# DeMeshNet: Blind Face Inpainting for Deep MeshFace Verification

Shu Zhang    Ran He    Tieniu Tan

National Laboratory of Pattern Recognition, CASIA

Center for Research on Intelligent Perception and Computing, CASIA

{shu.zhang,rhe,tnt}@nlpr.ia.ac.cn

## Abstract

*MeshFace photos have been widely used in many Chinese business organizations to protect ID face photos from being misused. The occlusions incurred by random meshes severely degenerate the performance of face verification systems, which raises the MeshFace verification problem between MeshFace and daily photos. Previous methods cast this problem as a typical low-level vision problem, i.e. blind inpainting. They recover perceptually pleasing clear ID photos from MeshFaces by enforcing pixel level similarity between the recovered ID images and the ground-truth clear ID images and then perform face verification on them.*

*Essentially, face verification is conducted on a compact feature space rather than the image pixel space. Therefore, this paper argues that pixel level similarity and feature level similarity jointly offer the key to improve the verification performance. Based on this insight, we offer a novel feature oriented blind face inpainting framework. Specifically, we implement this by establishing a novel DeMeshNet, which consists of three parts. The first part addresses blind inpainting of the MeshFaces by implicitly exploiting extra supervision from the occlusion position to enforce pixel level similarity. The second part explicitly enforces a feature level similarity in the compact feature space, which can explore informative supervision from the feature space to produce better inpainting results for verification. The last part copes with face alignment within the net via a customized spatial transformer module when extracting deep facial features. All the three parts are implemented within an end-to-end network that facilitates efficient optimization. Extensive experiments on two MeshFace datasets demonstrate the effectiveness of the proposed DeMeshNet as well as the insight of this paper.*

## 1. Introduction

Benefitting from recent advancements in deep representation learning, there have been remarkable improvements in deep face recognition (verification in particu-
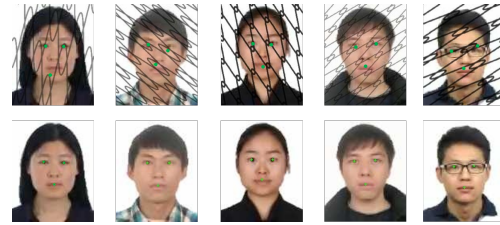


Figure 1. MeshFaces (first row) refer to the ID photos corrupted by randomly generated mesh-like lines or watermarks. The corruptions significantly degenerate the performance of facial landmark detection and facial feature extraction, thus leading to poor verification accuracy.

lar) [28, 29, 31]. In real life applications, face verification between ID photos and daily photos (FVBID) [37] is gaining traction because it uses a face image from an ID photo as gallery and thus does not require the probe to be registred in advance.

When FVBID is applied to real-world scenarios, such as automated custom control and VIP recognition in commercial banks, the ID photo in an identity card may potentially be misused or illegally distributed. Therefore, ID photos are often deliberately corrupted by mesh-like lines or watermarks for privacy protection when used by some business organizations, *e.g.* banks and hotels. For convenience, we denote this type of corrupted ID photo as MeshFace. As shown in Fig. 1, MeshFaces incur catastrophic influence to face recognition systems [2, 13]. Directly verifying Mesh-Faces against daily photos leads to very poor accuracy [36]. Therefore, these corruptions raise a novel and challenging problem called MeshFace verification which deals with face verification between MeshFaces and daily photos.

Some efforts have been made to address this challenging problem. Zhang et al. [36] propose a multi-task residual learning CNN for this problem. They propose to learn a non-linear transformation with SRCNN [7] based architecture to recover clear ID photos from MeshFaces. Then, the recovered clear ID photos are used for face verification. They treat the recovery of clear ID photos from Mesh-Faces as *blind face inpainting* because the position of cor-
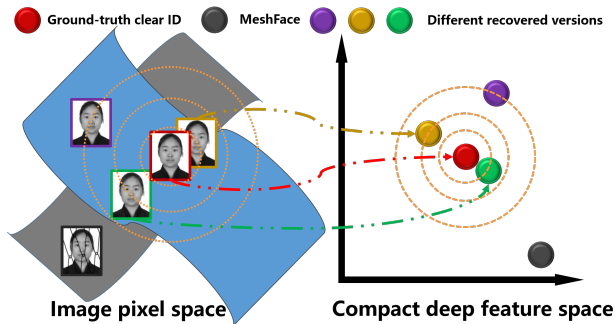
Figure 2. Recovered ID with higher PSNR may be further away from the ground-truth clear ID in the compact deep feature space. For instance, the green and yellow sample. Best viewed in color.

ruptions is unknown during testing phase. In a related vein of research, many contemporary works [22, 25, 26, 34] have shown that CNN is very effective in solving hole filling (non-blind inpainting) problems [4] because this data-driven learning method can exploit the structure of natural images to predict occluded parts.

Improved verification performance is observed after blind face inpainting in [36] because using an occlusion free image will greatly improve the accuracy of face detection and alignment. However, in their work, the performance gap between using their inpainted ID and the clear ID is still very large. On one hand, this is because SRCNN is less powerful in modeling corruption distributions and recovering the exact image content. As a result, the difference in image content between the recovered ID photos and the ground-truth clear ID photos are too large (as illustrated by the red and the purple sample in Fig.2).

