per subject) for the probe set. The recognition using DAMs consists of first reconstructing the “clean” faces using DAMs followed by a classification based on LDA and the nearest neighbor classifier. The cosine distance is used for the matching score. Table 4 shows the rank-1 recognition results obtained by different models. DAMs outperform other methods and illustrate their effectiveness in handling the face occlusions.

### 5.4 Facial Age Estimation

Besides some other previous age estimation approaches (Fu and Huang, 2008; Luu et al, 2009), we employ our proposed DAMs to this problem to further demonstrate their robustness and effectiveness.

**Evaluation on Reconstructed Images:** Since texture is an important factor to predict a person’s age given his facial

**Table 5** The MAEs (years) of different methods against impulsive noise.

| Methods | No noise    | Noise range |             |             |             |
|---------|-------------|-------------|-------------|-------------|-------------|
|         |             | 25          | 50          | 100         | 150         |
| AAMs    | 6.14        | 6.15        | 6.11        | <b>6.13</b> | 6.47        |
| DAMs    | <b>5.67</b> | <b>5.81</b> | <b>5.56</b> | 6.14        | <b>6.18</b> |

**Table 6** The MAEs (years) of different methods against low-resolution testing faces.

| Methods | Magnification factor $\alpha$ |             |             |             |
|---------|-------------------------------|-------------|-------------|-------------|
|         | 2                             | 4           | 6           | 8           |
| Bicubic | 5.96                          | 6.95        | 7.15        | 7.21        |
| AAMs    | 6.13                          | 6.33        | 6.44        | 6.69        |
| DAMs    | <b>5.91</b>                   | <b>6.00</b> | <b>6.11</b> | <b>6.21</b> |

image, this experiment will evaluate how good the reconstructed image is as well as how much aging information is retained by the model.

To make this task more challenging, we add noise to the testing facial image and then predict the age of that person using “clean” reconstructed face from DAMs. For the evaluation system, we modified the age estimation systems presented in (Luu et al, 2009, 2011a, 2010, 2011b; Duong et al, 2011) with three-group classification in the first step (youths, adults, and elders) before constructing three Support Vector Regression (SVR) based aging functions. Then we train this age estimator with 802 images from FG-NET. The remaining 200 images were used for testing. To generate noisy testing images, all pixels of facial images were mixed with uniform noise ranged within  $[-r, r]$ .

A similar experiment is set up as follows: given the low-resolution testing face, the system will predict the age of that person using his high-resolution reconstructed face. The Mean Absolute Errors (MAEs) of different methods against noise and low-resolution testing faces are represented in Table 5 and Table 6, respectively. From these results, in both cases, the smallest error is achieved with DAMs model. Therefore, our proposed model produces better reconstructed results under the effects of noise and low-resolution factor.

**Evaluation on Model Features:** Besides the ability of generalizing the faces, DAMs can produce a higher level representation for both facial shape and texture. Therefore, instead of using pixel values, we extracted the model parameters as described in Section 3.5 and evaluated them with the age estimation system. For the AAMs features, the number of features for shape and texture was chosen so that 93% of variations are retained. Table 7 lists the MAEs of four different inputs: reconstructed image of DAMs (**DAMs-Rec**) and AAMs (**AAMs-Rec**), model parameters extracted from AAMs (**AAMs-Mod**) and DAMs (**DAMs-Mod**) as well as other age estimation methods. Not surprisingly, our DAMs feature achieves the lowest MAEs as compared with AAMs features. Notice that, in this experiment, although DAMs are not tuned toward the aging labels, the features extracted

**Table 7** Comparison of age estimation results on FG-NET database with four different features and other age estimation approaches.

| Inputs                      | MAEs (years) |
|-----------------------------|--------------|
| DAMs-Mod                    | <b>4.67</b>  |
| AAMs-Mod                    | 4.81         |
| DAMs-Rec                    | 5.67         |
| AAMs-Rec                    | 6.14         |
| DLF-CNN (Wang et al, 2015b) | 4.26         |
| CA-SVR (Chen et al, 2013)   | 4.67         |
| PLO (Li et al, 2012)        | 4.82         |

by DAMs still achieve quite competitive results to others. We believe that with better age estimator, i.e. deep learning based age estimator, and the use of aging labels during training DAMs, the MAE would be reduced significantly. We leave this as our future work of DAMs.

