# High-fidelity Face Tracking for AR/VR via Deep Lighting Adaptation

Lele Chen[1,2]     Chen Cao[1]     Fernando De la Torre[1]     Jason Saragih[1]     Chenliang Xu[2]     Yaser Sheikh[1]

[1] Facebook Reality Labs     [2] Univeristy of Rochester

Figure 1. High-fidelity face tracking results using our method. From left to right: Input image captured by iPhone, normal map, fully-lit avatar (avatar before relighting), relit avatar (avatar after relighting), and relit avatar under different viewpoints. Please notice the specular highlight changes on the avatars under different viewpoints.

## Abstract

*3D video avatars can empower virtual communications by providing compression, privacy, entertainment, and a sense of presence in AR/VR. Best 3D photo-realistic AR/VR avatars driven by video, that can minimize uncanny effects, rely on person-specific models. However, existing person-specific photo-realistic 3D models are not robust to lighting, hence their results typically miss subtle facial behaviors and cause artifacts in the avatar. This is a major drawback for the scalability of these models in communication systems (e.g., Messenger, Skype, FaceTime) and AR/VR. This paper addresses previous limitations by learning a deep learning lighting model, that in combination with a high-quality 3D face tracking algorithm, provides a method for subtle and robust facial motion transfer from a regular video to a 3D photo-realistic avatar. Extensive experimental validation and comparisons to other state-of-the-art methods demonstrate the effectiveness of the proposed framework in real-world scenarios with variability in pose, expression, and illumination. Our project page can be found at https://www.cs.rochester.edu/u/lchen63. Please visit https://www.youtube.com/watch?v=dtz1LgZR8cc for more visual results.*

## 1. Introduction

Currently, video conferencing (*e.g.*, Zoom, Skype, Messenger) is the best 2D available technology for internet communication. To allow for more advance levels of communication and sense of presence, Augmented Reality (AR) and Virtual Reality (VR) technologies aim to build 3D personalized avatars, and superimpose virtual objects in the real space. If successful, this new form of face-to-face interaction will allow extended remote work experiences that can improve productivity, reducing cost and stress of commuting, have a huge impact on the environment, and overall improving the work/life balance.

Today most real-time systems for avatars in AR are cartoon-like (*e.g.*, Apple Animoji, Tiktok FaceAnimation, Hyprsense, Loom AI); on the other hand, digital creators in movies have developed uncanny digital humans using advanced computer graphics technology and person-specific (PS) models (*e.g.*, Siren). While some of these avatars can

1

be driven in real-time from egocentric cameras (*e.g.*, Doug character made by digital domain), building the PS model is an extremely time-consuming and hand-tuned process that prevents democratizing this technology. This paper contributes toward this direction, and it proposes new algorithms to robustly and accurately drive 3D video-realistic avatars from monocular cameras to be consumed by AR/VR displays (see Fig. 1).

Model-based photo-realistic 3D face reconstruction/animation from a video has been a core area of research in computer vision and graphics in the last thirty years [3, 4, 6, 8, 14, 16, 23, 28, 33, 49, 44, 59, 2]. While different versions of morphable models or active appearance models have provided good facial animation results, the existing 3D models do not provide the quality that is needed it for a good immersive viewing experience in AR/VR. In particular, the complex lighting, motion, and other in-the-wild conditions do result in artifacts in the avatar due to poor decouple of rigid and non-rigid motion, as well as, no accurate texture reconstruction. To tackle this problem, we build on recent work on Deep Appearance Model (DAM) [27] that learns a person-specific model from a multi-view capture setup. [27] can render photo-realistic avatars in a VR headset by inferring 3D geometry and view-dependent texture from egocentric cameras .

This paper extends DAM [27] with a new deep lighting adaptation method to recover subtle facial expressions from monocular videos in-the-wild and transfer them to a 3D video-realistic avatar. The method is able to decouple rigid and non-rigid facial motions, as well as, shape, appearance and lighting from videos in-the-wild. Our method combines a prior lighting model learned in a lab-controlled scenario and adapts it to the in-the-wild video, recovering accurate texture and geometric details in the avatar from images with complex illuminations and challenging poses (*e.g.* profile). There are two main contributions of our work. First, we provide a framework for fitting a non-linear appearance model (based on a variational auto-encoder) to in-the-wild videos. Second, we propose a new lighting transfer module to learn a global illumination model. Experimental validation shows that our proposed algorithm with deep lighting adaptation outperforms state-of-the-art methods and provides robust solutions in realistic scenarios.

## 2. Related Work

**3D Face Tracking.** 3D morphable face models (3DMM) [3, 4, 6, 14, 16, 23, 33] and Active Appearance Models (AAMs) [8, 28, 49] have been extensively utilized for learning facial animations from 3D scans and face tracking. These methods produce texture and geometry through the idea of analysis-by-synthesis, where a parametric face model is iteratively adapted until the synthesized face matches the target image. For example,

by leveraging the photometric error in both shape and texture, AAMs [8] have shown strong efficiency and expressibility to register faces in images. More recently, Deep Appearance Model (DAM) [27] extends the AAMs with deep neural networks in-place of linear generative functions. DAM learns the latent presentation of geometry and texture using a conditional variational autoencoder [19] and is able to reconstruct a high-fidelity view-dependent avatar with the aid of the multi-view camera system.

**In-the-Wild Face Reconstruction.** Face reconstruction under an in-the-wild scenario is known as a challenging problem since the rigid, lighting, and expression are unknown. For example, the surface information presented by a single image [15, 24, 34, 38, 42, 43, 57, 48, 54, 58, 56, 45, 17, 12] or even an image collection [35, 36, 25, 37, 50, 41, 46] is limited. Thus, achieving high-fidelity face reconstruction from in-the-wild imagery usually relies on prior knowledge like 3DMMs [3], FLAME [23], or DAM [27]. For example, instead of directly regressing the geometry, MoFA [43] uses a CNN-based image encoder to extract the semantically meaningful parameters (*e.g.*, facial expression, shape, and skin reflectance) from a single 2D image and then uses the parametric model-based decoder to render the output image. Tran *et al.* [48] propose a weakly supervised model that jointly learns a nonlinear 3DMM and its fitting algorithm from 2D in-the-wild image collection. Gecer *et al.* [13] train a facial texture generator in the UV space with self-supervision as their statistical parametric representation of the facial texture. Lin *et al.* [24] propose a GCN-based network to refine the texture generated by a 3DMM-based method with facial details from the input image. Yoon *et al.* [54] propose a network (I2ZNet) to learn a latent vector $z$ and head pose for DAM from a single image to reconstruct the texture and geometry.

**Lighting Estimation for In-the-Wild Imagery.** Most existing face reconstruction works [24, 26, 39, 42, 43, 47, 48, 29] estimate the illumination using spherical harmonics (SH) basis function. For instance, Tewari *et al.* [42] regress the illumination parameters from the input image, and the rendering loss is computed after combining the estimated illumination and skin reflectance in a self-supervised fashion. In this way, it is hard to analyze the quality of the estimated illumination and skin albedo. Moreover, low-order spherical harmonics cannot produce hard shadows cast from point light sources, which will decrease the face reconstruction quality in many real-world scenarios. I2ZNet [54] proposes a MOTC module to convert the color of the predicted texture to the in-the-wild texture of the input image. The MOTC can be viewed as a color correction matrix that corrects the white-balance between the two textures. However, the MOTC can only model the low-frequent lighting information, which will decrease the face registration performance if the lighting environment is complicated. Mean-

while, the surface information presented by a single image is limited, which leads to some artifacts such as over-smoothing and incorrect expression. To explicitly model the lighting, we propose a physics-based lighting model to learn the high-frequent lighting pattern (*e.g.*, shading, and brightness) from data captured in a lab-controlled environment. Besides, a domain adaptation schema is proposed to bridge the domain mismatch between lab and wild environments.