On the other hand, existing works often assume that using recovered ID photos with higher PSNR are more likely to achieve better verification performance [36]. But in fact, face similarity is compared in a compact feature space rather the image pixel space. Moreover, it has been shown that CNN can be potentially 'fooled' by adding even a tiny amount of noise to the original input [9, 23, 30]. That is, when facial features are extracted by a CNN, two perceptually indistinguishable face images (*e.g.* the ground-truth clear ID and its recovered version) may still have very large feature level differences (as shown by the green and yellow sample in Fig.2). It is a common belief that a large intra-class feature distance will generally deteriorate face verification performance. Therefore, this suggests that treating blind face inpainting as a typical low-level vision problem by only enforcing the pixel level similarity can hardly guarantee an improvement in verification performance.

To address the aforementioned problems, we present DeMeshNet to take verification performance into account when dealing with the blind face inpainting problem. DeMeshNet is trained on a large scale dataset of Mesh-Face/clear ID photo pairs to learn a non-linear transformation to recover clear ID photos from MeshFaces. Note that DeMeshNet aims to improve the MeshFace verification performance rather than simply to recover perceptually pleasing clear ID photos. Therefore, we refer to DeMeshNet as a feature oriented blind face inpainting framework. We briefly introduce each part of the DeMeshNet and our contributions in the following paragraphs.

In the first part, we enforce the *pixel level similarity* to explore the structure of face images so as to recover uncorrupted ID photos from MeshFaces. Specifically, we propose to adopt a fully convolutional network (FCN) [21] with a weighted Euclidean loss to minimize the pixel differences of the ground-truth clear ID/recovered ID photo pairs. Extra supervision from corruption positions is further exploited to accurately model the corruption distribution.

In the second part, we enforce a *feature level similarity* between the ground-truth clear ID photo and the recovered ID photo pairs in a compact deep feature space. This feature space is spanned by a pre-trained CNN, and distance in it directly corresponds to a measure of face similarity. This part will force the inpainted image to have a smaller distance to the ground-truth clear ID in the deep feature space, which will in turn facilitate accurate verification. Moreover, we propose to measure the feature level similarity with the reverse Huber loss function so that the feature level similarity can be more efficiently optimized when the differences are very small.

In the third part, we employ a *customized spatial transformer module* [15] to align and crop the face region for accurate feature extraction within the network. It is essential to take alignment into account because the MeshFace which is the input of DeMeshNet and the aligned face which is the input of the feature extraction sub-net are different in sizes, scales and orientations (as shown in Fig. 3).

All the three parts are implemented within an end-to-end network that facilitates efficient optimization with gradient back propagation. Extensive experimental results on two MeshFace datasets demonstrate that DeMeshNet achieves the best verification accuracy and outperforms previous work by a large margin. Furthermore, we thoroughly evaluate different configurations of DeMeshNet to gain insight into the factors for such significant improvements.

## 2. Approach

### 2.1. Overview

In this section, we present an overview of the proposed DeMeshNet. We cast the proposed feature oriented blind face inpainting problem as a dense regression problem, which aims to regress a perceptually pleasing and verification favorable clear ID photo from a MeshFace $X$. For convenience, we refer to the ground-truth clear ID photo as
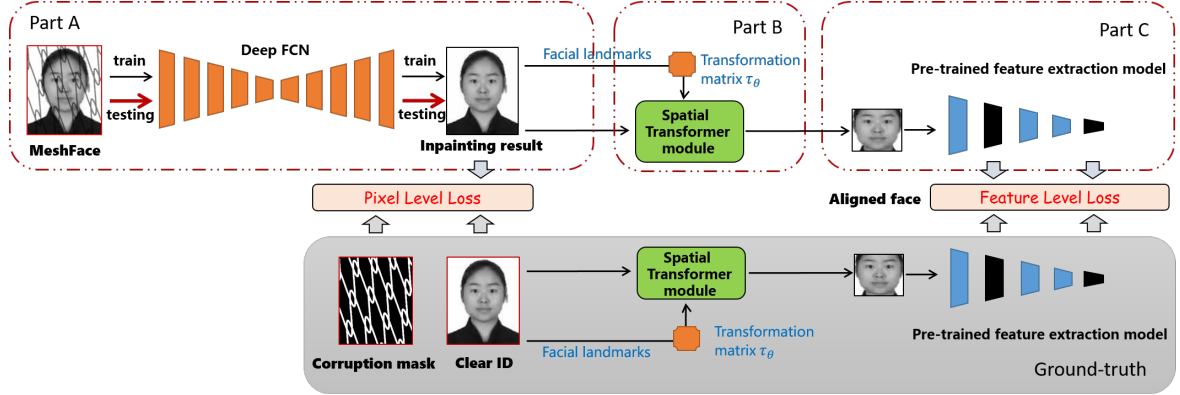
Figure 3. Conceptual diagram of DeMeshNet. Black solid lines stand for the training phase and red solid lines represent the testing phase. Note that the corruption mask and clear ID are only used in the training phase to provide ground-truth for the pixel and feature level loss. Feature extraction model is pre-trained and has fixed parameters during training. In the testing phase, DeMeshNet takes a MeshFace as input and produce an inpainted photo with the learned deep FCN.

the *target*, termed as $Y$ and the recovered ID photo from our blind face inpainting model $\psi$ as the *prediction*, termed as $\psi(X)$. The prediction will be used for face verification against clear daily photos.