### 5.5 Shape Fitting in DAMs

Besides some other previous shape fitting approaches (Alabort-i Medina et al (2014); Alabort-i Medina and Zafeiriou (2014, 2015, 2017); Antonakos et al (2015); Tzimiropoulos and Pantic (2017)), we employ our proposed DAMs to this problem on LFPW database to further demonstrate their robustness and effectiveness. The model configurations are kept the same as in previous sections except it is now trained with 811 training images of LFPW. For evaluation and comparison, we use the average distance of each landmark to its ground truth position normalized by face size as in (Tzimiropoulos and Pantic, 2013). Moreover, in order to remove the effect of face detection error during fitting step, we use the bounding boxes provided in (Sagonas et al, 2013) for initialization. Then we simply place the mean shape with 68 landmarks inside the face’s bounding box and start the fitting process. We compare our method with two other fitting strategies, i.e. AAMs and RCPR (Burgos-Artiztu et al, 2013), and present the results in Table 8. The Cumulative Error Distribution (CED) curves are showed in Figure 14.

**Comparisons Against AAMs Based Approaches:** We conduct an ablation study to compare between our DAMs and other AAMs based fitting approaches in Table 9. In these comparisons, we also include the feature-based AAMs (Antonakos et al, 2015), i.e. represent texture with HoG and Dense SIFT features, as well as hybrid approaches, i.e. Supervised Descent Method combined with AAMs (SDM + AAMs), proposed in (Antonakos et al, 2016). Notice that to make a fair comparison between our DAMs and other approaches, we use a single-resolution pyramidal scheme for the AAMs model configurations. The Project out Forward Compositional (POFC) fitting algorithm is used in this experiment. For holistic AAMs, the number of appearance parameters is fixed at 50. The dimension of 2D shape pa-

**Table 8** The fitting errors using different methods against LFPW database.

| Methods        | Fitting Error |
|----------------|---------------|
| Initialization | 0.0618        |
| Fast-SIC       | 0.0391        |
| RCPR           | 0.0505        |
| DAMs (Ours)    | 0.0398        |

**Table 9** The fitting errors against other feature-based AAMs methods and hybrid approach (SDM + AAMs).

| Methods         | Fitting Error |
|-----------------|---------------|
| HoG AAMs        | 0.0372        |
| Dense SIFT AAMs | 0.035         |
| SDM + AAMs      | 0.0407        |
| DAMs (Ours)     | 0.0398        |

**Table 10** The fitting errors against different initial bounding boxes.

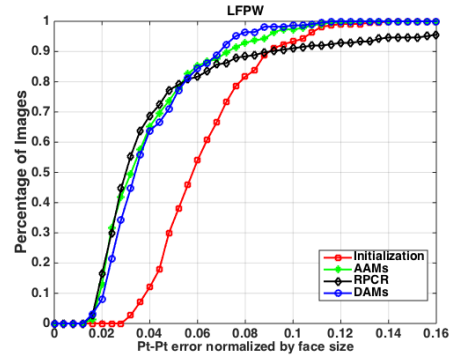
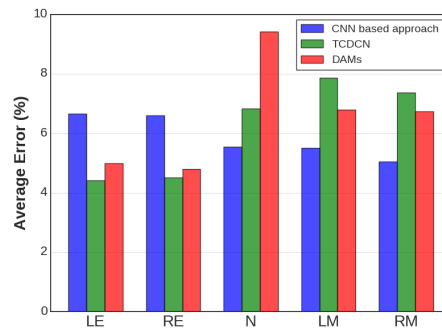
| Face Detection Methods                    | DAMs Fitting Error |
|-------------------------------------------|--------------------|
| Initialization from (Sagonas et al, 2013) | 0.0398             |
| MTCNN (Zhang et al, 2016a)                | 0.0383             |
| CMS-RCNN (Zhu et al, 2017)                | 0.0395             |

rameters is set to 12. The results again show that our DAMs achieve a comparable fitting error with feature-base AAMs.