## 3. Adaptive Lighting Model

This section describes existing work on DAMs [27] (Sec. 3.1), construction of the lighting model with light-stage data (Sec. 3.2), and the adaptation of the model for in-the-wild settings in Sec. 3.3.

### 3.1. Deep Appearance Models

Our work is based on the face representation described in DAM [27]. It uses a variational auto-encoder (VAE) [20] to jointly represent the 3D facial geometry and appearance that are captured from a multi-view capture light-stage. The decoder, $D$, can generate instances of a person's face by taking, as input, a latent code $\mathbf{z}$, which encodes the expression, and a vector $\mathbf{v}^v$ that represents viewing direction as a normalized vector pointing from the center of the head to the camera $v$:

$$\hat{\mathbf{M}}, \hat{\mathbf{T}}^v \leftarrow D(\mathbf{z}, \mathbf{v}^v) \ . \tag{1}$$

Here, $\hat{\mathbf{M}}$ denotes the 3D face mesh (geometry) and $\hat{\mathbf{T}}^v$, the view-dependent texture.

In this work, we assume the availability of a pre-trained DAM decoder of a subject, and propose a system to fit the model to images by estimating the rigid and non-rigid motion (i.e., facial expression). A major challenge in implementing such a system is how to account for illumination differences between the high controlled studio lighting system where the avatar was captured, and the in-the-wild captures where the avatar is deployed. In the following, we will describe an adaptive lighting model that extends the original DAM formulation to enable high precision tracking in uncontrolled and complex environments.

### 3.2. Lighting Model

In order to incorporate a generative model of lighting into the DAM formulation, we extend the capture system in [27] to include 460 controllable lights that are synchronized with the multi-view camera system. The captured sequence was extended to include a portion where non-overlapping groups of approximately 10 lights were turned on, interleaved with fully lit frames that were used for tracking. This data was used to build a relightable face model using the scheme illustrated in Figure 2.

Our formulation is inspired by the light-varying residual proposed by Nestmeyer *et al.* [31], where illumination vari-
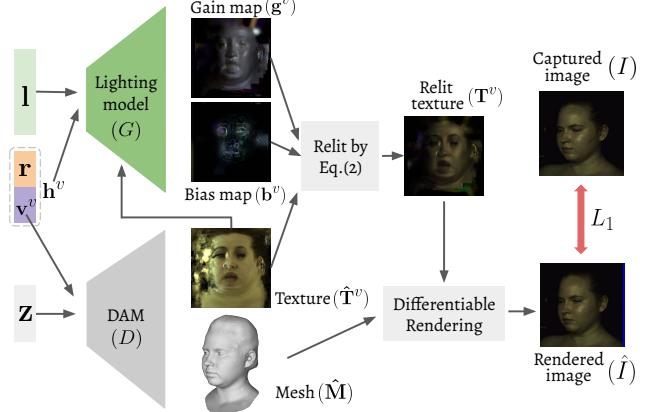


Figure 2. Training the lighting model on the light-stage data. We update the lighting model $G$ and per-frame expression code $\mathbf{z}$ while fixing the other parameters.

ations are represented using gain and bias maps, $\mathbf{g}$ and $\mathbf{b}$, each matching DAM's texture dimensions ($H \times W \times 3$):

$$\mathbf{T}^v = \hat{\mathbf{T}}^v \odot \mathbf{g}^v + \mathbf{b}^v, \tag{2}$$

where $\odot$ denotes the element-wise product, $\mathbf{T}^v$ is the relit texture, and $\hat{\mathbf{T}}^v$ is the DAM avatar's original texture (fully-lit). The gain and bias maps depend on the lighting, head pose[1], viewpoint, and expression. These inputs, represented by $\mathbf{l}$, $\mathbf{h}^v$, and $\hat{\mathbf{T}}^v$, are processed by MLPs and are spatially repeated for concatenation with the DAM's texture followed by additional convolution operations to produce the final relit texture. This lighting model is, thus, defined as follows:

$$\mathbf{g}^v, \mathbf{b}^v \leftarrow G(\mathbf{l}, \mathbf{h}^v, \hat{\mathbf{T}}^v; \phi) \ , \tag{3}$$

where $\phi$ denotes the weights of the network. Further details about $G$ are in Sec. 4.1 and in the supplementary material.

Since the goal of our lighting model is to enable accurate registration in uncontrolled scenarios, we do not require a lighting representation that is geometrically interpretable, only one that can span the space of facial illuminations. As such, we represent lighting conditions using a one vector that specifies which of the lights, in each color channel, are active for a given training frame. We use a binary vector of 150 dimensions, which comprises the 50 lighting groups with three color channels each. The fully-lit frames are encoded as the all-one vector. In combination with the continuously parameterized head pose, this representation allows for continuous and smoothly varying illumination synthesis on the face that can model complex effects such as how shadows move as the subject's head rotates in the scene.

---

[1]Here, the rigid head pose consists of two parts: rigid rotation $\mathbf{r} \in \mathbb{R}^3$ and camera viewpoint vector $\mathbf{v}^v \in \mathbb{R}^3$. Similar to [27], we assume that the viewpoint vector is relative to the rigid head orientation that is estimated from the tracking algorithm.

**Algorithm 1** Lighting Model Adaptation

**Input:** lighting model $G$ with weights $\phi$, $K$ key frames with initial face parameters $\{(I, \tilde{\mathbf{p}})_k\}$, camera viewpoint vector $\mathbf{v}^v$
**Output:** adapted lighting model $G$ and lighting code $\mathbf{l}$
**Initialization:** set $\mathbf{l}$ to zeros, $\phi$ to pre-trained weights by Sec. 3.2, face parameters $\{\mathbf{p}_k\}$ to $\{\tilde{\mathbf{p}}_k\}$
**for** number of iterations **do**
    **# Fitting** $\mathbf{l}$
    **for** number of iterations **do**
        Unfreeze $\mathbf{l}$, freeze $\phi$ and $\{\mathbf{p}_k\}$
        Calculate $\mathcal{L}_{pix}$ using Eq. 7
        $\mathbf{l} \leftarrow Adam\{\mathcal{L}_{pix}\}$
    **# Fitting** $\phi$ **and** $\{\mathbf{p}_k\}$
    **for** number of iterations **do**
        Freeze $\mathbf{l}$, unfreeze $\phi$ and $\{\mathbf{p}_k\}$
        Calculate $\mathcal{L}_{pix}$ using Eq. 7
        $\phi, \{\mathbf{p}_k\} \leftarrow Adam\{\mathcal{L}_{pix}\}$

To train $G$, we take a pre-trained DAM decoder and fix its weights while minimizing the reconstruction error over all camera views in the subject's sequence, that is:

$$\mathcal{L}_{\text{render}}(\phi, \mathbf{Z}) = \sum_{t,v} \left\| \left( I_t^v - \mathcal{R}(\mathbf{T}_t^v, \hat{\mathbf{M}}_t) \right) \odot m^v \right\|_1 , \quad (4)$$

where the minimization is performed over the lighting model's weights, $\phi$, as well as the expression codes for each frame $\mathbf{Z} = \{\mathbf{z}_t\}$. In Equation 4, $m^v$ is a foreground mask from view $v$ and $\mathcal{R}$ is a differentiable rasterization function [32]. To ensure stable convergence, we employ $L_2$-shrinkage on the expression codes, $\mathbf{Z}$, and use the technique in [52] to obtain a tracked mesh, $\mathbf{M}_t$, from the fully-lit frames, that are used to geometrically constrain the optimization, resulting in the following total objective:

$$\mathcal{L} = \mathcal{L}_{\text{render}} + \lambda_{\text{geo}} \sum_t \|\mathbf{M}_t - \hat{\mathbf{M}}_t\|^2 + \lambda_{\text{reg}}\|\mathbf{Z}\|^2. \quad (5)$$

Regularization weights of $\lambda_{\text{geo}} = 1.0$ and $\lambda_{\text{reg}} = 0.1$ were chosen for all experiments in Sec. 4 by cross validation.

