

# Learning Inverse Rendering of Faces from Real-world Videos

Yuda Qiu<sup>†1,2</sup>, Zhangyang Xiong<sup>†1,3</sup>, Kai Han<sup>4</sup>, Zhongyuan Wang<sup>3</sup>,  
Zixiang Xiong<sup>5</sup>, and Xiaoguang Han<sup>\*1,2</sup>

<sup>1</sup> Shenzhen Research Inst. of Big Data

<sup>2</sup> The Chinese University of Hong Kong, Shenzhen

<sup>3</sup> Wuhan University <sup>4</sup>University of Oxford <sup>5</sup>Texas A&M University

**Abstract.** In this paper we examine the problem of inverse rendering of real face images. Existing methods decompose a face image into three components (albedo, normal, and illumination) by supervised training on synthetic face data. However, due to the domain gap between real and synthetic face images, a model trained on synthetic data often does not generalize well to real data. Meanwhile, since no ground truth for any component is available for real images, it is not feasible to conduct supervised learning on real face images. To alleviate this problem, we propose a weakly supervised training approach to train our model on real face videos, based on the assumption of consistency of albedo and normal across different frames, thus bridging the gap between real and synthetic face images. In addition, we introduce a learning framework, called IlluRes-SfSNet, to further extract the residual map to capture the global illumination effects that give the fine details that are largely ignored in existing methods. Our network is trained on both real and synthetic data, benefiting from both. We comprehensively evaluate our methods on various benchmarks, obtaining better inverse rendering results than the state-of-the-art.

## 1 Introduction

Inverse rendering aims at estimating the components of an image in its formation process. An image is often decomposed to three components, namely, albedo (reflectance properties), normal (shape attributes), and illumination [2], [4], [38], [44]. Inverse rendering has important applications in image analysis (e.g., scene segmentation and material recognition) and editing (e.g., photo relighting).

In this paper, we consider inverse rendering of human face images because they belong to the most important class of images in vision and recognition tasks. Promising results have been achieved under constrained scenarios (e.g.,

---

<sup>†</sup> Equal contribution.

\* Corresponding author: hanxiaoguang@cuhk.edu.cn

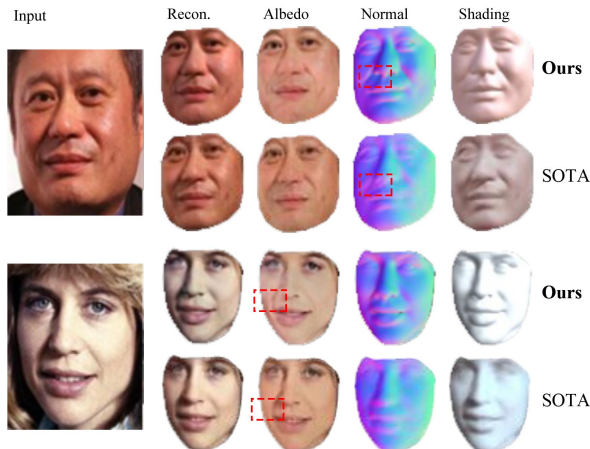


Fig. 1: Decomposing real world faces into albedo, normal, and illumination. The first row of each sample shows our results, and the second one is SfSNet,[34], the state-of-the-art work on inverse face rendering. (Best viewed in PDF with zoom.)

on the well designed lighting stages[11][9][13]). However, inverse rendering of in-the-wild face images remains an open problem due to complex variations in human face appearances, illumination conditions, and shadows. Moreover, the lack of ground-truth decomposition components makes this task highly ill-posed. To tackle this problem, some hand-crafted priors (e.g., [19], [4]) have been introduced for each component to guide the decomposition. Unfortunately, these priors make strong assumptions on face attributes (e.g., strong shape and reflectance priors for Caucasians) that are tailored for specific tasks. Thus they do not generalize well to in-the-wild face images [7].

Owing to the success of deep learning in many tasks, deep convolutional neural networks (CNNs) have recently been employed to address the problem of inverse rendering. For example, [38,34] employed CNNs to deal with this problem by training on synthetic data. After training on synthetic face image datasets, the network is applied on real face images to obtain different image modeling components. However, training only on synthetic data does not generalize well to real data since real face images contain much richer facial variations that cannot be captured by synthetic data (e.g., faces wearing glasses and/or makeup, faces with beard, etc). It is thus very desirable to make use of real data for training. To this end, [34] proposed to train on real face images with pseudo labels. However, results reported in [34] are far from satisfactory because the pseudo labels obtained from a network pre-trained on synthetic data cannot model complex variations in real face images.

Instead of using static independent real images, images from different viewpoints or video sequences of the same scene have been used to reduce the gap

between real and synthetic data for intrinsic decomposition of indoor and/or outdoor scenes (e.g., [18], [44], [25]).

Inspired by the above approach, in this work we consider the use of images from real face videos<sup>3</sup> to bridge the domain gap because consistency constraints can be derived from different frames of the same person and they provide a stronger error measure in learning than the simple image reconstruction loss. Our approach is motivated by the key observation that the albedo and normal maps of two face frames in a same-person video only differ in pose and expression and that consistency of these maps can be derived. Specifically, we transform the components between frames to align albedo and normal: the transformation of albedo is the displacement of corresponding pixels, while for normal it consists of displacement and change of directions. Note that such transformations cannot be achieved by traditional (e.g. optical flow) algorithm due to the texture-less nature of albedo and directional shifts in normal. After aligning one frame with another, we leverage the consistency constraint between their albedo and normal maps to regularize the decomposition procedure, leading to weakly supervised learning on real face images. This way, the domain gap between real and synthetic face images is drastically reduced. Our experiments corroborates that the alignment (though not perfect) obtained by *AlignNet<sub>a</sub>* and *AlignNet<sub>n</sub>* can effectively improve the inverse rendering performance on real data.

We also note that existing methods do not take into account high frequency details such as shadow and highlights, making the inverse rendering results less realistic. Instead of considering the inverse rendering problem as decomposing an image into three components, we treat it as decomposing an image into four components (albedo, normal, illumination, and residual). The advantage of the additional residual map component in our new approach is that it leads to more realistic face images that maintain high frequency details induced by global illumination effects such as shadow and highlights. We propose a learning framework called *IlluRes-SfSNet* to perform face image decomposition (into albedo, normal, illumination, and residual). We train *IlluRes-SfSNet* on both real and synthetic data<sup>4</sup>, benefiting from both, and compute the albedo, normal, illumination components in a similar way to existing works (e.g., the *SfSNet* [34]). The residual map is obtained by subtracting the image with global illumination and the one that only contains local illumination. We run extensive experiments to evaluate *IlluRes-SfSNet* on various benchmarks, obtaining much better results than the state-of-the-art.

To summarize, our main contributions are threefold:

- The conceptual significance of introducing a consistency assumption of albedo and normal of the same person in a video. This consistency assumption led to the idea of weakly supervised training of our neural network model, hence bridging the domain gap between synthetic and real images.

---

<sup>3</sup> We do not employ multi-view images because they do not provide additional lighting constraints.

<sup>4</sup> We create two synthetic datasets (in addition to an existing one) to train our model.

- Our proposed dual CNN models  $AlignNet_a$  and  $AlignNet_n$  as workhorses for learning the alignment between albedo and normal maps of different frames in a same-person face video. After pretraining on synthetic data, they are applied to real data to align the predicted albedo and normal maps, enabling weakly supervised learning on real data.
- An IlluRes-SfSnet learning framework that seamlessly integrates SfSnet [34] and Illumination Residual Net to predict a residual map, in addition to albedo, normal, and illumination, for better inverse rendering results than the state-of-the-art.