**Comparisons Against CNN Based Approaches:** We also compare our method with another CNN based landmark detection approach (Sun et al, 2013), and Tasks-Constrained Deep Convolutional Network (TCDCN) (Zhang et al, 2016b). Since these approaches results in a detection of five landmark points (i.e. left eye, right eye, nose, and two mouth corners), we only compute the average distance of these landmarks for comparison. Our DAMs approach achieves the fitting error of 0.028 while the error of the CNN based system and TCDCN are 0.026 and 0.027, respectively. The average detection errors (Sun et al, 2013) of all five landmarks are also presented in Figure 15. These results show that DAMs achieve comparable accuracy to other face alignment methods.

#### The Sensitivity Against Different Shape Initial Bounding Boxes:

In order to evaluate the sensitivity of our DAMs approach to the size of detected bounding boxes during fitting step, we employ various face detection techniques including the initial bounding boxes provided by (Sagonas et al, 2013), MTCNN (Zhang et al, 2016a), CMS-RCNN (Zhu et al, 2017) to obtain different types of bounding boxes. Then, we place the mean shape inside these bounding boxes and start the fitting process. Table 10 presents the fitting accuracy of our DAMs fitting approach when different bounding boxes are used. From these results, one can see that our fitting approach is quite robust to the initial location and scale of the mean shape.

**Fig. 14** Cumulative Error Distribution (CED) curves of LFPW database.**Fig. 15** The average errors (%) of CNN based approach (Sun et al, 2013), TCDCN (Zhang et al, 2016b), and DAMs of different landmarks: left eye (LE), right eye (RE), nose (N), left and right mouth corners (LM and RM).**Table 11** Computational time of DAMs and AAMs in three stages: training the models, fitting and reconstruct an image.

| Stages | Training  | Fitting | Reconstruction |
|--------|-----------|---------|----------------|
| DAMs   | 12.87 hrs | 17.5 s  | 0.53 s         |
| AAMs   | 564.06 s  | 2.28 s  | 0.023 s        |

## 5.6 Computational Costs

The computational costs of DAMs, i.e. training, fitting and reconstruction stages are discussed in this section. Both Helen and LFPW databases are combined to use in this evaluation. The numbers of training and testing images are 2811 and 554, respectively. The method is implemented in Matlab environment and runs in a system of Core i7-2600 @3.4GHz CPU, 8.00 GB RAM. The shape contains 68 landmarks and the appearance is represented in a vector of 9652 dimensions. Each layer was trained using Contrastive Divergence learning in 600 epochs. It is noted that the current version is implemented without using parallel processing. The computational costs of DAMs and AAMs are shown in Table 11.



## 6 Conclusions

This paper has introduced novel Deep Appearance Models that have abilities of generalizing and representing faces in large variations. With the deep structured models for shapes and textures, the proposed approach was shown to achieve remarkable improvements in both facial reconstruction and facial age estimation tasks compared with PCA-based AAMs model. Moreover, the new model can produce a more robust face shape and texture representation based on their high-level relationships. Experimental results in several applications such as facial super-resolution, face off-angle reconstruction, occlusion removal and facial age estimation have shown the potential of the model in dealing with large variations.

**Acknowledgements** This work was supported in part by the Natural Sciences and Engineering Research Council (NSERC) of Canada.