We model the highly non-linear function for dense regression as a FCN as illustrated in part A of Fig. 3. Part B shows the customized spatial transformer module. The following part is a pre-trained CNN which is utilized to compute the feature representation of the aligned face region. It should be noted that parameters in both the spatial transformer module and the pre-trained CNN are fixed during training. The learnable parameters in the FCN are optimized through minimizing a unified loss function that jointly models the pixel and feature level similarities between the prediction and the target pairs. No identity information is needed to train such a blind inpainting network.

The pixel level loss helps to obtain perceptually pleasing inpainting results and serves as a means to capture the distribution difference between actual face texture and mesh-like corruptions. And the feature level loss explores supervision in a compact feature space to provide regularizations to the network training. Thus, the network's prediction will not only have similar appearance but also have similar feature representation to that of the target. Specifically, we develop a weighted Euclidean loss to model the pixel level similarity and employ a reverse Huber function [18] to characterize the feature level loss on the spatial transformed face region. Combining the pixel level loss and the spatial transformed feature level loss, we define the unified loss function for DeMeshNet as follows:

$$
\begin{aligned}
L_i &= l_{pixel} + l_{feature} = \\
&||\psi(X_i) - Y_i||_F^2 + \lambda_1||M_i \odot (\psi(X_i) - Y_i)||_F^2 + \\
&\lambda_2 \sum_{j=1}^{2} RH(\phi_j(\psi(ST(X_i))) - \phi_j(ST(Y_i)))
\end{aligned}
\tag{1}
$$

where $ST$ denotes the spatial transformation implementation that samples an aligned $128 \times 128$ face region from the original input solely based on the positions of facial s, $RH$ is the reverse Huber function, $\lambda_1$ and $\lambda_2$ are the balance parameters which are empirically set to 1 throughout the paper. We postpone explanations of other symbols to later sections when we meet them.

For simplicity, we omit the regularization term on the parameters of FCN (weight decay) in Equation 1, which is used to reduce overfitting when optimizing our network. The objective function can be efficiently optimized by gradient back propagation in an end-to-end manner. We will elaborate each of the three parts in the following subsections.

## 2.2. Pixel Level Regression Network

### 2.2.1 Pixel Level Loss

Blind face inpainting is naturally characterized as a pixel-wise regression problem. In the first part of DeMeshNet, it learns a highly non-linear transformation by optimizing a well-designed pixel level loss. For blind face inpainting, although positions of the corruption are not provided during the testing phase, we can still make use of this information when training DeMeshNet. Specifically, we propose to implicitly exploit this extra supervision by introducing a weighted Euclidean loss function as below:

$$
l_{pixel} = ||\psi(X_i) - Y_i||_F^2 + \lambda||M_i \odot (\psi(X_i) - Y_i)||_F^2 \tag{2}
$$

where $X_i$, $Y_i$ and $\psi(X_i)$ are a corrupted input, target and prediction, respectively. $M_i$ is the binary mask with a value of 1 indicating the pixel is corrupted and a value of 0 otherwise. $\odot$ is the element-wise product operation. Therefore, the second term in the loss function only measures the Euclidean loss on corrupted areas, emphasizing

losses on those areas with a weighting parameter $\lambda$. This loss function helps the network to learn the distributions of corrupted pixels better by exploiting extra supervision from the corruption positions. Experimental results demonstrate that it also helps to generate predictions with higher PSNR.

### 2.2.2 Network Architecture

Since FCN has achieved outstanding performance in dense prediction tasks like depth prediction [8] and semantic segmentation [21], it motivates us to use FCN as the non-linear function to improve the blind face inpainting performance.

The main difference between FCN and the architecture in [36] is the introduction of down-sampling and up-sampling layers in FCN. This simple adjustment has enabled FCN to admit much deeper layers and to expand the receptive fields with the same amount of computational cost. The expanded receptive fields are critical to blind inpainting as it can enclose more contextual information for identifying the corrupted areas.

In this work, we use the network architecture of SegNet as proposed in [1]. It is feasible to adopt other architectures such as ResidualNet [12] and Deconvolution Network [24], but this is beyond the scope of this paper. The input and output to the network are MeshFaces and clear ID photos respectively. Gray scale images of size $220 \times 178$ are used for input and output throughout the paper.

## 2.3. Feature Level Regression Network

### 2.3.1 Feature Level Loss

As aforementioned in Section 1, images with very small Euclidean distance may have large feature distance. This problem is raised in [30] and has been shown to severely deteriorate classification performance because of the enlarged intra-class feature distance. In fact, the enlarged intra-class feature distance can also influence the verification task at hand. To improve the verification performance, it won't be enough to only enforce pixel level similarity. Therefore, in the second part, we explicitly enforce the target and the prediction to have a small distance in the compact feature space computed by the pre-trained CNN $\phi$.