Our current work opens doors to development of exciting applications such as relighting and albedo editing, for which we show some results in the end of the paper.

## 2 Related work

**Inverse rendering of images:** Decomposing an image into its intrinsic components is a long-standing and challenging task in computer vision. The most popular forms are intrinsic images [15][5][39][7], which define the decomposition layers as reflectance and shading (the function of shape and illumination). Recently, SIRFS [4] showed that further recovering surface normal and lighting from shading can improve the performance of decomposition. Since it is impossible to learn the decomposition without any constraints on the intrinsic components, classical algorithms for inverse rendering usually rely on a sophisticated design of priors, such as sparsity of reflectance [37][31], user strokes [8][35], and RGB-D settings [22][3][10]. These priors lead to promising decomposition in specific applications but do not generalize well to in-the-wild images. Deep learning has also been applied to inverse rendering. Given the lack of real ground truth data, some works perform supervised learning on synthetic datasets[29][12][23], but they still suffer from poor performance on real data. Our paper focuses on inverse rendering of human face images, which commonly serve as major objects in real world photos.

**Inverse rendering of human face:** Since 3D morphable model (3DMM) of faces was proposed in 1999 [6], it has served as a statistical shape prior for face inverse rendering [21][43][17][24]. Tewari *et al.* [41] combined deep learning- and model-based capture in an end-to-end network to infer intrinsic components from a single input image, with a well designed differential parametric decoder; they [40] further developed an algorithm to learn the facial shape and reflectance variations from unconstrained images, without a pre-existing shape identity or albedo model. However, these works are still based on parametric models on shape and reflectance, which resulted in loss of details.

Another branch of research focuses on model-free structures [38][34]. Sen Gupta *et al.* [34] designed a training paradigm called SfSNet to learn fine-scale separation of albedo and normal; they first trained a simple network on synthetic data to obtain coarse estimations for real images and then trained another decomposition architecture with residual blocks, on both synthetic and real data;

they claim that SfsNet can outperform state-of-the-art algorithms for inverse rendering of faces, but SfsNet’s performance is restricted by the coarsely estimated labels of real data. In addition, SfsNet only models local illumination of images, ignoring the spatially-varying components. This leads to artifacts on albedo and normal map, when there are obvious cast shadows or mutual illumination. We propose a weakly supervised learning on real face videos, avoiding the dependence on labels of real data. Moreover, we design an IlluRes-SfsNet framework to extract spatially-varying illumination information.

**Inverse rendering of in-the-wild images:** Due to the lack of dataset, multi-image based approaches have been introduced, with consistency of scene variables in images being used to constrain the solution, especially the similarity in albedo of the same object. Researchers in [20][27][25] trained their network with a set of images of a scene under varying illuminations but the same viewpoint, to help disambiguate albedo and shading, and the author of [44] further removed the constraint on the viewpoint of inputs by warping albedo before measuring similarity.

### 3 Method

In this section, we introduce our method for inverse rendering of in-the-wild face images. Similar to [34], we also consider human faces as Lambertian surfaces. Given the albedo  $\mathbf{A}_{p \times q \times 3}$ , normal  $\mathbf{N}_{p \times q \times 3}$ , and lighting  $\mathbf{L}_{9 \times 3}$ <sup>5</sup>, the face image  $\mathbf{I}_{p \times q \times 3}$  can be rendered by

$$\mathbf{I} = \mathbf{A} \circ f(\mathbf{N}, \mathbf{L}), \quad (1)$$

where  $f(\cdot)$  denotes the function that renders the shading image from normal and lighting, and  $\circ$  denotes element-wise multiplication.

However, the above image formation model does not consider global illumination effects such as shadow, highlights, light interactions, etc<sup>6</sup>. Thus, the face images rendered using the above model are often less realistic. In order to alleviate this problem, we propose to add a residual map to account for the global illumination effects. The enhanced image formation model can then be written as

$$\mathbf{I}_g = \mathbf{I}_l + \mathbf{R}, \quad (2)$$

where  $\mathbf{I}_g$  denotes the image considering global illumination,  $\mathbf{I}_l = \mathbf{A} \circ f(\mathbf{N}, \mathbf{L})$  stands for the image only considering local illumination, and  $\mathbf{R}$  is the residual map that compensates global illumination effects.

Our objective is to decompose the face images into  $\mathbf{A}$ ,  $\mathbf{N}$ ,  $\mathbf{L}$ , and  $\mathbf{R}$ , which can be used to render realistic face images. To achieve this goal, we first make

<sup>5</sup> Following [34], we also use the first-three order spherical harmonic coefficients to encode lighting, and we repeat three times for each color channel here.

<sup>6</sup> Here we follow the same definitions as in [32] for global illumination and local illumination.

use of synthetic data. We create two synthetic datasets and expand one dataset from existing synthetic datasets. Each data sample contains  $\mathbf{I}_g$ ,  $\mathbf{I}_l$ ,  $\mathbf{A}$ ,  $\mathbf{N}$ , and  $\mathbf{L}$ , allowing fully supervised learning. However, training on synthetic data alone is not enough to decompose real face images faithfully, since there is a huge domain gap between real and synthetic images. Therefore, it is essential to include real face data in our training. Unfortunately, there is no ground-truth  $\mathbf{I}_l$ ,  $\mathbf{A}$ ,  $\mathbf{N}$ , and  $\mathbf{L}$  available for real face images. To get around the problem, we propose a weakly supervised learning method using real face videos.

In the sequel, we cover supervised learning in IlluRes-SfSNet in Section 3.1 first, we then describe weakly supervised learning by further introducing *AlignNet<sub>a</sub>* and *AlignNet<sub>n</sub>* in Section 3.2.

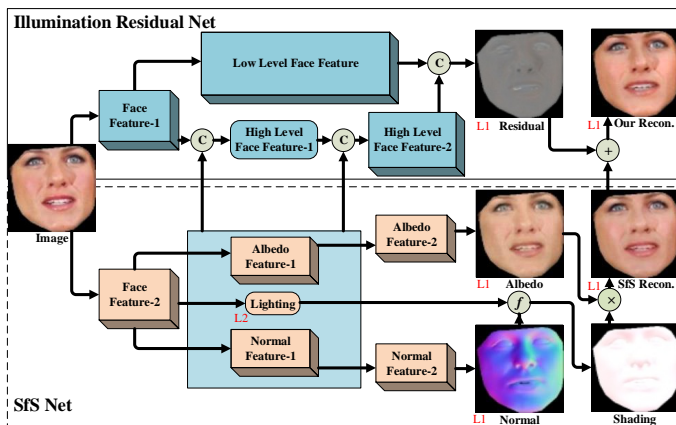


Fig. 2: The IlluRes-SfSNet architecture that decomposes an image into albedo, normal, lighting, and residual. It consists of two subnets: SfSNet and Illumination Residual Net, with SfSNet decomposing a face image into albedo, normal, and lighting, and Illumination Residual Net predicting the residual map that accounts for the global illumination effects such as shadow, highlights, etc. The two subnets are seamlessly combined with each other.