## References

- Amberg B, Blake A, Vetter T (2009) On compositional image alignment, with an application to active appearance models. In: CVPR, IEEE, pp 1714–1721
- Anderson R, Stenger B, Wan V, Cipolla R (2013) Expressive visual text-to-speech using active appearance models. In: CVPR, IEEE, pp 3382–3389
- Antonakos E, Alabort-i Medina J, Tzimiropoulos G, Zafeiriou S (2014) Hog active appearance models. In: ICIP, IEEE, pp 224–228
- Antonakos E, Alabort-i Medina J, Tzimiropoulos G, Zafeiriou SP (2015) Feature-based lucas–kanade and active appearance models. *IEEE Transactions on Image Processing* 24(9):2617–2632
- Antonakos E, Snape P, Trigeorgis G, Zafeiriou S (2016) Adaptive cascaded regression. In: Image Processing (ICIP), 2016 IEEE International Conference on, IEEE, pp 1649–1653
- Belhumeur PN, Jacobs DW, Kriegman D, Kumar N (2011) Localizing parts of faces using a consensus of exemplars. In: CVPR, IEEE, pp 545–552
- Burgos-Artizzu XP, Perona P, Dollár P (2013) Robust face landmark estimation under occlusion. In: ICCV, IEEE, pp 1513–1520
- Chen K, Gong S, Xiang T, Loy C (2013) Cumulative attribute space for age and crowd density estimation. In: CVPR, pp 2467–2474
- Cootes TF, Taylor CJ (2006) An algorithm for tuning an active appearance model to new data. In: BMVC, pp 919–928
- Cootes TF, Edwards GJ, Taylor CJ (1998) Interpreting Face Images using Active Appearance Models. In: FG, pp 300–305
- Cootes TF, Edwards GJ, Taylor CJ (2001) Active appearance models. *TPAMI* 23(6):681–685
- Ding C, Tao D (2015) Robust face recognition via multimodal deep face representation. *Multimedia, IEEE Transactions on* 17(11):2049–2058
- Dong C, Loy CC, He K, Tang X (2014) Learning a deep convolutional network for image super-resolution. In: ECCV, Springer, pp 184–199
- Donner R, Reiter M, Langs G, Peloschek P, Bischof H (2006) Fast active appearance model search using canonical correlation analysis. *TPAMI* 28(10):1690
- Duong CN, Quach KG, Luu K, Le HB, Jr KR (2011) Fine tuning age-estimation with global and local facial features. In: Intl. Conf. on Acoustics, Speech and Signal Processing (ICASSP), IEEE
- Duong CN, Luu K, Gia Quach K, Bui TD (2015) Beyond principal components: Deep boltzmann machines for face modeling. In: CVPR, pp 4786 – 4794
- Edwards GJ, Cootes TF, Taylor CJ (1998) Face recognition using active appearance models. In: ECCV, Springer, pp 581–595
- Eslami SA, Heess N, Williams CK, Winn J (2014) The shape boltzmann machine: a strong model of object shape. *IJCV* 107(2):155–176
- Ferrari C, Lisanti G, Berretti S, Del Bimbo A (2016) Effective 3d based frontalization for unconstrained face recognition. In: Pattern Recognition (ICPR), 2016 23rd International Conference on, IEEE, pp 1047–1052
- Fu Y, Huang TS (2008) Human age estimation with regression on discriminative aging manifold. *Multimedia, IEEE Transactions on* 10(4):578–584
- Gao S, Zhang Y, Jia K, Lu J, Zhang Y (2015) Single sample face recognition via learning deep supervised autoencoders. *TIFS* 10(10):2108–2118
- Ge Y, Yang D, Lu J, Li B, Zhang X (2013) Active appearance models using statistical characteristics of gabor based texture representation. *JVCIR* 24(5):627–634
- Gross R, Matthews I, Baker S (2005) Generic vs. person specific active appearance models. *Image and Vision Computing* 23(12):1080–1093
- Haase D, Rodner E, Denzler J (2014) Instance-weighted transfer learning of active appearance models. In: CVPR, IEEE, pp 1426–1433
- Hassner T, Harel S, Paz E, Enbar R (2015) Effective face frontalization in unconstrained images. In: CVPR, pp 4295 – 4304
- Hou X, Li SZ, Zhang H, Cheng Q (2001) Direct appearance models. In: CVPR, IEEE, vol 1, pp I–828 – I–833
- Huang GB, Lee H, Learned-Miller E (2012) Learning hierarchical representations for face verification with convolutional deep belief networks. In: CVPR, IEEE, pp 2518–2525