Let $\phi_j(\psi(X_i))$ be the activations of the $j$th layer of the pre-trained face model $\phi$. We intend to improve the verification performance by minimizing the residual $r = \phi(\psi(x_i)) - \phi(y_i)$ at each position of an image. To efficiently back-propagate the errors when the residual is very small, we employ the reverse Huber loss function [18] to measure the feature level difference, its formulation is:

$$RH(r) = \begin{cases} |r| & |r| > c \\ \frac{r^2 + c^2}{2c} & |r| \le c \end{cases} \qquad (3)$$

The reverse Huber loss is equivalent to L1 norm when the residual $r$ is in the interval of $[-c, c]$ and equals to a
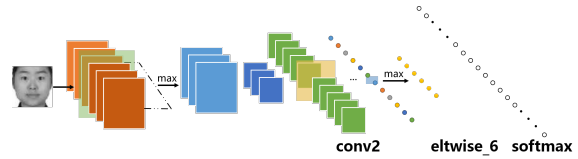


Figure 4. Schematic architecture for face feature extraction. It uses Max-Feature-Map nonlinearities instead of ReLU.

transformed L2 norm otherwise. This loss experimentally works better than L2 norm. Note that when $c$ is smaller than 1, the derivative of the L1 norm is greater than that of the L2 norm, which will speed up error back-propagation when the residual is very tiny. Like in [18], we use a dynamic threshold for $c$. That is, in each batch-minimization step, $c$ is set to be at $20\%$ of maximum residual in that batch.

Since our objective is to improve verification performance, we impose the reverse Huber loss on the output of $eltwise\_6$ in the pre-trained model $\phi$, which is a 256-dim feature vector used for similarity comparison. Drawing inspirations from [17], where feature loss is used for super-resolution, we also impose the reverse Huber loss on the early layers ($conv2$ in our case) of the pre-trained face model $\phi$ to include deeper supervision [20]. Therefore, our final loss function for feature level regression is:

$$l_{feature} = \sum_{j=1}^{2} RH(\phi_j(\psi(X_i)) - \phi_j(Y_i)) \qquad (4)$$

Feature loss has recently been considered in the literature of super-resolution and sketch inversion for better visual performance [10, 17, 19]. But it should be noted that in this paper, the feature level loss is proposed from a totally different perspective. Instead of pursing a perceptually pleasing image transformation results, we want to deal with the large intra-class feature distance problem and improve verification performance.

### 2.3.2 Pre-trained CNN for Face Feature Extraction

Both training DeMeshNet and evaluate the verification performance need a pre-trained CNN to compute the facial features for face images. We use the architecture proposed in [33] for facial feature extraction because it is computationally efficient in both time and space. Fig. 4 briefly illustrates its architecture, which uses Max-Feature-Map nonlinearities instead of ReLU and thus can generate dense features at its output. The network takes aligned gray-scale face images of size $128 \times 128$ as input and returns a 256-dim feature (output of $eltwise\_6$). Alignment is conducted by transforming two facial landmarks (i.e., centers of two eyes) to $(32, 32)$ and $(96, 32)$ with a similarity transformation.
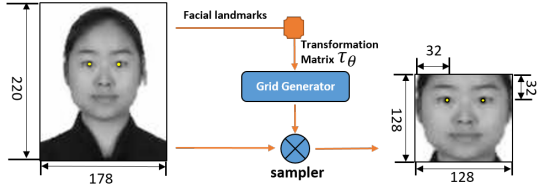
Figure 5. The customized spatial transformer module used in our model. It uses the locations of facial landmarks to calculate a transformation matrix for face alignment.

We train a base model on the *purified* MS-Celeb-1M dataset [11] (the original data is very noisy, we purify them before use). One single base model achieves a verification accuracy of $98.80\%$ on the LFW [14] benchmark, which is very competitive among already published works [6, 32]. We further finetune the base model with triple-let loss on a large-scale ID-daily photo dataset collected from the web and use the finetuned model $\phi$ to extract compact facial features in DeMeshNet. Note that this model's parameters should stay fixed during the training of DeMeshNet because we only use it for feature computation and no identity related supervision is adopted to finetune $\phi$ while learning the blind inpainting FCN.

## 2.4. Spatial Transformer For Face Alignment

Face alignment is essential for feature extraction. It helps to improve verification performance by providing a normalized input. The pre-trained model $\phi$ admits $128 \times 128$ aligned face region to compute the 256-dim feature. But DeMeshNet takes $220 \times 178$ un-aligned MeshFace as input and outputs a prediction with the same size. Therefore, in order to compute the 256-dim feature in the fully connected layer $eltwise\_6$, we must implement face alignment within the network.

Moreover, unlike the image classification models that are trained with multi-scale natural images from ImageNet [27], the pre-trained facial feature extraction model $\phi$ only takes single-scale, well-aligned face regions for training. This means that even we only compute the features from conv layers like in [10, 17, 19], we will still need to align the MeshFaces first to acquire an accurate feature representation.

We incorporate a customized spatial transformer module [15] between the pixel level regression sub-net and the feature level regression sub-net to sample an aligned $128 \times 128$ face region from the $220 \times 178$ prediction according to the facial landmarks. This procedure is illustrated in Fig. 3 and detailed in Fig. 5.

The spatial transformer module comprises of a localization network, a grid generator and a sampler. Since the similarity transformation $\tau_\theta$ for face alignment is uniquely determined by the coordinates of two eye centers, we do not

need to learn it through the localization network as in [15]. $\tau_\theta$ is parameterized by $\theta = [a, b, 1; -b, a, 1]$ and can be determined with the following equation:

$$
\begin{pmatrix} x_l & x_r \\ y_l & y_r \\ 1 & 1 \end{pmatrix} = \begin{bmatrix} a & b & 1 \\ -b & a & 1 \end{bmatrix} \begin{pmatrix} -0.5 & 0.5 \\ -0.5 & -0.5 \\ 1 & 1 \end{pmatrix} \quad (5)
$$

where $(x_l, y_l)$, $(x_r, y_l)$ and $(-0.5, -0.5)$, $(0.5, -0.5)$ are the *normalized coordinates* (normalized to $[-1, 1]$) of two eye centers in the original image and the aligned face image respectively.