### 3.1 Supervised learning on synthetic data

Our learning framework, called IlluRes-SfSNet, for supervised learning on synthetic data is shown in Fig.2. In general, our model takes a single face image as input and decomposes it into its corresponding albedo, normal, lighting, and residual map. IlluRes-SfSNet consists of two subnets: SfSNet for decomposing a face image into normal map  $\mathbf{N}$ , albedo  $\mathbf{A}$ , and lighting  $\mathbf{L}$ , and Illumination Residual Net for extracting the global illumination effects to compensate details that cannot be obtained by SfSNet. The final face image  $\mathbf{I}_g$  is rendered by adding the residual map  $\mathbf{R}$  from Illumination Residual Net to the face image  $\mathbf{I}_l$  with

local illumination obtained from SfsNet. Note that the two subnets are seamlessly combined with each other: features learned in SfsNet are combined with those learned in Illumination Residual Net for better residual map estimation, and the gradients for learning  $\mathbf{I}_g$  are back-propagated through both SfsNet and Illumination Residual Net. In this manner, information learned by both subnets mutually enhance each other. IlluRes-SfsNet is trained in two stages on synthetic data. In the first stage, we train SfsNet by minimizing

$$L_s = \lambda_l \|\mathbf{I}_l - \bar{\mathbf{I}}_l\| + \lambda_a \|\mathbf{A} - \bar{\mathbf{A}}\| + \lambda_n \|\mathbf{N} - \bar{\mathbf{N}}\| + \lambda_h \|\mathbf{L} - \bar{\mathbf{L}}\|^2, \quad (3)$$

where  $\mathbf{X}$  and  $\bar{\mathbf{X}}$  represent the ground truth and the prediction, respectively, and  $\lambda_{\{l,a,n,h\}}$  the weights of each loss term. In the second stage, we train the whole IlluRes-SfsNet with the loss function

$$L_g = \lambda_g \|\mathbf{I}_g - \bar{\mathbf{I}}_g\| + \lambda_r \|\mathbf{R} - \bar{\mathbf{R}}\| + L_s, \quad (4)$$

where  $\lambda_g$  and  $\lambda_r$  represent the weights for image reconstruction and residual regression losses.

### 3.2 Weakly supervised learning on real data

As mentioned earlier, although our IlluRes-SfsNet can improve the details by considering global illumination effects, training on synthetic data alone is still not enough to bridge the gap between real and synthetic data. Due to the lack of ground-truth decomposition in real image datasets, it is not possible to directly train on real face images. Simply adopting reconstruction loss on real images does not help, since reconstruction loss alone does not provide any constraint on decomposition of real face images. Instead, we propose to train on real face videos to mitigate this problem. Although the ground truth is still not available, we can subtly make use of the consistency among video frames of the same identity, thus allowing weakly supervised learning on real face data.

Specifically, we observe that two different frames in the face video of the same identity only differ in pose and expression, which indicates that albedo/normal maps between them are subject to a transformation. However, this is not the case for lighting, since the shadings between different frames may vary a lot. The transformation between different albedo/normal maps can be easily learned using fully supervised method [46]. Since albedo and normal are independent from lighting and lighting in real data is much more complicated than synthetic data, we conjecture that the domain gap between albedo and normal of synthetic and real data is smaller than that between images of synthetic and real images. Therefore, we propose to learn the transformation between albedo and normal of different frames from the same person using synthetic data; this way the trained model can be reliably transferred to real images.

Note that, similar to SfsNet [34], IlluRes-SfsNet is model free. However, the key distinction between them is that there is no need for coarsely estimated labels of real data in IlluRes-SfsNet. Instead, it exploits consistency of albedo and normal in synthetic and real data.

Instead of using a single CNN to align albedo and normal maps of different face images, we use two CNNs, namely,  $AlignNet_a$  and  $AlignNet_n$ , for albedo and normal, respectively, to account for context information. In particular, the transformation between albedo of different images of the same person is mainly geometric transformation. Despite of geometric transformation, normal directions of the same point will change with different poses and expressions. This explains why we adopt two CNNs instead of one.

To enforce the network to learn the geometric transformations for albedo and normal, we propose to train  $AlignNet_a$  and  $AlignNet_n$  using face contours, which mainly contain the geometric transformation, together with the albedo and normal maps. Using albedo and normal maps alone may result in trivial solutions such as color transformation. The face contour  $\mathbf{C}$  can be easily obtained by detecting facial landmarks on a face image and connecting the detected landmarks on each parts such as eyes, nose, and mouth [1].

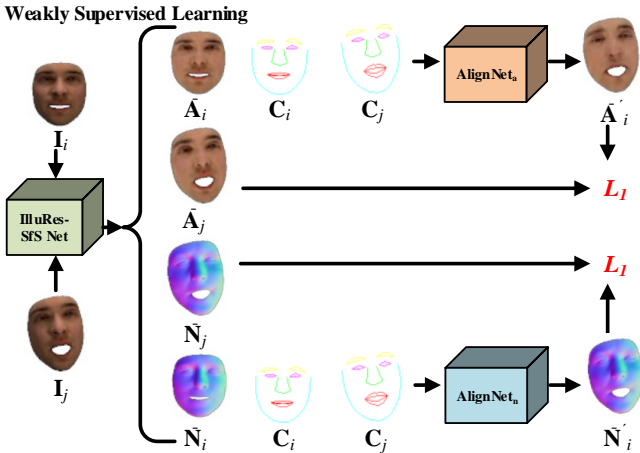


Fig. 3: Weakly supervised learning of  $AlignNet_a$  and  $AlignNet_n$  in IlluRes-SFSNet on real video data.

Fig. 3 shows our weakly supervised learning pipeline on real data. Here, we only introduce weakly supervised learning of  $AlignNet_a$  based on consistency of albedo as the process for  $AlignNet_n$  on normal is very similar. Consider two frames  $\mathbf{I}_i$  and  $\mathbf{I}_j$  from the video of a particular identity, our pre-trained IlluRes-SFSNet first predicts albedo, normal, lighting, and residual map for each of them. By taking the predicted albedo  $\bar{\mathbf{A}}_i$  of the  $i$ -th frame, together with contours  $\mathbf{C}_i$  and  $\mathbf{C}_j$  of the  $i$ -th and  $j$ -th frame, respectively, as inputs,  $AlignNet_a$  first transforms  $\bar{\mathbf{A}}_i$  to  $\bar{\mathbf{A}}'_i$ . Ideally,  $\bar{\mathbf{A}}'_i$  and albedo  $\bar{\mathbf{A}}_j$  of the  $j$ -th frame should be identical under the consistency constraint.



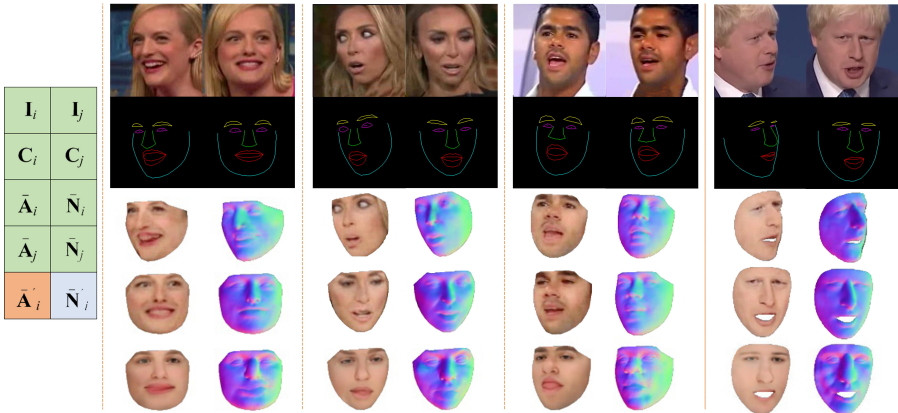


Fig. 4: Example outputs of  $AlignNet_a$  and  $AlignNet_n$  in IlluRes-SfSNet on real videos.