### 3.3. Registration In-the-Wild

With the lighting model described in Sec. 3.2, we have a personalized face model that can synthesize realistic variations in expression as well as lighting suitable for an analysis-by-synthesis approach to register the model to in-the-wild video with uncontrolled illumination. However, to achieve robust and precise results, special care needs to be applied in how registration is performed. Our algorithm comprises three steps (Fig. 3) outlined below, each designed to address initialization, accuracy, and computational efficiency, respectively.

**Step 1: Initialization.** To avoid registration terminating in poor local minima, we initialize the pose and expression parameters by matching against facial keypoints in the image found via an off-the-shelf detector (*e.g.*,[22]). Specifically, face landmarks, $\{L_i\}$, describing facial features such as eye corners, lip contours, and the face silhouette correspond to fixed vertices in the face model's geometry, $\{\ell_i\}$. The initial face parameters, $\tilde{\mathbf{p}} = [\tilde{\mathbf{r}}, \tilde{\mathbf{t}}, \tilde{\mathbf{z}}]$, are then found by minimizing the reprojection error over these landmark points in all camera views, $v$, for every frame:

$$\mathcal{L}_{\text{land}}(\tilde{\mathbf{p}}) = \sum_{v,i} \left\| \Pi_v \left( \tilde{\mathbf{r}}\tilde{\mathbf{M}}^{(\ell_i)} + \tilde{\mathbf{t}} \right) - L_i^v \right\|^2 , \quad (6)$$

where $\Pi_v$ is the projection operator based on camera parameters that are assumed to be available. The face mesh, $\tilde{\mathbf{M}}$, is calculated using Eq. 1 with the expression code $\tilde{\mathbf{z}}$. Due to the landmarks' sparsity and detection noise, minimizing $\mathcal{L}_{\text{land}}$ results in only a rough alignment of the face in each frame (*e.g.*, Fig. 4 (b)). Nonetheless, it places the model within the vicinity of the solution, allowing the more elaborated optimization procedure described next to converge.

**Step 2: Lighting Model Adaptation.** Although the lighting model described in Sec. 3.2 equips us with the ability to synthesize variations in facial illumination, using the light-stage to simulate the total span of lighting variations encounter in the wild remains challenging. Effects such as nearfield lighting, cast shadows, and reflections from nearby objects are commonly observed in uncontrolled settings. To account for variations not spanned by our lighting model, $G$, in addition to solving for the lighting parameters, $\mathbf{l}$, we simultaneously fine-tune the model's weights, $\phi$, to obtain a better fit to in-the-wild images. Specifically, we minimize the following loss over a collection of $K$ frames [2] from the target environment:

$$\mathcal{L}_{\text{pix}}(\mathbf{l}, \phi, \{\mathbf{p}_k\}) = \sum_k \|r_k \odot w_k\|_1 + \lambda_{\triangle} \|\triangle r_k \odot w_k\|_1 , \quad (7)$$

where $r_k = I_k - \hat{I}(\mathbf{p}_k)$ is the reconstruction residual, $w_k$ is the foreground mask, and $\mathbf{p}_k = [\mathbf{r}_k, \mathbf{t}_k, \mathbf{z}_k]$. Here, $\triangle$ denotes the image Laplacian operator that makes the loss more robust to residual differences due to illumination and generally improves results.

**Step 3: Face Tracking.** The procedure described previously can generate accurate estimates of facial expression, but requires batch processing to adapt the lighting model simultaneously over several frames. However, once the lighting model has been adapted, the parameters $\mathbf{p} = [\mathbf{r}, \mathbf{t}, \mathbf{z}]$ for any new frames can be estimated independently of the lighting model $G$. Thus, in practice, we adapt the lighting model using only a small subset of $K$ frames and estimate

---

[2]We can select K frames out of the testing sequence for adaptation if the lighting of the K frames is the same as the testing sequence.

Figure 3. The pipeline of in-the-wild registration. We estimate the initial tracking parameters in **step 1**, adapt the lighting model and tracking parameters $l, \phi, \{\mathbf{p}_k\}$ with $K$ reference frames in **Step 2**, and further optimize the tracking parameters in **step 3**.
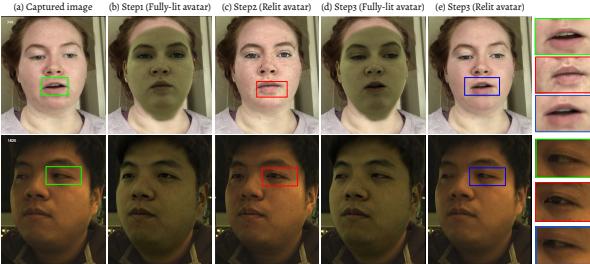


Figure 4. Result of each step in Sec. 3.3. From left to right: (a) captured image, (b) avatar output by **step 1**, (c) relit avatar after **step 2**, (d) fully-lit avatar after step 3, and (e) relit avatar after **step 3**. The last column shows the zoom-in of corresponding color rectangle, please notice the lip shapes and gaze directions.

$\mathbf{p}$ with the updated model $G$. To further improve accuracy, in addition to the optimizing the loss in Eq. 7, similarly to [5], we use dense optical flow [21] between the rendered model and the image to further constrain the optimization. It is computed over all projected mesh vertices in all camera views as follows:

$$\mathcal{L}_{\text{flow}}(\mathbf{p}) = \sum_{v,i} \left\| \left( \mathbf{r}\mathbf{M}^{(i)} + \mathbf{t} \right) - \Pi_v \left( \tilde{\mathbf{r}}\tilde{\mathbf{M}}^{(i)} + \tilde{\mathbf{t}} \right) - \mathbf{d}_i^v \right\|^2, \tag{8}$$

where mesh $\mathbf{M}$ and $\tilde{\mathbf{M}}$ are calculated using Eq. 1 with the latent face code $\mathbf{z}$ and $\tilde{\mathbf{z}}$ respectively, and $(\tilde{\mathbf{r}}, \tilde{\mathbf{t}}, \tilde{\mathbf{z}})$ are initial parameters from **Step 1**. We optimize the per-frame face parameters $\mathbf{p} = \{\mathbf{r}, \mathbf{t}, \mathbf{z}\}$ by minimizing the total loss $\mathcal{L} = \mathcal{L}_{\text{pix}} + \lambda_{\text{flow}}\mathcal{L}_{\text{flow}}$, where $\lambda_{\text{flow}} = 3.0$ was chosen via cross validation. Fig. 4 (d-e) shows some examples of results obtained through this process, demonstrating accurate alignment and reconstruction of lip shape and gaze direction that were absent in **Step 1**.