In the forward pass, a sampling grid is firstly determined with the given $\tau_\theta$. A sampling grid is a set of points with continuous coordinates. Sampling an input image according to this sampling grid will generate a transformed output. By defining the output pixels to lie on a regular grid $G = G_i$ of pixels $G_i = (x_i^t, y_i^t)$, the sampling grid $\tau_\theta(G)$ is given by the point-wise transformation:

$$
\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \tau_\theta(G) = \begin{bmatrix} a & b & 1 \\ -b & a & 1 \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix} \quad (6)
$$

where $(x_i^s, y_i^s)$ are the source coordinates in the input feature map, and $(x_i^t, y_i^t)$ are the target coordinates of the regular grid. Since the coordinates in the sampling grid are continuous numbers, a bilinear kernel is applied to those positions to produce the corresponding pixel values in the output:

$$
Q = max(0, 1 - |x_i^s - m|)max(0, 1 - |y_i^s - n|) \quad (7)
$$

$$
P_{(x_i^t, y_i^t)} = \sum_{n}^{H} \sum_{m}^{W} P_{(x_n^s, y_m^s)} Q \quad (8)
$$

where $H$ and $W$ are the height and width of the input image respectively and $P_{(x_i^t, y_i^t)}$ represents the pixel value of $(x_i^t, y_i^t)$. To allow the feature loss defined in the last subsection to be back-propagated from the output of the spatial transformer module to the input image, we give the gradients with respect to the input image as follows:

$$
\frac{\partial(P(x_i^t, y_i^t))}{\partial(P(x_m^s, y_n^s))} = \sum_{n}^{H} \sum_{m}^{W} Q \quad (9)
$$

The gradients of the feature loss defined earlier can be easily flowed back to the input image using chain rule with Equation 9. Note that the gradients with respect to the sampling grid coordinates $(x_i^s, y_i^s)$ are not derived, because the transformation parameters are not learned in the customized spatial transformer module.

## 3. Experiments

In this section, we experimentally evaluate the proposed framework. We begin by introducing the datasets for training and testing. Then we specify the baseline methods and
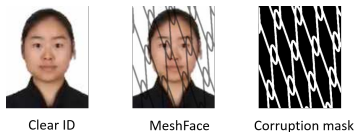
Figure 6. An image triplet sample from the training dataset [36].



Figure 7. Sample image pairs from SV1000. Note that this dataset is very hard as the variations in illumination, pose and hair style are very significant.



Figure 8. Sample image pairs from SYN500. Images of some Chinese celebrities are collected from the Internet.

implementation details. At last, we present detailed algorithmic evaluation, as well as comparison with other methods.

### 3.1. Datasets

All the compared models are trained on the dataset as in [36] that contains over $500,000$ data triplets of $11,648$ individuals. Each data triplet consists of a MeshFace, its clear version and a corruption mask (as illustrated in Fig. 6). Facial landmarks (two eye centers) are detected using Intraface [35] to aid the spatial transformer module. Data triplets of 400 individuals are sampled for validation and testing (200 each) and all the other individuals are used for training.

Besides the SYN500 used in [36], we collect another dataset with 1000 MeshFace/daily photo pairs from 1000 individuals named SV1000 to evaluate the MeshFace verification performance. Daily photos in SV1000 are captured under surveillance cameras. As shown in Fig. 7, this dataset not only contains more individuals but also presents more variations in the daily photo, which makes it more challenging than SYN500 (shown in Fig. 8).

We develop a protocol for evaluation of verification performance on these two datasets. Specifically, face comparison is conducted between all the possible recovered clear ID/daily photo pairs in the compact feature space (spanned by model $\phi$) with cosine distance. For a dataset with $N$ data pairs, $N^2$ comparisons are conducted in total. To exclude influences from metric learning methods, no supervised learning methods, *e.g.* joint bayesian [3], are employed on the extracted features.

### 3.2. Baselines and Implementation Details

Although many algorithms have been proposed for non-blind inpainting, few have been developed to address the blind inpainting problem, even less for the blind face inpainting problem addressed in this paper. We implement the multi-task CNN (MtNet) [36] as a baseline. This method employs architecture that resembles the SRCNN [7] and use

multi-task learning to make use of the information of corruption position in the training phase.

To give a detailed evaluation of each part of the proposed DeMeshNet, we also compare it with various configurations. The compared configurations include FCN with Euclidean pixel level loss (FCNE), FCN with weighted pixel level loss (FCNW) and feature loss FCN without spatial transformer module (FCNF). All three configurations use FCN as the backbone for the blind face inpainting task. Both FCNE and FCNW only adopts the pixel level loss during training, but FCNW introduces an implicit supervision from the corruption mask with a weighted loss in addition to the Euclidean loss. FCNF takes both weighted pixel loss and whole-image feature level loss into consideration. But the feature level differences are only computed at the output of $conv2$, using whole-image as input to the pre-trained feature extraction network $\phi$. Like in [17], we don't implement face alignment within the network.