We thus train  $AlignNet_a$  using the loss function

$$L_a = \|\bar{\mathbf{A}}_j - \phi_a(\mathbf{C}_i, \mathbf{C}_j, \bar{\mathbf{A}}_i)\|, \quad (5)$$

where  $\phi_a$  transforms  $\bar{\mathbf{A}}_i$  to  $\bar{\mathbf{A}}'_i$  by applying the transformation between  $\mathbf{C}_i$  and  $\mathbf{C}_j$  on  $\bar{\mathbf{A}}_i$ . Similarly,  $AlignNet_n$  is trained with the loss function

$$L_n = \|\bar{\mathbf{N}}_j - \phi_n(\mathbf{C}_i, \mathbf{C}_j, \bar{\mathbf{N}}_i)\|, \quad (6)$$

where  $\phi_n$  transforms  $\bar{\mathbf{N}}_i$  to  $\bar{\mathbf{N}}'_i$  by applying the transformation between  $\mathbf{C}_i$  and  $\mathbf{C}_j$  on  $\bar{\mathbf{N}}_i$ .

In this way, IlluRes-SfSNet can then be jointly trained with  $AlignNet_a$  and  $AlignNet_n$  on real face videos by minimizing

$$L = \lambda_g \|\mathbf{I}_g - \bar{\mathbf{I}}_g\| + \lambda_{ab} L_a + \lambda_{no} L_n, \quad (7)$$

where  $\lambda_{ab}$  and  $\lambda_{no}$  are weights for  $AlignNet_a$  and  $AlignNet_n$  losses. In practice, we first pre-train  $AlignNet_a$  and  $AlignNet_n$  using synthetic data with ground-truth  $\mathbf{A}_{\{i,j\}}$  and  $\mathbf{N}_{\{i,j\}}$ , and then conduct weakly supervised learning on real face videos. Note that we align albedo and normal maps on two directions, namely, from  $i$  to  $j$  and from  $j$  to  $i$ , in order to have stronger consistency.

We present a few example outputs of  $AlignNet_{\{a,n\}}$  in Fig. 4. Note that although the prediction of  $AlignNet_{\{a,n\}}$  is not perfect, our experiments verify that the consistency constraint on real face video can still improve the inverse rendering performance, which corroborates our assumption that coarse alignment on albedo and normal is enough to perform weakly supervised learning on real face images.

## 4 Experimental results

### 4.1 Datasets and implementation details

**Synthetic data:** To conduct supervised training for our model, we make use of three synthetic datasets. Among them, one is an existing dataset and the other two are created by us based on other datasets. We first use the synthetic dataset provided by SfsNet [34], denoted as *SfsNet-syn*. This dataset contains 12,500 identities. 15 images are rendered for each identity under varying illuminations and poses. The images are rendered by fitting the 3DMM [6]. In total, there are 187,500 images in this dataset. However, this dataset does not take expression variations and global illuminations into consideration. Therefore we create two other synthetic datasets *Expre-syn* and *Caric-syn* to compensate for this. For *Expre-syn*, we randomly choose 114 video clips from VoxCeleb [28] and each clip corresponds to one identity. We fit the Basel Face Model [30] for each frame to approximate realistic poses and expressions under four random viewpoints and five random illuminations using Mitsuba [16]. Six different expressions are fitted for each identity. In total we obtain 13,680 images in *Expre-syn*. Since in practice there can be some exaggerated expressions that do not appear in VoxCeleb, we further create *Caric-syn* to increase expression variations with 100 virtual identities. For each identity, we render four distinct poses with six caricature expressions under five random lighting conditions, resulting in 12,000 synthetic face images. Our data generation process inherently considers global illumination. For each face image  $\mathbf{I}_g$ , we also obtain the corresponding albedo  $\mathbf{A}$ , normal  $\mathbf{N}$ , lighting  $\mathbf{L}$ , residual map  $\mathbf{R}$ , and face image  $\mathbf{I}_l$  with local illumination. A summary of the synthetic data is given in Table 1. We use images for one identity from each of *Expre-syn* and *Caric-syn* for testing and all the rest together with *SfsNet-syn* for training.

In addition, to train AlignNet, we further aggregate albedo and normal in *Expre-syn* and *Caric-syn*. Specifically, we add various textures into synthetic albedos to help *AlignNet<sub>a</sub>* focus on learning displacements. We create 50 texture styles and randomly transfer them onto face albedos. We also perform shape deformation on synthetic face models to avoid strong face priors in *AlignNet<sub>n</sub>*. We design 20 different deformations on faces and build a linear combination basis. The deformations are generated using a sketch-based modelling system [14]. In total, we have 540 different albedos under four viewpoints and six expressions for *AlignNet<sub>a</sub>* training, and 540 identities with four viewpoints and six deformation styles (produced by randomly weighting the linear deformation bases) for *AlignNet<sub>n</sub>*. Note that we only use this dataset to train AlignNet.

**Real data:** We extract 114 real video clips from the 300VW dataset [36]. Each clip corresponds to one identity (with four viewpoints and five illuminations). There are 3,462 images in total. This dataset contains rich but non-exaggerated expressions. It is used for weakly supervised training only. For qualitative evaluations, we use face images in CelebA [26] by picking 500 images with various face attributes (e.g., bread, glasses, and age).

For quantitative evaluation on normal recovery, we adopt Photoface [45], which is created under various harsh illumination conditions and captures the ground truth of face normal.

**Implementation details:** Our network architecture is implemented in TensorFlow. The input images are all of size  $128 \times 128$ . We train our model using a batch size of four using the learning rate of  $5e-4$  with a decay rate of 0.98 and a decay step of five. As a preparation procedure, we pre-train  $AlignNet_a$  and  $AlignNet_n$  with synthetic data, with a learning rate of  $5e-4$  for 30 epochs using a batch size of eight. We then train our inverse rendering model in three stages: first, we train our IlluRes-SfSNet for 20 epochs on the combined synthetic dataset of SfSNet-syn, *Caric-syn*, and *Expre-syn*; second, we train IlluRes-SfSNet on real data by using the pre-trained  $AlignNet_a$  and  $AlignNet_n$  to provide weak supervision for 20 epochs; and third, we jointly fine-tune the entire network that includes IlluRes-SfSNet,  $AlignNet_a$  and  $AlignNet_n$  on both real and synthetic data. During the first stage of training, for each epoch we randomly sample 30,000 images from *SfSNet-syn* and use all the training data from *Caric-syn* (13,560 images) and *Expre-syn* (11,880 images) to avoid the training being biased by *SfSNet-syn*.

During the second stage of training, we use the real videos containing 3,487 images, *Caric-syn*, and *Expre-syn* on the fixed  $AlignNet_{\{a,n\}}$  to train IlluRes-SfSNet. In each epoch, we randomly sample 30,000 pairs of real images and 16,272 pairs of *Expre-syn* images. At the last stage, we use synthetic and real data to fine tune the entire network.

Table 1: Dataset construction with controlled variations for each set. Note that for *SfSNet-syn* dataset, each identity has 15 images rendered under varying illuminations and poses.