## 4. Experiments

In this section, we conduct quantitative and qualitative experiments to show the performance of the proposed face tracking framework on *in-the-wild* videos. Sec. 4.1 explains dataset and implementation details. Sec. 4.2 compares our method to state-of-the-art methods, and Sec. 4.3 describes the ablation studies.

### 4.1. Experimental Settings

**Dataset Collection.** We recorded our light-stage data in a calibrated multi-view light-stage consisting of 40 machine vision cameras capable of synchronously capturing HDR images at $1334 \times 2048$ / 90 fps and a total of 460 white LED lights. We flash a group of LEDs (at most 10) per frame and instruct our subjects to make diverse expressions with head movements. There are 50 different lighting patterns and one fully-lit pattern. We record a total of 13 minutes video sequence of one subject (see supplementary videos).

The in-the-wild video test were gathered using the frontal camera of an iPhone. We captured videos for 10 subjects. We collected around 5 video clips for each subject, performing different facial expressions and head movements, under various lighting conditions and environments.

**Implementation Details.** The light-stage training step (Sec. 3.2) and the lighting adaptation step cost 36 hours and 4mins, respectively, on an NVIDIA DGX machine. In all our experiments, we used the Adam optimizer [18] to optimize the losses. In order to cover more lighting space during the training of the lighting transfer module, we augmented the RGB channels of the light-stage lighting color with randomly selected scales, which are also used to scale the lighting code $l$ in the training data. The architecture of our lighting decoder $G$ is as follows: we first encode the input head pose $\mathbf{h}^v$ and $l$ with two MLPs to 256 dimensions. After concatenating the two latent features, we pass it to a fully-connected layer and a convolution layer followed by four transposed convolutions with each layer. The fully-lit texture $\hat{\mathbf{T}}^v$ is encoded by four convolution layers with each layer followed by a down-sampling layer to texture feature. Then we concatenate the texture feature and lighting feature and pass it to two separate branches consists of two transpose convolution. The two branches output gain map $\mathbf{g}_t^v$ and bias map $\mathbf{b}_t^v$ with the resolution of $256 \times 256$, and we up-sample them 4 times using bilinear interpolation to the same resolution as texture $\hat{\mathbf{T}}^v$. Please refer to supplementary materials for details. While training the lighting model on the light-stage data in Sec. 3.2, the rendering loss is optimized with an initial learning rate of $1e^{-3}$, which is decreased by a quarter after every 10 epochs.

During the registration of the *in-the-wild* videos in Sec. 3.3, we fit the DAM code $\mathbf{z}_k$ in **step 1** with 1000 it-
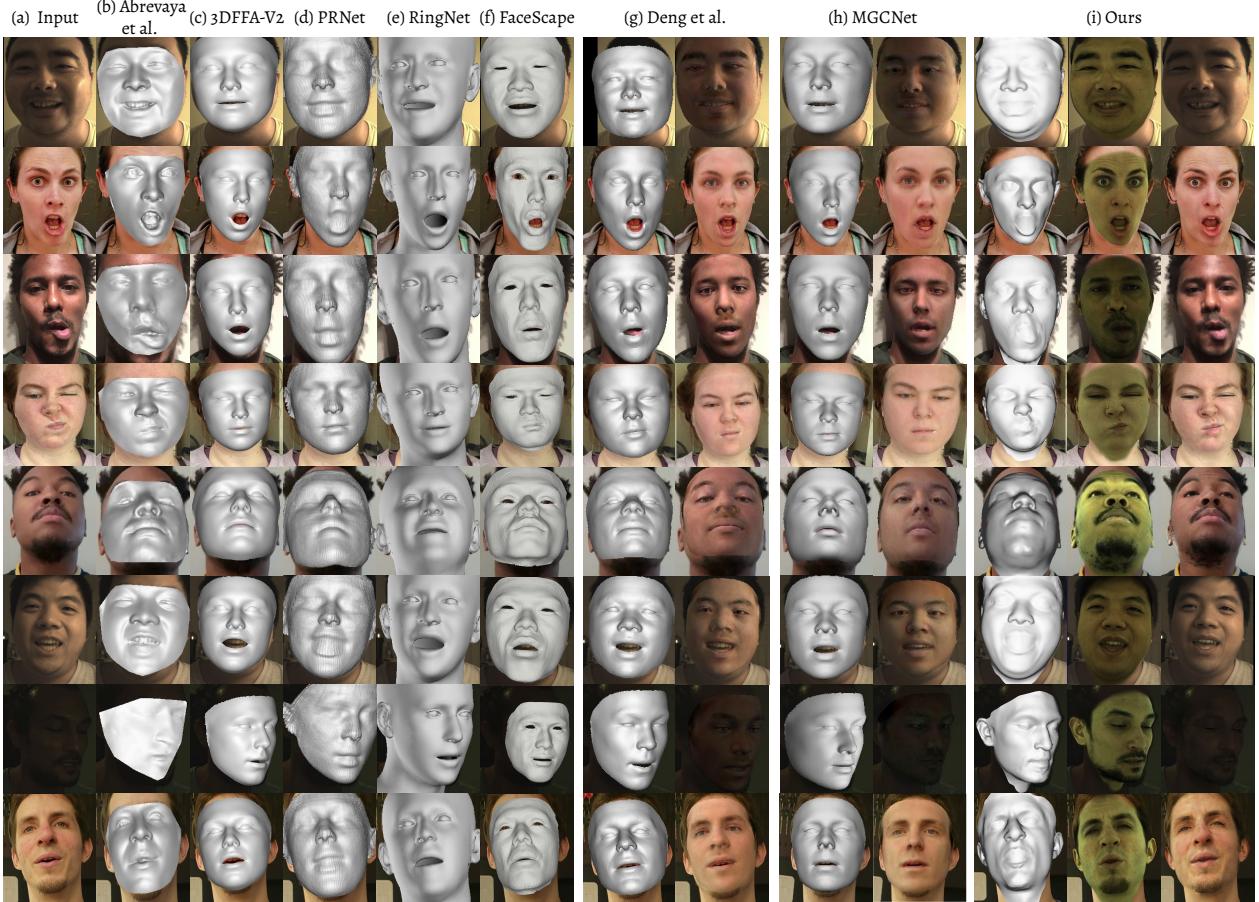
Figure 5. Qualitative comparison. We suggest to view it using a monitor for better visual quality. We selected 8 subjects from the test set with different lighting conditions, facial expression, and head motion. From left to right: (a) Captured image, (b) Abrevaya *et al.* [1], (c) 3DDFA$v$2 [15], (d) PRNet [11], (e) RingNet [37], (f) FaceScape [53], (g) Deng *et al.* [10], (h) MGCNet [40], and (i) our method.
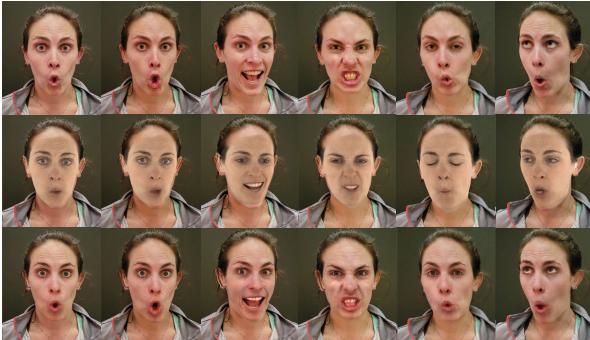