All the compared models are trained on the training set with photo pairs of size $220 \times 178$, gray scale images are used in all the experiments. For all the compared network structure, training is carried out using Adam [5] with a batch size of 30. The learning rate is set to $10^{-4}$ initially, and decreased by a factor of 10 each $40k$ iterations. The training process takes approximately $160k$ iterations to converge. For the proposed approach, feature level loss is computed at layer $conv2$ and $eltwise\_6$ of the pre-trained face model $\phi$. For FCNF, only $conv2$ is used for computing the feature loss as the computation of $fc$ layer $eltwise\_6$ requires the input to be of size $128 \times 128$. All the MeshFace verification experiments use the pre-trained face model $\phi$ for facial feature extraction. All the experiments are conducted with the Caffe framework [16] on a single GTX Titan X GPU.

### 3.3. Evaluation of Verification Results

In this section, we conduct MeshFace verification experiments on two datasets, i.e., SYN500 and SV1000. Recovered ID photos are used for FVBID according to the aforementioned protocol. We report ROC curves in Fig. 9. TPR@FPR=1% (true positive rate when false positive rate is %1), TPR@FPR=0.1% and TPR@FPR=0.01% are reported in Table 1 for closer inspection. We also present the face verification performances with ground-truth clear ID photos and MeshFaces (denoted as *Clear* and *Corrupted*) for fair comparison.

As expected, when using the MeshFaces for verification,

Table 1. Verification performance on SYN500/SV1000 and inpainting results on the testing set. RMSE illustrates the feature distance between the recovered ID photos and the ground-truth clear ID photos, while PSNR indicates the pixel distance. Note that smaller RMSE consistently indicates better verification performance, but higher PSNR doesn't guarantee that.

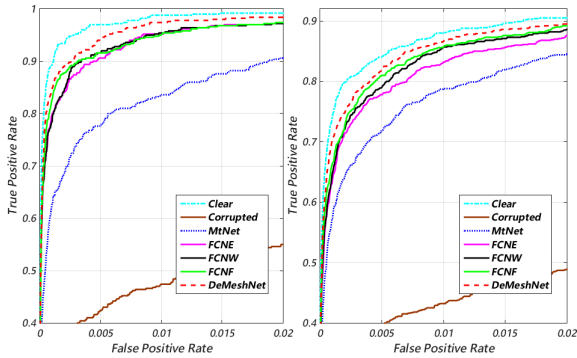| Method | TPR@FPR=1% | TPR@FPR=0.1% | TPR@FPR=0.01% | PSNR | RMSE |
|---|---|---|---|---|---|
| MtNet | 83.60% / 78.80% | 62.80% / 57.50% | 36.80% / 35.40% | 29.89 | 55.47 |
| Clear | **98.80% / 88.10%** | **89.40% / 74.30%** | **67.40% / 53.60%** | - | 0 |
| Corrupted | 47.40% / 43.20% | 33.80% / 28.50% | 18.40% / 18.20% | 20.69 | 112.63 |
| FCNE | 95.20% / 83.20% | 79.80% / 63.90% | 53.80% / 43.00% | 35.11 | 49.52 |
| FCNW | 95.40% / 85.70% | 78.40% / 64.90% | 52.80% / 44.40% | **35.31** | 48.22 |
| FCNF | 95.20% / 85.80% | 82.80% / 66.40% | 54.40% / 46.30% | 25.26 | 38.19 |
| DeMeshNet_E | 96.60% / 86.30% | 83.90% / 70.00% | 54.80% / 46.50% | 29.28 | 35.77 |
| DeMeshNet | **97.40% / 86.70%** | **84.80% / 70.70%** | **55.20% / 47.00%** | 29.16 | **34.57** |



Figure 9. ROC curves for SYN500 and SV1000.

the accuracy suffers a severe drop on both datasets due to large detection and alignment errors. After processing the MeshFace with blind face inpainting models, face verification performance with the recovered ID photos has seen a great improvement. Owning to the deeper FCN architecture and expanded receptive fields, all the FCN based models outperform the baseline model MtNet [36] by a large margin on both datasets.

As shown in Fig. 9, feature loss based models (DeMeshNet, FCNF) perform better than the models that only seek a visually pleasing inpainting results (FCNE, FCNW). This suggests that by enforcing a feature level similarity during training, predictions from DeMeshNet lie in a low-dimensional feature space that is closer to the ground-truth clear ID photos than predictions from pixel-level only networks. This is validated in the next section where we calculate the RMSE (rooted mean square error) between the features of ground-truth clear IDs and recovered IDs.

We further investigate the role of the spatial transformer module in our model by comparing DeMeshNet with FCNF which uses the image of original size ($220 \times 178$) as input to the feature loss component. We find that DeMeshNet consistently performs better than FCNF. This is because FCNF takes the whole image, which is different from the aligned face region in both scale and orientation, for fea-

ture extraction. Unlike the ImageNet [27] models that are trained with multi-scale natural images, the pre-trained face model $\phi$ only takes single-scale, well-aligned face regions for training. The scale and orientation differences have led to the performance degeneration. Therefore, it is significant to take face alignment into account when optimizing the feature level loss as done in DeMeshNet.

It should be noted that for blind inpainting models, there is an upper limit for their verification performance, which is the verification performance with ground-truth clear ID photos. From Table 1, we can observe that the gap between DeMeshNet and clear ID at TPR@FPR=1% is very small (1.4% for both datasets), validating the outstanding performance of DeMeshNet.