Dataset	Iden.	Per Identity				Total
		Illu.	Pose	Expr.	Exag.	
<i>SfSNet-syn</i>	12.5k	15	15	\	\	<b>187.5k</b>
<i>Expre-syn</i>	114	5	4	6	\	<b>13.68k</b>
<i>Caric-syn</i>	100	5	4	\	6	<b>12.00k</b>
Real	114	\	\	25~35	\	<b>3462</b>

## 4.2 Comparison with state of the art

**Inverse rendering:** We compare our full model, IlluRes-SfSNet-Align, with NeuralFace[38] and SfSNet[34]. Note that since different strategies are used for learning from real data in [38][34], we directly use their released pre-trained models to reproduce results for inclusion in Fig.5. It can be seen that IlluRes-SfSNet-Align and SfSNet generate more realistic decomposition than NeuralFace. The regions in red rectangles show our method could alleviate the ambiguity on face shape, which are caused by makeup, such as the moustache in the first sample and the eyeliner in the second sample. Our method decomposes them correctly.

Moreover, compared with SfsNet, our model provides extra discernible shadows on albedo while capturing more fine-scale details in the normal map.

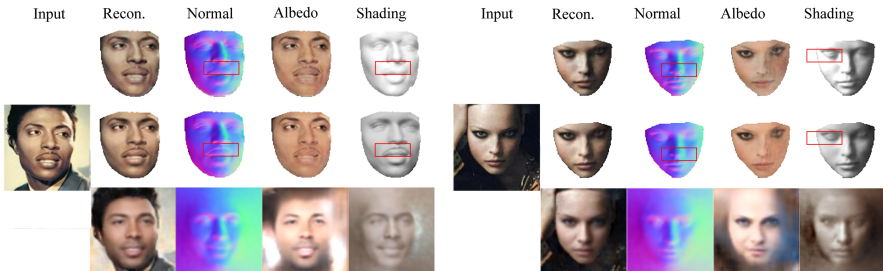


Fig. 5: Inverse rendering results on two samples. The first row of each sample shows our decomposition results. The following rows come from SfsNet[34] and NeuralFace[38]. (Best viewed in PDF with zoom.)

**Shape estimation:** Next we compare the quality of shape estimation among llures-SfsNet-Align, 3DMM and Pix2Vertex [33]. The dataset used for quantitative testing is Photoface [45], which provides ground truth of face normal maps under various illumination conditions. Following [42][34], we randomly pick 100 identities from a total of 454. We evaluate the performance by mean angular error of the normals and the percentage of pixels at various angular error thresholds and report them in Table 2. The higher percentage of the fourth column means a model could capture more low frequency shape information, while data of smaller angle measures the performance on estimating fine-scale structure. It can be seen that 3DMM achieves the highest percentage in the fourth column but quite small value at the second one, which means it captures most low frequency shape while losing lots of detailed information. Pix2Vertex captures much more details in the images but fails to recover shape information from harsh lighting and unusual expression, as the lowest value in the fourth column suggested. Fig.6 compares estimated normals of the showcases in [34]. The normal images agrees with the Table 2 that our model outperforms SfsNet on normal estimation.

Table 2: Comparison on normal reconstruction error on the Photoface dataset. Lower is better for column 1, and higher is better for the percentage of pixels at specific thresholds.

Algorithm	Mean $\pm$ std	$<20^\circ$	$<25^\circ$	$<30^\circ$
3DMM	$31.9 \pm 11.6$	3.7%	53.2%	87.3%
Pix2Vertex[33]	$36.3 \pm 6.2$	23.6%	35.4%	46.1%
SfsNet[34]	$30.7 \pm 8.6$	40.5%	54.2%	64.5%
Ours	<b><math>25.3 \pm 6.3</math></b>	<b>43.6%</b>	<b>57.6%</b>	<b>68.8%</b>

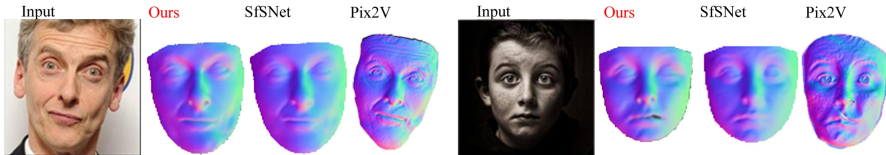


Fig. 6: Normal recovery results on the showcases of SfsNet[34]. The first two columns are results from inverse rendering while the last column is from face reconstruction algorithm.

### 4.3 Ablation Studies

**IlluRes-SfsNet v. SfsNet:** To demonstrate the effectiveness of weakly supervised learning in Illuses-SfsNet, we use a total of five models. They are

**Sfs-pretrain:** The authors of [34] provide this model which is trained by using Sfs-supervision. They first use a ‘skip-net’ to obtain  $\mathbf{A}$  and  $\mathbf{N}$  of real images, then combine synthetic and real data with ‘ground-truth’ to train SfsNet. As Sfs-pretrain has been trained with ‘Sfs-supervision’, real images can be reconstructed. However, SfsNet’s structure cannot handle global illumination effects regardless of how training is done. These effects always exist on albedo and/or normal.

**Sfs:** We use synthetic data to train Sfs. In some situations this model outputs reasonable normal or albedo maps. Since it cannot be trained on real images, its overall reconstruction results are not satisfactory.

**Sfs-Align:** Adding *AlignNet* to allow weakly supervised training in Sfs. This model can effectively decompose and reconstruct real images. Its performance is similar to that of Sfs-pretrain.

**Illuses-Sfs-pre:** We train Illures-Sfs with synthetic data without using *AlignNet*. IlluRes-Sfs-pre successfully compensates for defects in Sfs by separating the residual portion. Without weakly supervised learning, this network cannot bridge the gap between real and synthetic data, hence its reconstruction results are not competitive.

**IlluRes-SfsNet:** Adding *AlignNet* to Illuses-Sfs-pre. This model separates the residual through IlluResNet to capture fine image details. Moreover, *AlignNet* enables the network to effectively decompose and reconstruct real images.

Fig.7 compares decomposition and reconstruction results of the above five models. In the comparison of decomposition, our model is optimal, which proves that our *AlignNet* and Residual net are effective.

### 4.4 Application

Based on the inverse rendering results of our method, we are able to develop a wide range of applications. Here we show two of examples, namely, face relighting and albedo editing.



Fig. 7: Five models and four examples prove our *AlignNet* and *IlluResNet* are effective. (Best viewed in PDF with zoom.)



Fig. 8: The left of the figure is case of photo relighting. (a) are the source images and (b) are the results corresponding to the light information in the small photos. The right of the figure is case of albedo editing. (c) input photo, (d) result of albedo editing, (e) result of face photo editing. The corresponding albedo and modified one are shown in the small images.

To relight photo, source and target images are sent into the well-trained *IlluRes-SfSNet-Align*, decomposed into corresponding albedo, normal, and lighting information. The lighting of source image is replaced by the target one and the novel combination of components are used to generate the relighted image. Fig.8 shows the results.

Similarly, in albedo editing, the source photo is fed into our full model and produces the decomposition components. After modifying the albedo map, a new photo could be generated. We demonstrate a sample of beard editing in Fig.8. The result shows editing on albedo could generate more realistic result than directly editing on face photo. For example, we can see that the cheek in Fig.8(e) contains obvious artifacts, while the result of our model in Fig.8(d) is much more visually pleasing.