Figure 6. Visual comparison between our method and I2ZNet [54] on testing video frames. From top to bottom: captured image, I2ZNet [54], and our method.

erations. The initial learning rate is $0.1$, and we decrease it by half after every 500 iterations. In **step 2**, the $K$ frames are uniformly sampled according to the value of $\tilde{\mathbf{r}}$ to cover diverse head movements. We fit the lighting code $\mathbf{l}$ with 250 iterations with a learning rate of $1e^{-2}$, then update $G$'s network parameter $\phi$ and face parameters $\{\mathbf{p}_k\}$ with 500

iterations with a learning rate of $1e^{-3}$. We alternatively update $\mathbf{l}$, $\phi$, and $\{\mathbf{p}_k\}$ for $4$ times. In **step 3**, we update the face parameters $\{\mathbf{p}_k\}$ with the same hyper-parameters as **step 1**.

**Evaluation Metrics.** We used a variety of perceptual measures to quantitatively compare the registered image against the ground-truth in-the-wild image. Besides the pixel-level $L_2$ distance, we adopted PSNR and structural similarity (SSIM) [51] for human perceptual response. To evaluate the realism of the output avatar, we computed the cosine similarity (CSIM) between embedding vectors of the state-of-the-art face recognition network [9] suggested by [55, 7] for measuring identity mismatch between the input image and reconstructed avatar.

## 4.2. Comparison with state-of-the-art methods

To demonstrate the effectiveness of our lighting model and the face tracking quality, we compared our algorithm against the following set of related methods using the pre-trained models provided by the authors: Abrevaya *et al.* [1], 3DDFA$v$2 [15], PRNet [11], RingNet [37], FaceScape [53],

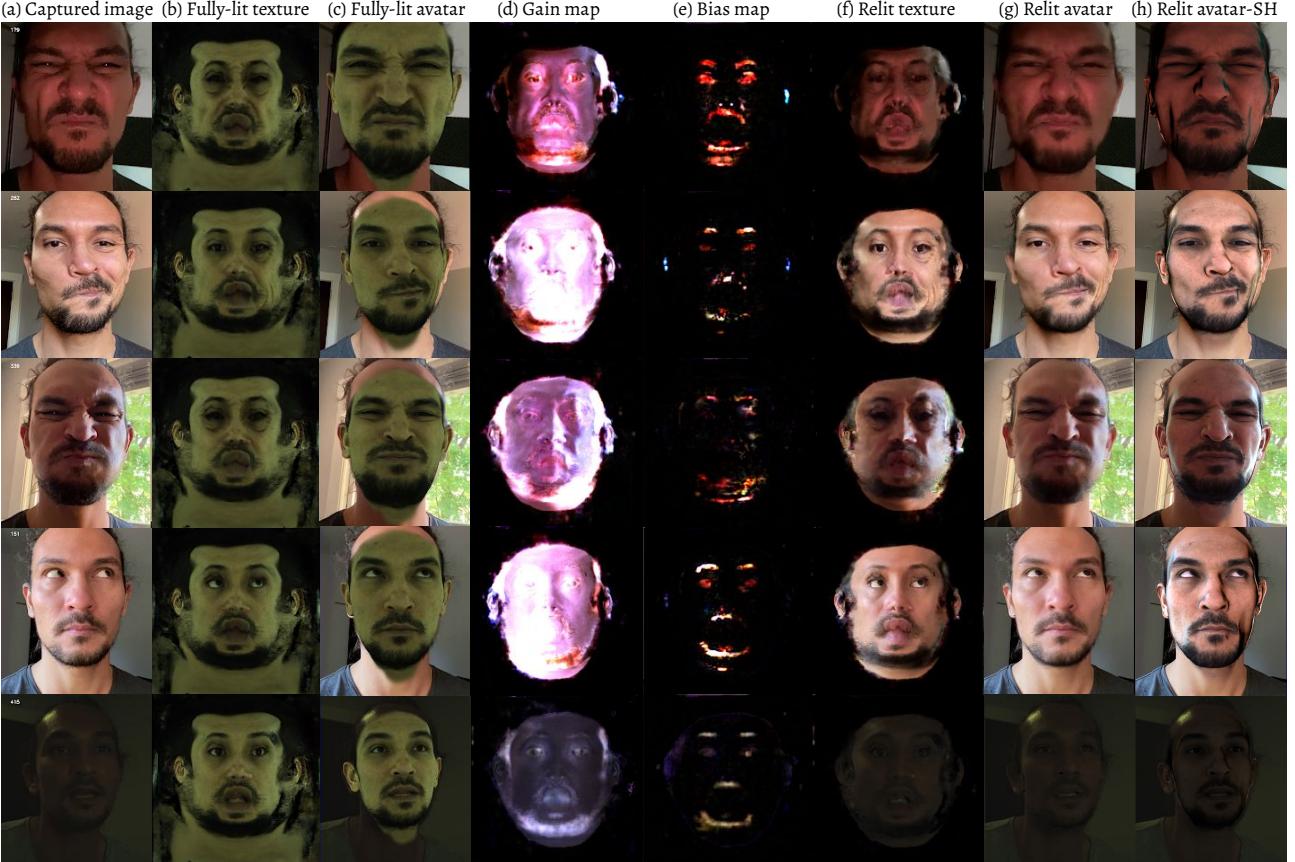| (a) Captured image | (b) Fully-lit texture | (c) Fully-lit avatar | (d) Gain map | (e) Bias map | (f) Relit texture | (g) Relit avatar | (h) Relit avatar-SH |

Figure 7. Visual results under different lighting conditions. Besides our final relit avtar (g), we also show the tracked avatar (h) using spherical harmonics (SH) as illumination model. Our method can handle different lighting conditions well.

Deng *et al*. [10], MGCNet [40], and I2ZNet [54]. As the baseline method, we used the same face registration method proposed in Sec. 3, but use the standard spherical harmonics (SH) [30] illumination model, which is denoted as ours-SH. We adopt the same parameter setting as [43] for SH, and train a regression network to regress the input images to the 27 dimensional SH parameters.

**Qualitative Evaluation.** Fig. 5 shows the tracked geometry results for different methods [3]. We can observe that our proposed method is robust to diverse facial expressions, poses, and lighting conditions. For example, in the $3^{rd}$ and $4^{th}$ row, all other methods failed to describe the expression (*e.g*., lips) except our method and Abrevaya *et al*. [1]. In the $7^{th}$ row, all other methods can not output the correct head pose, while our method can still reconstruct high-quality geometry and texture under dark lighting conditions. We also show the reconstructed avatars by (g) [10], (h) [40], and (i) our method (fully-lit and relit avatar). Comparing with other methods, our method not only generates a more realistic avatar, but also considers the lighting details (*e.g*., shadows and specular highlights, see the relit avatar in $1^{th}$, $3^{rd}$,

---

[3]Note that our method relies on person-specific DAM.

and $6^{th}$ row). Furthermore, we compare our relit avatar with I2ZNet [54] in Fig. 6. Although I2ZNet is a person-specific model, our method produces better visual results.