### 3.4. Evaluation of Inpainting Results

In this section, we *qualitatively* and *quantitatively* evaluate the inpainting results on the testing set (200 individuals, 10000 photos). Firstly, we qualitatively evaluate the compared models by visual inspection of the inpainting results in Fig. 10. It is observed that MtNet fails to identify and recover some portions of the corruptions in these cases (cherry-picked to illustrate the point). In contrast, FCN based models can handle all the corruption areas very well because they can enclose more contextual information with expanded receptive fields. This demonstrates the improved capacity of FCN over SRCNN based architectures.

Regarding the details of the recovered ID photo, the models trained with only pixel level loss (FCNE, FCNW, MtNet) can better preserve the consistency of pixels and thus provide a smooth and clear photo which is more similar to the ground-truth. But the images recovered with models trained on feature level loss (FCNF in particular) contain many artifacts, making them visually less appealing. This is because the high-level features are robust to pixel level changes in the texture, shape and even color. However, images recovered from DeMeshNet looks much better than the ones from FCNF due to the introduction of the spatial transformer module within DeMeshNet.

Next, we quantitatively evaluate the models by measur-

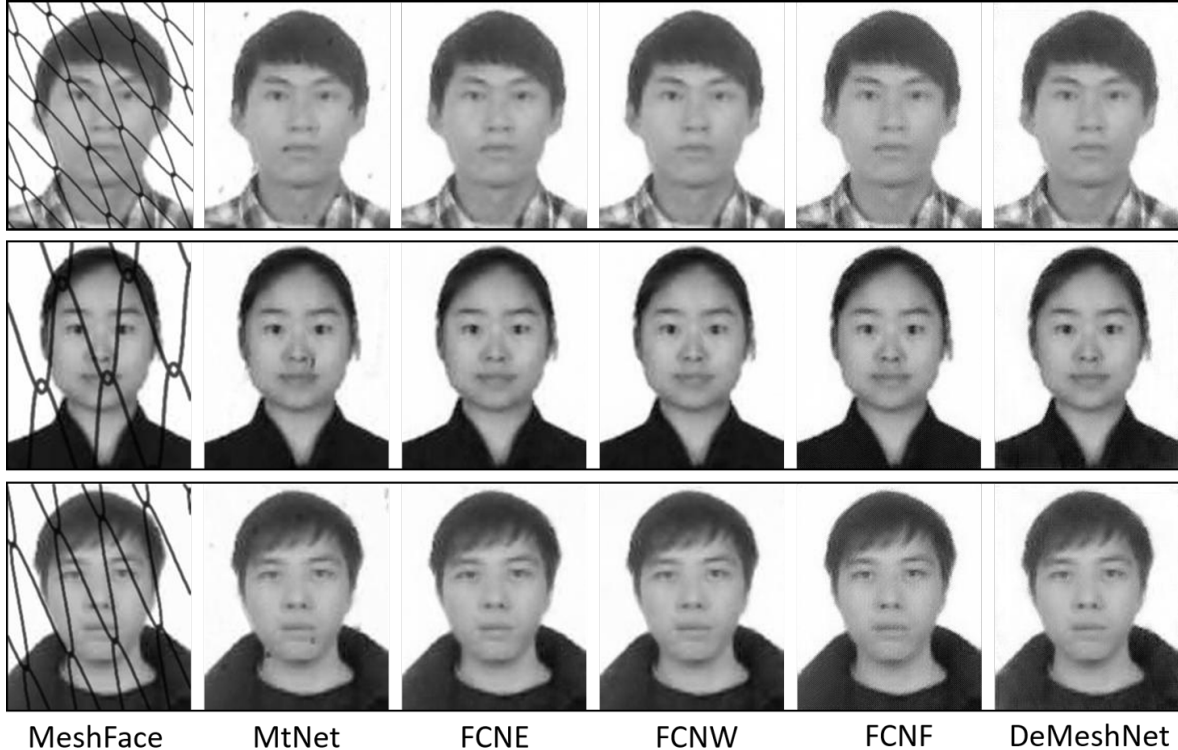| MeshFace | MtNet | FCNE | FCNW | FCNF | DeMeshNet |

Figure 10. Visual inspection of the inpainting results. Although these inpainting results all look very well and are quite similar, they will lead to entirely different verification rates because of the different RMSE in the feature space.

ing both pixel level and feature level ($eltwise\_6$) difference. The ground-truth clear ID photo is used as a baseline. Specifically, average PSNR and Euclidean distance between features (or rooted mean square error, RMSE) are shown in Table 1. Pixel level loss based models yield better PSNR as they explicitly optimize PSNR in their loss function. We also observe that FCNW performs slightly better than FCNF. This implies that exploiting extra supervision from the corruption position in the training phase is beneficial for acquiring visually pleasing inpainting results.

To validate the choice of the reverse Huber loss over the Euclidean loss on the feature level difference, we also implement a DeMeshNet_E that use Euclidean loss for the feature level loss. From Table 1, we can see that RMSE for DeMeshNet_E is slightly larger than DeMeshNet. This indicates that the reverse Huber loss is more effective at minimizing the feature level distance thanks to the imposed L1 norm when residuals are small.

Furthermore, from Table 1 we observe that smaller RMSE often means better verification performance, but higher PSNR doesn't guarantee smaller RMSE. This reveals that the models trained with only pixel level loss does suffer from the influence of the easily fooled nature of CNN [9, 23, 30]. Moreover, visually appealing inpainting results does not necessarily produce better verification re-

sults. By exploring supervision from the deep feature space, DeMeshNet can capture a distribution that is more robust to transformation in $\phi$ and thus provide a stable high-level representation in the compact deep feature space.