## 5 Conclusion

In this paper we propose a weakly supervised approach for inverse face rendering on real face videos, based on the assumption of consistency of albedo and normal of the same person in a video, bridging the domain gap between synthetic and real face. We propose  $AlignNet_{\{a,n\}}$  to align albedo and normal subspaces between different frames of a certain video clip. We empirically show that the alignment (though not perfect) obtained by  $AlignNet_{\{a,n\}}$  can provide enough constraints on frame consistency for weakly supervised learning on real images. Together with IlluRes-SfSNet, our framework has strong capability to disentangle normal and albedo into separate subspaces. Qualitative and quantitative evaluations show that our method outperforms state-of-the-art works on face inverse rendering.

## Acknowledgment

The work was supported in part by grants No. 2018YFB1800800, No. 2018B030338001, No. 2017ZT07X152, No. ZDSYS201707251409055 and in part by National Natural Science Foundation of China (Grant No.: 61902334 and 61629101).

## References

1. Faceplusplus. [www.faceplusplus.com](http://www.faceplusplus.com) (cited Nov 2019)
2. Aldrian, O., Smith, W.A.: Inverse rendering of faces with a 3d morphable model. *IEEE transactions on pattern analysis and machine intelligence* **35**(5), 1080–1093 (2012)
3. Barron, J.T., Malik, J.: Intrinsic scene properties from a single rgb-d image. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 17–24 (2013)
4. Barron, J.T., Malik, J.: Shape, illumination, and reflectance from shading. *IEEE transactions on pattern analysis and machine intelligence* **37**(8), 1670–1687 (2014)
5. Barrow, H.: *Recovering intrinsic scene characteristics* (1978)
6. Blanz, V., Vetter, T., et al.: *A morphable model for the synthesis of 3d faces*. (1999)
7. Bonneel, N., Kovacs, B., Paris, S., Bala, K.: Intrinsic decompositions for image editing. In: *Computer Graphics Forum*. vol. 36, pp. 593–609. Wiley Online Library (2017)
8. Bousseau, A., Paris, S., Durand, F.: User-assisted intrinsic images. In: *ACM Transactions on Graphics (TOG)*. vol. 28, p. 130. ACM (2009)
9. Bradley, D., Heidrich, W., Popa, T., Sheffer, A.: High resolution passive facial performance capture. In: *ACM transactions on graphics (TOG)*. vol. 29, p. 41. ACM (2010)
10. Chen, Q., Koltun, V.: A simple model for intrinsic image decomposition with depth cues. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 241–248 (2013)
11. Debevec, P., Hawkins, T., Tchou, C., Duiker, H.P., Sarokin, W., Sagar, M.: Acquiring the reflectance field of a human face. In: *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*. pp. 145–156. ACM Press/Addison-Wesley Publishing Co. (2000)
12. Fan, Q., Yang, J., Hua, G., Chen, B., Wipf, D.: Revisiting deep intrinsic image decompositions. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 8944–8952 (2018)
13. Ghosh, A., Fyfe, G., Tunwattanapong, B., Busch, J., Yu, X., Debevec, P.: Multi-view face capture using polarized spherical gradient illumination. In: *ACM Transactions on Graphics (TOG)*. vol. 30, p. 129. ACM (2011)
14. Han, X., Gao, C., Yu, Y.: Deepsketch2face: a deep learning based sketching system for 3d face and caricature modeling. *ACM Transactions on Graphics (TOG)* **36**(4), 126 (2017)
15. Horn, B.K.: Determining lightness from an image. *Computer graphics and image processing* **3**(4), 277–299 (1974)
16. Jakob, W.: Mitsuba renderer (2010), <http://www.mitsuba-renderer.org>
17. Kemelmacher-Shlizerman, I., Basri, R.: 3d face reconstruction from a single image using a single reference face shape. *IEEE transactions on pattern analysis and machine intelligence* **33**(2), 394–405 (2010)
18. Kim, K., Torii, A., Okutomi, M.: Multi-view inverse rendering under arbitrary illumination and albedo. In: *ECCV* (2016)
19. Kimmel, R., Elad, M., Shaked, D., Keshet, R., Sobel, I.: A variational framework for retinex. *International Journal of computer vision* **52**(1), 7–23 (2003)
20. Laffont, P.Y., Bazin, J.C.: Intrinsic decomposition of image sequences from local temporal variations. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 433–441 (2015)



21. Lee, J., Machiraju, R., Pfister, H., Moghaddam, B.: Estimation of 3d faces and illumination from single photographs using a bilinear illumination model (2005)
22. Lee, K.J., Zhao, Q., Tong, X., Gong, M., Izadi, S., Lee, S.U., Tan, P., Lin, S.: Estimation of intrinsic image sequences from image+ depth video. In: European Conference on Computer Vision. pp. 327–340. Springer (2012)
23. Lettry, L., Vanhoey, K., Van Gool, L.: Darn: a deep adversarial residual network for intrinsic image decomposition. In: 2018 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 1359–1367. IEEE (2018)
24. Li, C., Zhou, K., Lin, S.: Intrinsic face image decomposition with human face priors. In: European conference on computer vision. pp. 218–233. Springer (2014)
25. Li, Z., Snavely, N.: Learning intrinsic image decomposition from watching the world. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition pp. 9039–9048 (2018)
26. Liu, Z., Luo, P., Wang, X., Tang, X.: Large-scale celebfaces attributes (celeba) dataset
27. Ma, W.C., Chu, H., Zhou, B., Urtasun, R., Torralba, A.: Single image intrinsic decomposition without a single intrinsic image. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 201–217 (2018)
28. Nagrani, A., Chung, J.S., Zisserman, A.: Voxceleb: a large-scale speaker identification dataset. arXiv preprint arXiv:1706.08612 (2017)
29. Narihira, T., Maire, M., Yu, S.X.: Direct intrinsics: Learning albedo-shading decomposition by convolutional regression. In: Proceedings of the IEEE international conference on computer vision. pp. 2992–2992 (2015)
30. Paysan, P., Knothe, R., Amberg, B., Romdhani, S., Vetter, T.: A 3d face model for pose and illumination invariant face recognition. In: 2009 Sixth IEEE International Conference on Advanced Video and Signal Based Surveillance. pp. 296–301. Ieee (2009)
31. Rother, C., Kiefel, M., Zhang, L., Schölkopf, B., Gehler, P.V.: Recovering intrinsic images with a global sparsity prior on reflectance. In: Advances in neural information processing systems. pp. 765–773 (2011)
32. Schneider, A., Schonborn, S., Froben, L., Egger, B., Vetter, T.: Efficient global illumination for morphable models. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 3865–3873 (2017)
33. Sela, M., Richardson, E., Kimmel, R.: Unrestricted facial geometry reconstruction using image-to-image translation. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1576–1585 (2017)
34. Sengupta, S., Kanazawa, A., Castillo, C.D., Jacobs, D.W.: Sfsnet: Learning shape, reflectance and illuminance of faces in the wild. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 6296–6305 (2018)
35. Shen, J., Yang, X., Li, X., Jia, Y.: Intrinsic image decomposition using optimization and user scribbles. IEEE transactions on cybernetics **43**(2), 425–436 (2013)
36. Shen, J., Zafeiriou, S., Chrysos, G.G., Kossaifi, J., Tzimiropoulos, G., Pantic, M.: The first facial landmark tracking in-the-wild challenge: Benchmark and results. In: Proceedings of the IEEE International Conference on Computer Vision Workshops. pp. 50–58 (2015)
37. Shen, L., Yeo, C.: Intrinsic images decomposition using a local and global sparse representation of reflectance. In: CVPR 2011. pp. 697–704. IEEE (2011)
38. Shu, Z., Yumer, E., Hadap, S., Sunkavalli, K., Shechtman, E., Samaras, D.: Neural face editing with intrinsic image disentangling. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5541–5550 (2017)