Fig. 7 shows the visual results under different real-world lighting environments. The proposed method is robust to different unseen lighting conditions, and our face tracking system can output a high-quality avatar with the aid of our lighting model. Fig. 7(h) shows the tracking results using the SH illumination model. The reconstruction error between the captured frame and the avatar relit by SH model is large, and it decreases the face tracking performance.

We also show our avatar rendered from different viewpoints in Fig. 1 and Fig. 8. We can find that our method can output the high-fidelity avatar from any viewpoint. Our lighting model is conditioned on the camera viewpoints, so the gain map and bias map will be adjusted to match the lighting in the specific view. Please refer to the supplementary video for more visual results.

**Quantitative Evaluation.** Tab. 1 shows the quantitative results of MGCNet [40], Deng *et al*. [10], I2ZNet [54], and our method. Our method outperforms other methods not only at the human perception level, but also at the identity preserving ability (CSIM score). The high CSIM score

(a) Captured image (b) Normal map (c) Fully-lit texture (d) Fully-lit avatar (e) Gain map (f) Relit texture (g) Relit avatar
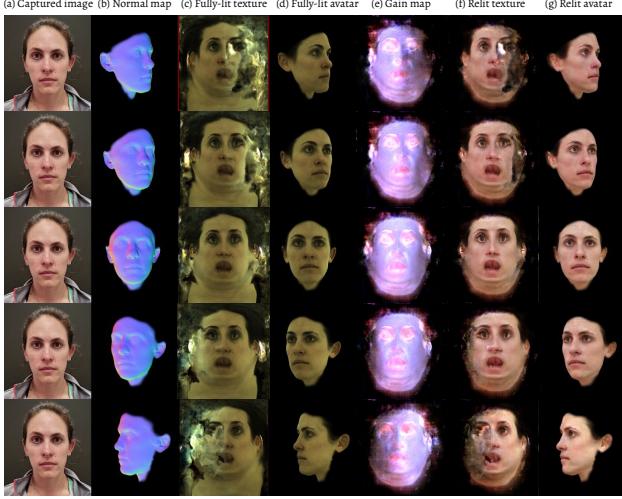
Figure 8. The visual results for different viewpoints. From left to right: captured image, normal map, fully-lit texture, fully-lit avatar, gain map, relit texture, and relit avatar. From top to bottom: different viewpoints. Please notice the changes of the gain map due to different viewpoints.

| Methods | Metrics | | | |
|---|---|---|---|---|
| | $L_2 \downarrow$ | SSIM↑ | PSNR↑ | CSIM↑ |
| MGCNet [40] | 24.56 | 0.86 | 34.47 | 0.305 |
| Deng *et al*. [10] | 27.33 | 0.86 | 33.81 | 0.434 |
| I2ZNet [54] | 20.49 | 0.927 | 35.03 | 0.592 |
| Ours-SH | 41.01 | 0.79 | 32.00 | 0.693 |
| Ours | **12.59** | **0.93** | **37.93** | **0.871** |

Table 1. The quantitative evaluation on the test set. ↓/↑ denote the lower/higher, the better. The top-1 scores are highlighted.

indicates that our reconstructed avatar produces high identity similarity. We can see that the $L_2$ and other perceptual scores of the avatar optimized with the SH illumination model result in higher error, although it preserves the identity information (CSIM).

## 4.3. Ablation Study

We have already shown the superiority of our face tracking algorithm along with the lighting model over other methods. To further demonstrate the effectiveness of different steps in Sec. 3, we make a comprehensive ablation study on a subset of the testing videos.

**The lighting model.** To evaluate the role of our lighting model, we test our face tracking system with the lighting model under two different settings: without pre-training on light-stage data (w/o Pre-train.) and without lighting adaptation on wild video frames (w/o Adapt.). We set $K = 48$ in the without pre-training experiment. Tab. 2 shows the quantitative evaluation results, and we can see that all the scores are improved with light-stage pre-training ($1^{st}$ and $5^{th}$ row). Although without adaptation ($2^{nd}$ row) performs

| Methods | Metrics | | | |
|---|---|---|---|---|
| | $L_2 \downarrow$ | SSIM↑ | PSNR↑ | CSIM↑ |
| w/o Pre-train. | 14.88 | 0.941 | 37.08 | 0.783 |
| w/o Adapt. | 18.03 | 0.938 | 35.59 | 0.593 |
| $K = 1$ | 18.58 | 0.900 | 35.48 | 0.550 |
| $K = 12$ | 12.30 | 0.950 | 37.32 | 0.842 |
| $K = 48$ | 11.35 | 0.950 | 37.66 | 0.809 |
| $K = $ All | 11.33 | 0.952 | 37.78 | 0.863 |

Table 2. Quantitative results of the ablation study.



(a) Captured image (b) w/o pre-training (c) w/o adaptation (d) K = 1 (e) K = 12 (f) K = 48 (g) K = All

Figure 9. Visual comparison of the ablation study.

slightly better than $K = 1$, with more reference frames for adaptation, the adapted lighting model performs much better than without adaptation. Fig. 9 shows the visual comparison. We can see that the pre-training on light-stage data and lighting adaptation both contribute to the final tracking results, where the pre-training on light-stage data enables the lighting interpretation ability, and the adaptation step enables the lighting model to generate accurate gain and bias map for target video frames.

**Influence of the value of $K$.** To evaluate the effect on the number of reference frames used in the lighting model adaptation step, we sample reference frames with different $K$ and keep other settings the same, and show the results in Tab. 2 and Fig. 9. We can find that with only 48 reference frames, the lighting model can be perfectly adapted to the target video and achieve comparable results in visual metrics ($L_2$, SSIM, and PSNR). If the amount of the selected reference frames is too small (*e.g.*, $K = 1$), the lighting model will be over-fitted to the selected frames and loss the lighting interpolation ability.

## 5. Conclusion

We present a new in-the-wild face tracking algorithm with an adaptive lighting model that can infer a high-fidelity 3D avatar. The proposed lighting model inherits the prior knowledge from the light-stage data, and it adapts to in-the-wild samples to produce high-quality re-lighting results for our face tracking. Results confirm that relatively few adaption samples (48) are enough to produce hyper-realistic results in the avatar. While the proposed method can generate photo-realistic avatars from videos in the wild, our lighting model assumes that the lighting source in the testing video is fixed, and it uses one single lighting code to represent the lighting in the whole video sequence. This is a limitation of the current model, and the model can produce undesirable results if the lighting is changing in the video. In the future,

we will explore on-line adaptation methods to address this limitation of the current work.