- Huiskes MJ, Thomee B, Lew MS (2010) New trends and ideas in visual concept detection: the mir flickr retrieval evaluation initiative. In: ICMR, ACM, pp 527–536
- Jeni LA, Cohn JF (2016) Person-independent 3d gaze estimation using face frontalization. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, pp 87–95
- Kan M, Shan S, Chang H, Chen X (2014) Stacked progressive auto-encoders (spae) for face recognition across poses. In: CVPR, pp 1883–1890
- Le V, Brandt J, Lin Z, Bourdev L, Huang TS (2012) Interactive facial feature localization. In: ECCV, Springer, pp 679–692
- Li C, Liu Q, Liu J, Lu H (2012) Learning ordinal discriminative features for age estimation. In: CVPR, IEEE, pp 2570–2577
- Li C, Zhou K, Lin S (2014) Intrinsic face image decomposition with human face priors. In: ECCV, Springer, pp 218–233
- Liu L, Xiong C, Zhang H, Niu Z, Wang M, Yan S (2016) Deep aging face verification with large gaps. *Multimedia, IEEE Transactions on* 18(1):64–75
- Luu K, Ricanek K, Bui TD, Suen CY (2009) Age estimation using active appearance models and support vector machine regression. In: BTAS, IEEE, pp 1–5
- Luu K, Bui TD, Suen CY, Ricanek K (2010) Spectral regression based age determination. In: Computer Vision and Pattern Recognition Workshops (CVPRW), IEEE
- Luu K, Bui TD, Suen CY (2011a) Kernel spectral regression of perceived age from hybrid facial features. In: International Conference on Automatic Face and Gesture Recognition and Workshops (FG), IEEE
- Luu K, Keshav Seshadri MS, Bui TD, Suen CY (2011b) Contourlet appearance model for facial age estimation. In: International Joint Conference on Biometrics (IJCB), IEEE
- Martinez A, Benavente R (1998) The ar face database. *Report technique* 24
- Matthews I, Baker S (2004) Active appearance models revisited. *IJCV* 60(2):135–164
- Alabort-i Medina J, Zafeiriou S (2014) Bayesian active appearance models. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 3438–3445
- Alabort-i Medina J, Zafeiriou S (2015) Unifying holistic and parts-based deformable model fitting. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 3679–3688
- Alabort-i Medina J, Zafeiriou S (2017) A unified framework for compositional fitting of active appearance models. *International Journal of Computer Vision* 121(1):26–64
- Alabort-i Medina J, Antonakos E, Booth J, Snape P, Zafeiriou S (2014) Menpo: A comprehensive platform for parametric image alignment and visual deformable models. In: Proceedings of the 22nd ACM international conference on Multimedia, ACM, pp 679–682
- Joan Alabort-i Medina SZ (2014) Bayesian active appearance models. In: CVPR, IEEE, pp 3438–3445
- Mollahosseini A, Mahoor MH (2013) Bidirectional warping of active appearance model. In: CVPRW, IEEE, pp 875–880
- Navarathna R, Sridharan S, Lucey S (2011) Fourier active appearance models. In: ICCV, IEEE, pp 1919–1926
- Ngiam J, Khosla A, Kim M, Nam J, Lee H, Ng AY (2011) Multimodal deep learning. In: ICML, pp 689–696
- Papandreou G, Maragos P (2008) Adaptive and constrained algorithms for inverse compositional active appearance model fitting. In: CVPR, IEEE, pp 1–8
- Pizarro D, Peyras J, Bartoli A (2008) Light-invariant fitting of active appearance models. In: CVPR, IEEE, pp 1–6
- Sagonas C, Tzimiropoulos G, Zafeiriou S, Pantic M (2013) A semi-automatic methodology for facial landmark annotation. In: CVPRW, IEEE, pp 896–903
- Sagonas C, Panagakis Y, Zafeiriou S, Pantic M (2015) Robust statistical face frontalization. In: Proceedings of the IEEE International Conference on Computer Vision, pp 3871–3879
- Salakhutdinov R, Hinton GE (2009) Deep boltzmann machines. In: Intl. Conf. on Artificial Intell. and Statistics, pp 448–455
- Salakhutdinov RR (2009) Learning in markov random fields using tempered transitions. In: NIPS, pp 1598–1606
- Saragih J, Goecke R (2007) A nonlinear discriminative approach to aam fitting. In: ICCV, IEEE, pp 1–8
- Srivastava N, Salakhutdinov R (2012) Multimodal learning with deep boltzmann machines. In: NIPS, pp 2222–2230
- Sun Y, Wang X, Tang X (2013) Deep convolutional network cascade for facial point detection. In: CVPR, pp 3476–3483
- Sun Y, Wang X, Tang X (2014) Deep learning face representation from predicting 10,000 classes. In: CVPR, pp 1891–1898
- Sung J, Kim D (2008) Pose-robust facial expression recognition using view-based 2D + 3D AAM. *TSMC* 38(4):852–866
- Taigman Y, Yang M, Ranzato M, Wolf L (2014) Deepface: Closing the gap to human-level performance in face verification. In: CVPR, pp 1701–1708
- Tang Y, Salakhutdinov R, Hinton G (2012a) Deep lambertian networks. In: ICML
- Tang Y, Salakhutdinov R, Hinton G (2012b) Robust boltzmann machines for recognition and denoising. In: CVPR, IEEE, pp 2264–2271
- Taylor GW, Sigal L, Fleet DJ, Hinton GE (2010) Dynamical binary latent variable models for 3d human pose tracking. In: CVPR, IEEE, pp 631–638