39. Tappen, M.F., Freeman, W.T., Adelson, E.H.: Recovering intrinsic images from a single image. In: Advances in neural information processing systems. pp. 1367–1374 (2003)
40. Tewari, A., Bernard, F., Garrido, P., Bharaj, G., Elgharib, M., Seidel, H.P., Pérez, P., Zollhofer, M., Theobalt, C.: Fml: face model learning from videos. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 10812–10822 (2019)
41. Tewari, A., Zollhofer, M., Kim, H., Garrido, P., Bernard, F., Perez, P., Theobalt, C.: Mofa: Model-based deep convolutional face autoencoder for unsupervised monocular reconstruction. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1274–1283 (2017)
42. Trigeorgis, G., Snape, P., Kokkinos, I., Zafeiriou, S.: Face normals. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 38–47 (2017)
43. Wang, Y., Zhang, L., Liu, Z., Hua, G., Wen, Z., Zhang, Z., Samaras, D.: Face relighting from a single image under arbitrary unknown lighting conditions. IEEE Transactions on Pattern Analysis and Machine Intelligence **31**(11), 1968–1984 (2008)
44. Yu, Y., Smith, W.A.: Inverserendernet: Learning single image inverse rendering. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3155–3164 (2019)
45. Zafeiriou, S., Hansen, M., Atkinson, G., Argyriou, V., Petrou, M., Smith, M., Smith, L.: The photoface database. In: CVPR 2011 WORKSHOPS. pp. 132–139. IEEE (2011)
46. Zakharov, E., Shysheya, A., Burkov, E., Lempitsky, V.: Few-shot adversarial learning of realistic neural talking head models. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 9459–9468 (2019)

## Appendix

### A Architectures of AlignNet and IlluResNet

In this section, we describe our network in detail. To better illustrate, we name layers in the figures with abbreviations:  $C$  as Convolution layer followed by Batch Normalization,  $CT$  as Transposed Convolution layer followed by Batch Normalization,  $AP$  as Average Pooling layer and  $Res$  as Residual Block. Each  $Res$  consists of BN - ReLU - C128 - BN - ReLU - C128. For net parameters,  $k$  represents kernel size and  $s$  is stride. The parameter  $u$  of CT means the size of output after upsampling.  $LR$  is the abbreviation of Leaky ReLU and  $R$  is ReLU.

We first propose AlignNet for albedo and normal, shown in Fig.9.  $\mathbf{I}_i$  is albedo or normal of source frame and  $\mathbf{I}'_i$  is the estimated component of target frame.  $\mathbf{C}_i$  and  $\mathbf{C}_j$  are the corresponding face contour of source frame  $i$  and target frame  $j$ . These two networks have the same structure while they are trained on different data.

Details of IlluResNet are provided in Fig.10.

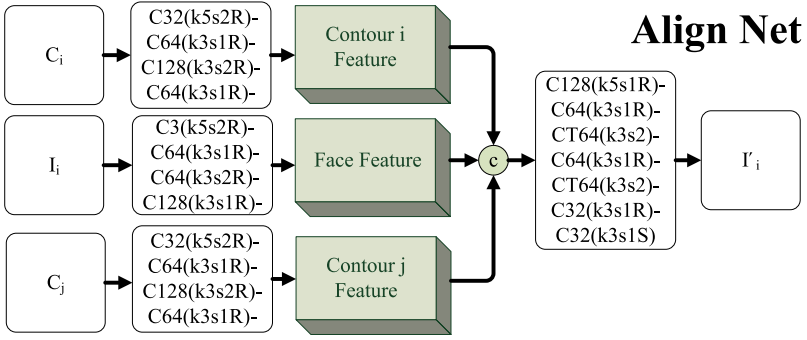


Fig. 9: AlignNet Architecture.

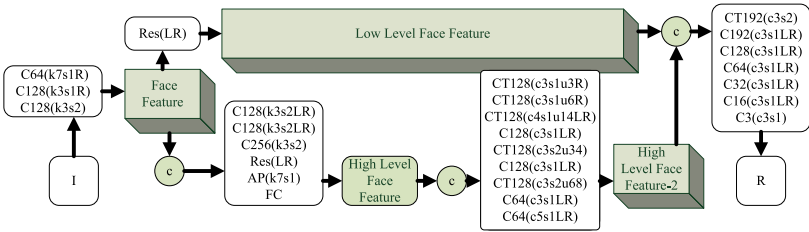


Fig. 10: IlluResNet Architecture.

## B More Qualitative Results

In Fig.11 and Fig. 12, we present inverse rendering results from our full model, IlluRes-SfSNet-Align, and SfSNet. The input images are sampled from CelebA[26].

In Fig.13, we further compare our method against SfSNet with a higher input resolution of  $256 \times 256$  (all the two networks are re-trained for such a new setting). With a higher resolution, it can be seen that our method can retain the details much better than SfSNet.

In Fig.14, we provide some samples for lighting transfer on faces.



Fig. 11: Inverse rendering. The first row of each sample is from IlluRes-SfSNet-Align, while the second one from SfSNet. In general, our method captures more details on normal, such as contour of nose and eyes, and wrinkles on face.