# References

[1] V. F. Abrevaya, A. Boukhayma, P. H. Torr, and E. Boyer. Cross-modal deep face normals with deactivable skip connections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4979–4989, 2020. 6, 7

[2] T. Beeler, F. Hahn, D. Bradley, B. Bickel, P. Beardsley, C. Gotsman, R. W. Sumner, and M. Gross. High-quality passive facial performance capture using anchor frames. *ACM Trans. Graph.*, 30:75:1–75:10, August 2011. 2

[3] V. Blanz and T. Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, pages 187–194, 1999. 2

[4] J. Booth, A. Roussos, A. Ponniah, D. Dunaway, and S. Zafeiriou. Large scale 3d morphable models. *International Journal of Computer Vision*, 126(2-4):233–254, 2018. 2

[5] C. Cao, M. Chai, O. Woodford, and L. Luo. Stabilized real-time face tracking via a learned dynamic rigidity prior. *ACM Transactions on Graphics (TOG)*, 37(6):1–11, 2018. 5

[6] C. Cao, Y. Weng, S. Zhou, Y. Tong, and K. Zhou. Facewarehouse: A 3d facial expression database for visual computing. *IEEE Transactions on Visualization and Computer Graphics*, 20(3):413–425, 2013. 2

[7] L. Chen, G. Cui, C. Liu, Z. Li, Z. Kou, Y. Xu, and C. Xu. Talking-head generation with rhythmic head motion. In *European Conference on Computer Vision*, pages 35–51. Springer, 2020. 6

[8] T. F. Cootes, G. J. Edwards, and C. J. Taylor. Active appearance models. *IEEE Transactions on pattern analysis and machine intelligence*, 23(6):681–685, 2001. 2

[9] J. Deng, J. Guo, N. Xue, and S. Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2019. 6

[10] Y. Deng, J. Yang, S. Xu, D. Chen, Y. Jia, and X. Tong. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. 6, 7, 8

[11] Y. Feng, F. Wu, X. Shao, Y. Wang, and X. Zhou. Joint 3d face reconstruction and dense alignment with position map regression network. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 534–551, 2018. 6

[12] P. Garrido, M. Zollhöfer, D. Casas, L. Valgaerts, K. Varanasi, P. Pérez, and C. Theobalt. Reconstruction of personalized 3d face rigs from monocular video. *ACM Transactions on Graphics (TOG)*, 35(3):1–15, 2016. 2

[13] B. Gecer, S. Ploumpis, I. Kotsia, and S. Zafeiriou. Ganfit: Generative adversarial network fitting for high fidelity 3d face reconstruction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1155–1164, 2019. 2

[14] T. Gerig, A. Morel-Forster, C. Blumer, B. Egger, M. Luthi, S. Schönborn, and T. Vetter. Morphable face models-an open framework. In *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*, pages 75–82. IEEE, 2018. 2

[15] J. Guo, X. Zhu, Y. Yang, F. Yang, Z. Lei, and S. Z. Li. Towards fast, accurate and stable 3d dense face alignment. *arXiv preprint arXiv:2009.09960*, 2020. 2, 6

[16] P. Huber, G. Hu, R. Tena, P. Mortazavian, P. Koppen, W. J. Christmas, M. Ratsch, and J. Kittler. A multiresolution 3d morphable face model and fitting framework. In *Proceedings of the 11th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications*, 2016. 2

[17] H. Kim, M. Zollöfer, A. Tewari, J. Thies, C. Richardt, and T. Christian. InverseFaceNet: Deep Single-Shot Inverse Face Rendering From A Single Image. In *Proceedings of Computer Vision and Pattern Recognition (CVPR 2018)*, 2018. 2

[18] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5

[19] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2

[20] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *2nd International Conference on Learning Representations*, 2014. 3

[21] T. Kroeger, R. Timofte, D. Dai, and L. Van Gool. Fast optical flow using dense inverse search. In B. Leibe, J. Matas, N. Sebe, and M. Welling, editors, *Computer Vision – ECCV 2016*, pages 471–488, Cham, 2016. Springer International Publishing. 5

[22] D. Lee, H. Park, and C. D. Yoo. Face alignment using cascade gaussian process regression trees. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4204–4212, 2015. 4

[23] T. Li, T. Bolkart, M. J. Black, H. Li, and J. Romero. Learning a model of facial shape and expression from 4d scans. *ACM Trans. Graph.*, 36(6):194–1, 2017. 2

[24] J. Lin, Y. Yuan, T. Shao, and K. Zhou. Towards high-fidelity 3d face reconstruction from in-the-wild images using graph convolutional networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5891–5900, 2020. 2

[25] C. Liu and X. Li. Superimposition-guided facial reconstruction from skull. *Arxiv*, arXiv:1810.00107, 2018. 2

[26] C. Liu, Z. Li, S. Quan, and Y. Xu. Lighting estimation via differentiable screen-space rendering. In *2020 IEEE Conference on Virtual Reality and 3D User Interfaces Abstracts and Workshops (VRW)*, pages 575–576, 2020. 2

[27] S. Lombardi, J. Saragih, T. Simon, and Y. Sheikh. Deep appearance models for face rendering. *ACM Transactions on Graphics (TOG)*, 37(4):1–13, 2018. 2, 3

[28] I. Matthews and S. Baker. Active appearance models revisited. *International journal of computer vision*, 60(2):135–164, 2004. 2

[29] A. Meka, R. Pandey, C. Häne, S. Orts-Escolano, P. Barnum, P. David-Son, D. Erickson, Y. Zhang, J. Taylor, S. Bouaziz, et al. Deep relightable textures: volumetric performance capture with neural rendering. *ACM Transactions on Graphics (TOG)*, 39(6):1–21, 2020. 2

[30] C. Müller. *Spherical harmonics*, volume 17. Springer, 2006. 7

[31] T. Nestmeyer, J.-F. Lalonde, I. Matthews, and A. Lehrmann. Learning physics-guided face relighting under directional light. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5124–5133, 2020. 3

[32] S. G. Parker, J. Bigler, A. Dietrich, H. Friedrich, J. Hoberock, D. Luebke, D. McAllister, M. McGuire, K. Morley, A. Robison, et al. Optix: a general purpose ray tracing engine. *Acm transactions on graphics (tog)*, 29(4):1–13, 2010. 4

[33] P. Paysan, R. Knothe, B. Amberg, S. Romdhani, and T. Vetter. A 3d face model for pose and illumination invariant face recognition. In *2009 Sixth IEEE International Conference on Advanced Video and Signal Based Surveillance*, pages 296–301. Ieee, 2009. 2

[34] E. Richardson, M. Sela, and R. Kimmel. 3d face reconstruction by learning from synthetic data. In *2016 fourth international conference on 3D vision (3DV)*, pages 460–469. IEEE, 2016. 2

[35] J. Roth, Y. Tong, and X. Liu. Unconstrained 3d face reconstruction. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2606–2615, 2015. 2

[36] J. Roth, Y. Tong, and X. Liu. Adaptive 3d face reconstruction from unconstrained photo collections. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(11):2127–2141, December 2016. 2

[37] S. Sanyal, T. Bolkart, H. Feng, and M. J. Black. Learning to regress 3d face shape and expression from an image without 3d supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7763–7772, 2019. 2, 6

[38] M. Sela, E. Richardson, and R. Kimmel. Unrestricted facial geometry reconstruction using image-to-image translation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1576–1585, 2017. 2

[39] S. Sengupta, A. Kanazawa, C. D. Castillo, and D. W. Jacobs. Sfsnet: Learning shape, reflectance and illuminance of facesin the wild'. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6296–6305, 2018. 2

[40] J. Shang, T. Shen, S. Li, L. Zhou, M. Zhen, T. Fang, and L. Quan. Self-supervised monocular 3d face reconstruction by occlusion-aware multi-view geometry consistency. *arXiv preprint arXiv:2007.12494*, 2020. 6, 7, 8

[41] A. Tewari, F. Bernard, P. Garrido, G. Bharaj, M. Elgharib, H.-P. Seidel, P. Pérez, M. Zollhöfer, and C. Theobalt. FML: Face Model Learning from Videos. In *Proceedings of Computer Vision and Pattern Recognition (CVPR 2019)*, 2019. 2