## 4. Conclusions

This paper addresses the MeshFace verification problem that verifies corrupted ID photos against clear daily photos. Specifically, we have proposed DeMeshNet that consists of three parts to blindly inpaint the MeshFace before conducting verification. The proposed DeMeshNet distinguishes itself from previous works by explicitly taking verification performance into consideration while recovering a clear ID photo. The training objective of DeMeshNet is motivated by the fact that minimizing pixel level differences alone cannot guarantee a small intra-class feature distance in the compact deep feature space, which is crucial for accurate face verification. By further incorporating a spatial transformer module, DeMeshNet can implement face alignment within the network, resulting in an end-to-end network. For optimizing DeMeshNet, a very well-performed facial feature extraction network has been trained in advance. Experimental results on two MeshFace datasets demonstrate that the proposed DeMeshNet outperforms previous work on verification performance.

# References

[1] V. Badrinarayanan, A. Handa, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for robust semantic pixel-wise labelling. *arXiv preprint arXiv:1505.07293*, 2015. 4

[2] X. P. Burgos-Artizzu, P. Perona, and P. Dollar. Robust face landmark estimation under occlusion. In *Proceedings of the IEEE International Conference on Computer Vision*, December 2013. 1

[3] D. Chen, X. Cao, L. Wang, F. Wen, and J. Sun. Bayesian face revisited: A joint formulation. In *Proceedings of European Conference on Computer Vision*, pages 566–579. Springer, 2012. 6

[4] A. Criminisi, P. Pérez, and K. Toyama. Region filling and object removal by exemplar-based image inpainting. *IEEE Transactions on image processing*, 13(9):1200–1212, 2004. 2

[5] K. Diederik and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6

[6] C. Ding and D. Tao. Robust face recognition via multimodal deep face representation. *IEEE Transactions on Multimedia*, 17(11):2049–2058, 2015. 5

[7] C. Dong, C. C. Loy, K. He, and X. Tang. Learning a deep convolutional network for image super-resolution. In *Proceedings of European Conference on Computer Vision*, 2014. 1, 6

[8] D. Eigen, C. Puhrsch, and R. Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in neural information processing systems*, pages 2366–2374, 2014. 4

[9] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014. 2, 8

[10] Y. Güçlütürk, U. Güçlü, R. van Lier, and M. A. van Gerven. Convolutional sketch inversion. *arXiv preprint arXiv:1606.03073*, 2016. 4, 5

[11] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. *arXiv preprint arXiv:1607.08221*, 2016. 5

[12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June 2016. 4

[13] R. He, W.-S. Zheng, and B.-G. Hu. Maximum correntropy criterion for robust face recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(8):1561–1576, 2011. 1

[14] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical Report 07-49, University of Massachusetts, Amherst, October 2007. 5

[15] M. Jaderberg, K. Simonyan, A. Zisserman, et al. Spatial transformer networks. In *Advances in Neural Information Processing Systems*, pages 2017–2025, 2015. 2, 5

[16] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014. 6

[17] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. *arXiv preprint arXiv:1603.08155*, 2016. 4, 5, 6

[18] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab. Deeper depth prediction with fully convolutional residual networks. *arXiv preprint arXiv:1606.00373*, 2016. 3, 4

[19] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. P. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi. Photo-realistic single image super-resolution using a generative adversarial network. *CoRR*, abs/1609.04802, 2016. 4, 5

[20] C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu. Deeply-supervised nets. *arXiv preprint arXiv:1409.5185*, 2014. 4

[21] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015. 2, 4

[22] X.-J. Mao, C. Shen, and Y.-B. Yang. Image restoration using convolutional auto-encoders with symmetric skip connections. *arXiv preprint arXiv:1606.08921*, 2016. 2

[23] A. Nguyen, J. Yosinski, and J. Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 427–436. IEEE, 2015. 2, 8

[24] H. Noh, S. Hong, and B. Han. Learning deconvolution network for semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, December 2015. 4

[25] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June 2016. 2

[26] J. S. Ren, L. Xu, Q. Yan, and W. Sun. Shepard convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 901–909, 2015. 2

[27] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. 5, 7

[28] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 815–823, 2015. 1

[29] Y. Sun, X. Wang, and X. Tang. Deep learning face representation from predicting 10,000 classes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1891–1898, 2014. 1

[30] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013. 2, 4, 8

[31] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1701–1708, 2014. 1

[32] D. Wang, C. Otto, and A. K. Jain. Face search at scale. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PP(99):1–1, 2016. 5

[33] X. Wu, R. He, and Z. Sun. A lightened cnn for deep face representation. *arXiv preprint arXiv:1511.02683*, 2015. 4

[34] J. Xie, L. Xu, and E. Chen. Image denoising and inpainting with deep neural networks. In *Advances in Neural Information Processing Systems*, pages 341–349, 2012. 2

[35] X. Xiong and F. De la Torre. Supervised descent method and its applications to face alignment. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013. 6

[36] S. Zhang, R. He, Z. Sun, and T. Tan. Multi-task convnet for blind face inpainting with application to face verification. In *Proceddings of the IEEE International Conference on Biometrics*, pages 1–8, June 2016. 1, 2, 4, 6, 7

[37] E. Zhou, Z. Cao, and Q. Yin. Naive-deep face recognition: Touching the limit of lfw benchmark or not? *arXiv preprint arXiv:1501.04690*, 2015. 1