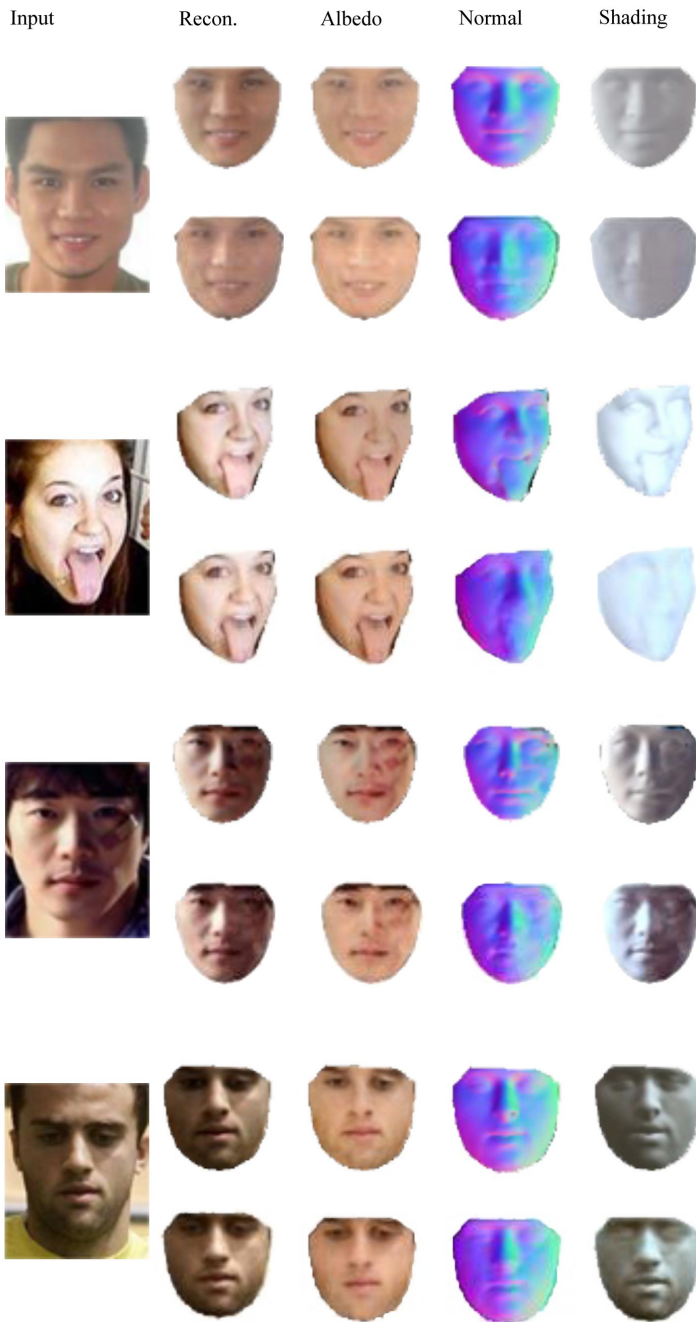


Fig. 12: Inverse rendering. The first row of each sample is from IlluRes-SfSNet-Align, while the second one from SfSNet.

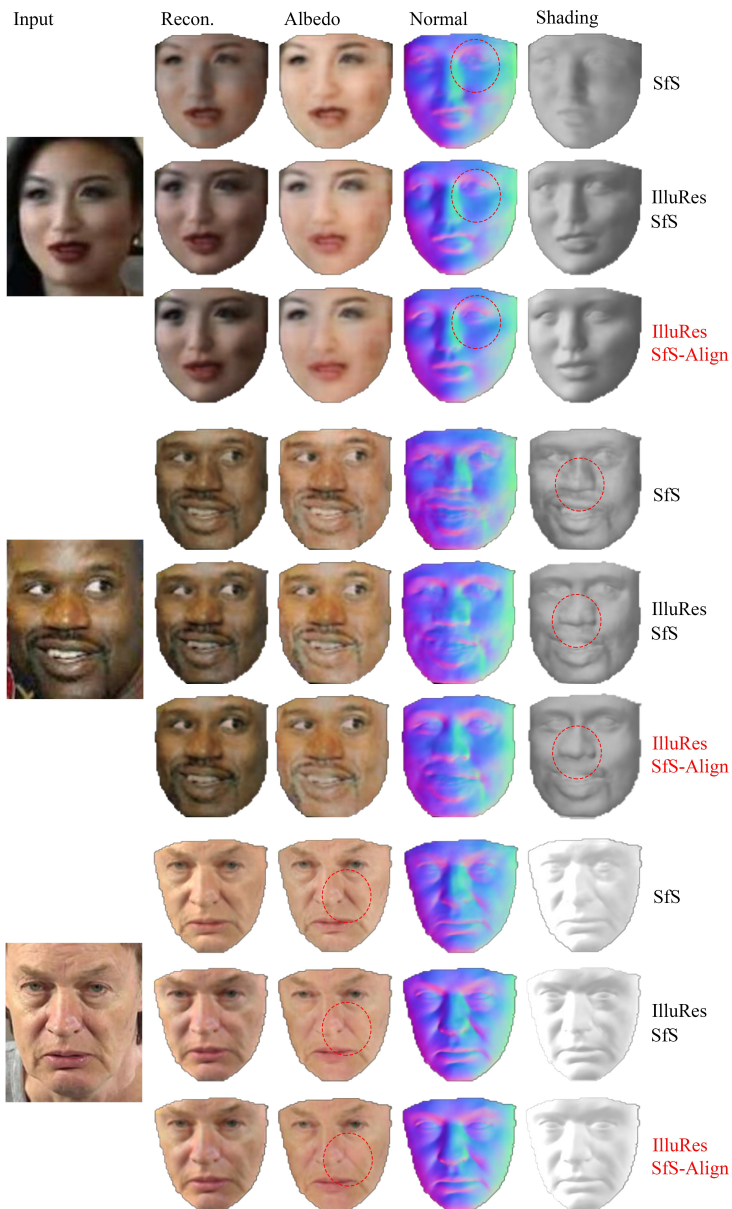


Fig. 13: Inverse rendering at  $256 \times 256$  resolution. The red circles highlight our improvements.

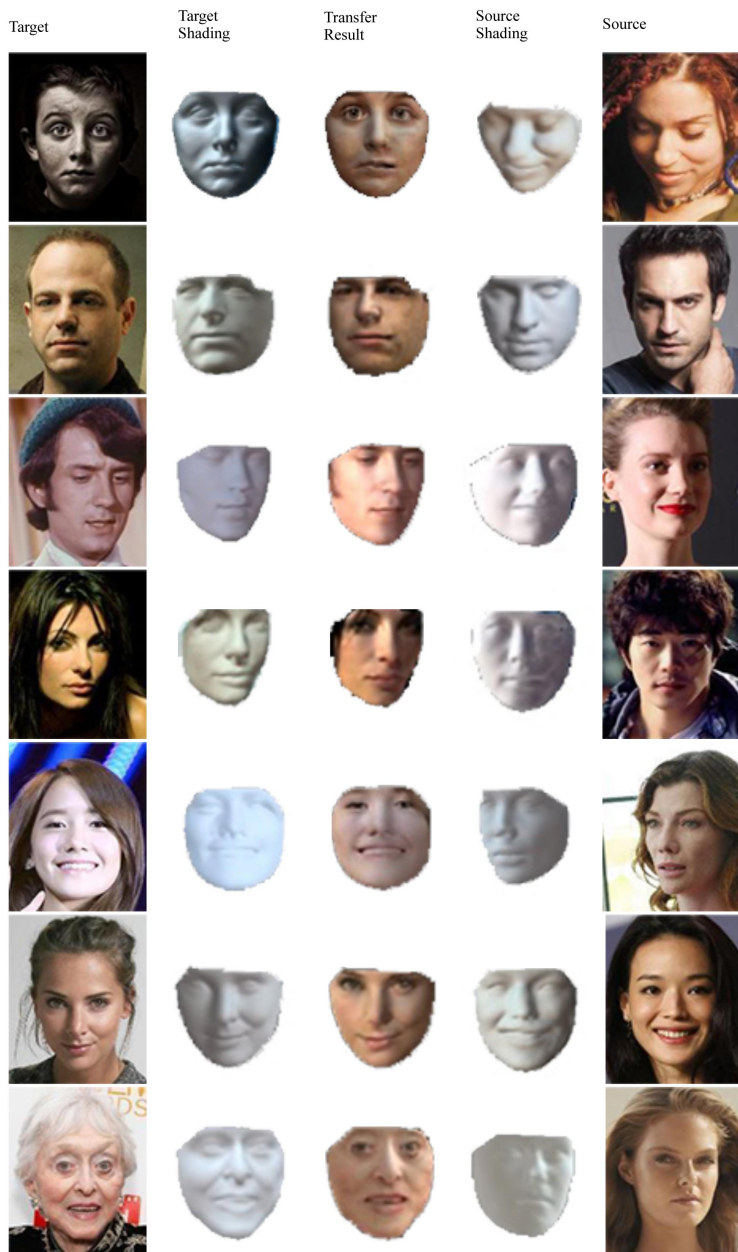


Fig. 14: Lighting transfer. Our method can transfer the lighting condition in source photo to target photo.