[42] A. Tewari, M. Zollhöfer, P. Garrido, F. Bernard, H. Kim, P. Pérez, and C. Theobalt. Self-supervised multi-level face model learning for monocular reconstruction at over 250 hz. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2549–2559, 2018. 2

[43] A. Tewari, M. Zollhofer, H. Kim, P. Garrido, F. Bernard, P. Perez, and C. Theobalt. Mofa: Model-based deep convolutional face autoencoder for unsupervised monocular reconstruction. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 1274–1283, 2017. 2, 7

[44] J. Thies, M. Zollhöfer, M. Nießner, L. Valgaerts, M. Stamminger, and C. Theobalt. Real-time expression transfer for facial reenactment. *ACM Transactions on Graphics (TOG)*, 34(6), 2015. 2

[45] J. Thies, M. Zollhöfer, M. Stamminger, C. Theobalt, and M. Niessner. Face2face: Real-time face capture and reenactment of rgb videos. *Commun. ACM*, 62(1):96–104, Dec. 2018. 2

[46] L. Thies, M. Zollhoefer, C. Richardt, C. Theobalt, and G. Greiner. Real-time halfway domain reconstruction of motion and geometry. In *Proceedings of the International Conference on 3D Vision (3DV)*, 2016. 2

[47] L. Tran, F. Liu, and X. Liu. Towards high-fidelity nonlinear 3d face morphable model. In *In Proceeding of IEEE Computer Vision and Pattern Recognition*, Long Beach, CA, June 2019. 2

[48] L. Tran and X. Liu. On learning 3d face morphable model from in-the-wild images. *IEEE transactions on pattern analysis and machine intelligence*, 2019. 2

[49] G. Tzimiropoulos, J. Alabort-i Medina, S. Zafeiriou, and M. Pantic. Generic active appearance models revisited. In *Asian Conference on Computer Vision*, pages 650–663. Springer, 2012. 2

[50] X. Wang, Y. Guo, B. Deng, and J. Zhang. Lightweight photometric stereo for facial details recovery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 740–749, 2020. 2

[51] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 6

[52] C. Wu, T. Shiratori, and Y. Sheikh. Deep incremental learning for efficient high-fidelity face tracking. *ACM Trans. Graph.*, 37(6), Dec. 2018. 4

[53] H. Yang, H. Zhu, Y. Wang, M. Huang, Q. Shen, R. Yang, and X. Cao. Facescape: A large-scale high quality 3d face dataset and detailed riggable 3d face prediction. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 6

[54] J. S. Yoon, T. Shiratori, S.-I. Yu, and H. S. Park. Self-supervised adaptation of high-fidelity face models for monocular performance tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4601–4609, 2019. 2, 6, 7, 8

[55] E. Zakharov, A. Shysheya, E. Burkov, and V. Lempitsky. Few-shot adversarial learning of realistic neural talking head models. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9459–9468, 2019. 6

[56] H. Zhou, J. Liu, Z. Liu, Y. Liu, and X. Wang. Rotate-and-render: Unsupervised photorealistic face rotation from single-view images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[57] H. Zhou, Y. Sun, W. Wu, C. C. Loy, X. Wang, and Z. Liu. Pose-controllable talking face generation by implicitly modularized audio-visual representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[58] X. Zhu, X. Liu, Z. Lei, and S. Z. Li. Face alignment in full pose range: A 3d total solution. *IEEE transactions on pattern analysis and machine intelligence*, 41(1):78–92, 2017. 2

[59] M. Zollhöfer. *Real-Time Reconstruction of Static and Dynamic Scenes*. PhD thesis, University of Erlangen-Nuremberg, 2015. 2
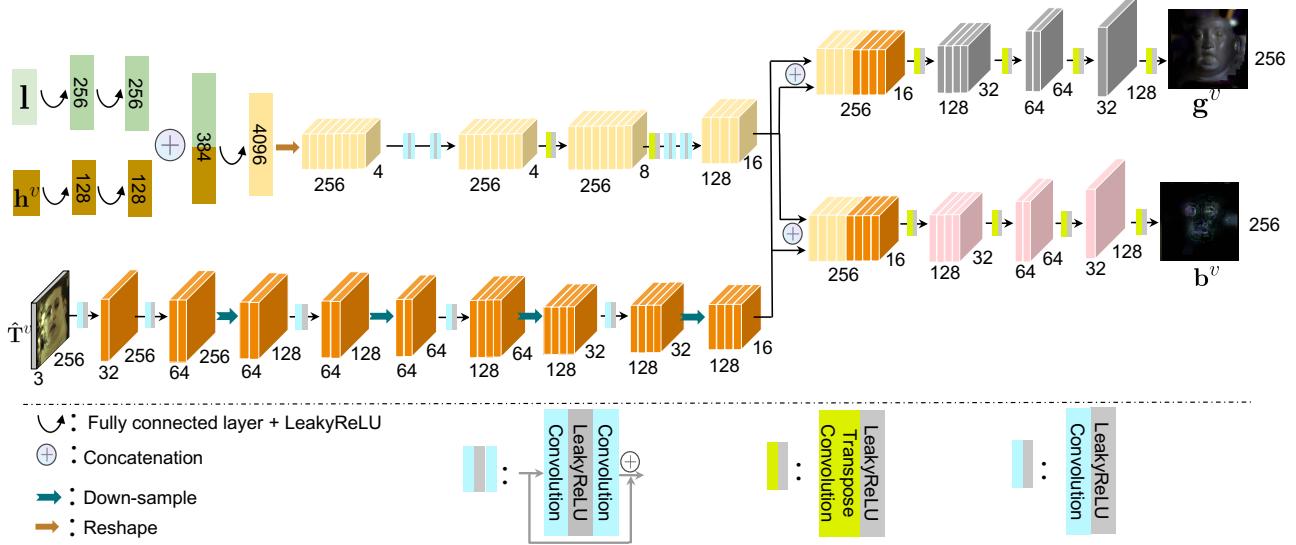
Figure 10. The detailed network structure of our lighting transfer network ($G$).

## Supplemental Materials

## A. Network Structure

We present the detailed network structure in Fig. 10.

## B. Inputs and Outputs

The lighting code $\mathbf{l}$ is a pre-defined vector when we train $G$ on light-stage data, and is a learnable vector when we refine $G$ on in-the-wild video frames. During training on light-stage data, the lighting direction is encoded by the position of the non-zero element in $\mathbf{l}$, and the lighting color is encoded by the value of the non-zero element in $\mathbf{l}$. The view-dependent head pose $\mathbf{h}^v \in \mathbb{R}^6 = \{\mathbf{r}_x, \mathbf{r}_y, \mathbf{r}_z, \mathbf{v}_x^v, \mathbf{v}_y^v, \mathbf{v}_z^v\}$, where $\mathbf{r} = \{\mathbf{r}_x, \mathbf{r}_y, \mathbf{r}_z\}$ and $\mathbf{v}^v = \{\mathbf{v}_x^v, \mathbf{v}_y^v, \mathbf{v}_z^v\}$ are rigid head rotation and viewpoint vector, respectively. The fully-lit texture $\hat{\mathbf{T}}^v$ is obtained from DAM decoder, and we down-sample it to the size of $3 \times 256 \times 256$.

The outputs are the gain and bias map $\mathbf{g}^v, \mathbf{b}^v$, and we upsample the output $\mathbf{g}^v, \mathbf{b}^v$ back to the size of $3 \times 1024 \times 1024$ by bilinear interpolation.