- Tzimiropoulos G, Pantic M (2013) Optimization problems for fast aam fitting in-the-wild. In: ICCV, IEEE, pp 593–600
- Tzimiropoulos G, Pantic M (2017) Fast algorithms for fitting active appearance models to unconstrained images. *International journal of computer vision* 122(1):17–33
- Van Der Maaten L, Hendriks E (2010) Capturing appearance variation in active appearance models. In: CVPRW, IEEE, pp 34–41
- Wang B, Feng X, Gong L, Feng H, Hwang W, Han JJ (2015a) Robust pose normalization for face recognition under varying views. In: Image Processing (ICIP), 2015 IEEE International Conference on, IEEE, pp 1648–1652
- Wang X, Guo R, Kambhamettu C (2015b) Deeply-learned feature for age estimation. In: WACV, IEEE, pp 534–541
- Wang Z, Bovik AC (2009) Mean squared error: love it or leave it? a new look at signal fidelity measures. *Signal Processing Magazine, IEEE* 26(1):98–117
- Wu Y, Wang Z, Ji Q (2013) Facial feature tracking under varying facial expressions and face poses based on restricted boltzmann machines. In: CVPR, IEEE, pp 3452–3459
- Xing J, Niu Z, Huang J, Hu W, Yan S (2014) Towards multi-view and partially-occluded face alignment. In: CVPR, pp 1829–1836
- Yang CY, Liu S, Yang MH (2013) Structured face hallucination. In: CVPR, IEEE, pp 1099–1106
- Yang J, Wright J, Huang TS, Ma Y (2010) Image super-resolution via sparse representation. *TIP* 19(11):2861–2873
- Yildirim I, Kulkarni TD, Freiwald WA, Tenenbaum JB (2015) Efficient analysis-by-synthesis in vision: A computational framework, behavioral tests, and comparison with neural representations. In: CogSci
- Zhai H, Liu C, Dong H, Ji Y, Guo Y, Gong S (2015) Face verification across aging based on deep convolutional networks and local binary patterns. In: IScIDE, Springer, pp 341–350
- Zhang K, Zhang Z, Li Z, Qiao Y (2016a) Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Signal Processing Letters* 23(10):1499–1503
- Zhang Z, Luo P, Loy CC, Tang X (2016b) Learning deep representation for face alignment with auxiliary attributes. *IEEE transactions on pattern analysis and machine intelligence* 38(5):918–930
- Zhu C, Zheng Y, Luu K, Savvides M (2017) Cms-rnn: contextual multi-scale region-based cnn for unconstrained face detection. In: Deep Learning for Biometrics, Springer, pp 57–79
- Zhu J, Hoi SC, Lyu MR (2006) Real-time non-rigid shape recovery via active appearance models for augmented reality. In: ECCV, Springer, pp 186–197
- Zhu Z, Luo P, Wang X, Tang X (2013) Deep learning identity-preserving face space. In: CVPR, pp 113–120
- Zhu Z, Luo P, Wang X, Tang X (2014) Multi-view perceptron: a deep model for learning face identity and view representations. In: NIPS, pp 217–225