

SAFA: Structure Aware Face Animation

Qiulin Wang
JD Technology

wangqiulin5@jd.com

Lu Zhang
JD Technology

zhanglu32@jd.com

Bo Li
JD Technology

prclibo@gmail.com

Abstract

Recent success of generative adversarial networks (GAN) has made great progress on the face animation task. However, the complex scene structure of a face image still makes it a challenge to generate videos with face poses significantly deviating from the source image. On one hand, without knowing the facial geometric structure, generated face images might be improperly distorted. On the other hand, some area of the generated image might be occluded in the source image, which makes it difficult for GAN to generate realistic appearance. To address these problems, we propose a structure aware face animation (SAFA) method which constructs specific geometric structures to model different components of a face image. Following the well recognized motion based face animation technique, we use a 3D morphable model (3DMM) to model the face, multiple affine transforms to model the other foreground components like hair and beard, and an identity transform to model the background. The 3DMM geometric embedding not only helps generate realistic structure for the driving scene, but also contributes to better perception of occluded area in the generated image. Besides, we further propose to exploit the widely studied inpainting technique to faithfully recover the occluded image area. Both quantitative and qualitative experiment results have shown the superiority of our method. Code is available at <https://github.com/Qiulin-W/SAFA>.

1. Introduction

Face animation refers to the task to generate a face video from a source face image to reenact the face in a given driving video, possibly of different identity from the source. The source image provides the appearance and identity information and the driving video determines face poses and translational motions. Generally, a generated face video with high quality should meet multiple requirements including but not limited to the following:

(1) **pose-preserving**: The generated face video need accurately follow the specific pose, expression and motion of

the driving video.

(2) **identity-preserving**: The generated face video should preserve the identity information of the source face.

(3) **realistic**: The generated face images should possess realistic face shape, texture and visual sharpness.

(4) **occlusion aware**: When an area in the generated image is occluded in the source image, the model should be aware of the occlusion and recover the missing part. Such occlusion might occur on the foreground face when a head turns or on the background when a head undergoes translational move.

With the rapid development of Generative Adversarial Networks (GANs) [9], recent deep learning based methods [39, 15, 36, 42, 41, 29, 10, 7, 38] have shown promising results in generating face animation videos. However, it is still challenging to meet all the four requirements stated above especially when large face pose variation, translational/expressional motion and occlusion occur.

In this paper, we resolve the above challenges by integrating the scene structure knowledge in face animation procedure. The scene structure can be decomposed into foreground and background. The foreground refers to the head which can be further decomposed into face and other components like hair, beard and clothes. The face region is conventionally considered as a rigid body in computer vision while others are relatively soft. The background is static in the typical scenario of face animation. It is recognized to be difficult to learn this complex structure elaborately in the end-to-end manner and specific priors is essential for high quality generation. [29] proposes a First Order Motion Model (FOMM) to model foreground and background motion by affine motion and static motion, respectively. However, when face pose or expression changes significantly between source and driving images, 2D affine motion is insufficient to model such transforms.

To address this problem, we integrate the state-of-the-art 3D morphable model (3DMM) FLAME [16] into the structure model of FOMM to model the face region. FLAME models not only the face shape but also the expressional motion of eyelids, mouth and jaw, which provides strong prior knowledge of the face structure. This knowledge fur-

ther helps accurately distinguish the occluded area in the generated images. We further exploit the successful inpainting technique, contextual attention module [40], to recover the appearance feature at the perceived occluded area that is highly relevant to other unoccluded feature patches. In addition, we propose a novel Geometrically-Adaptive (DE)normalization (GADE) layer to integrate the 3DMM geometric embeddings with the recovered appearance, which is shown to further enhance the generation quality in terms of facial details.

We conduct extensive experiments to compare our method with state-of-the-art methods [36, 42, 41, 29]. Both qualitative and quantitative results show that our approach outperforms previous works, especially in synthesizing faces in case of large pose differences, translational/expressional motions and occlusions.

In summary, our contributions include:

- We combine the 2D affine motion model with 3DMM to elaborately model the scene structure in face images.
- We propose to exploit inpainting techniques to recover appearance feature at the perceived occluded area. To the best of our knowledge, this is the first combination of the two computer vision subjects of face animation and inpainting.
- We propose a novel geometrically-adaptive denormalization layer to make full use of the 3D face geometry information for generating realistic facial details.

2. Related Work

2.1. 3D morphable model (3DMM)

3D morphable model (3DMM) is a statistical model which transforms the shape and texture of a 3D face into a vector space representation [1]. Traditional 3DMMs [1, 22, 4, 2, 16] are low-dimensional linear models generated through principal component analysis (PCA). A 3D face is represented by a linear combination of those orthogonal bases with largest eigenvalues. The state-of-the-art linear model, FLAME (Faces Learned with an Articulated Model and Expressions) [16] is learned from about 3800 3D face scans. Apart from shape and expression, FLAME also models four articulated joints (jaw, neck, eyeballs) and pose-dependent corrective blendshapes, which enables more flexible and expressive full head modeling. In this paper, we adopt FLAME as a decoder to model 3D faces and use a convolutional neural network as an encoder to estimate 3DMM parameters from face images.

2.2. Face animation

Recent works on face animation can be divided into three categories: (1) motion-based methods, (2) landmark-based methods and (3) 3D-based methods.

Motion-based methods explicitly models a relative motion field from source to driving during the generation process. Early motion-based methods, e.g. X2Face [39], directly estimate a dense motion field by deep neural networks. Then the source face embedding is warped by the dense motion field to generate reenactment result. MonkeyNet [28] and its follow-up work, First Order Motion Model (FOMM) [29] uses abstract self-supervised 2D keypoints to estimate multiple local sparse motions and then combines them into a dense motion field to warp the source image. One-Shot Talking Heads [38] extends the 2D keypoints to the 3D space and estimates a dense motion field in 3D, which enables more expressive and flexible motion modeling. [30] replaces the keypoints with regions and models in-plane rotation and scaling via principal component analysis (PCA). However, it still lack the capability to model out-of-plane rotations and expressional motions. Among all the motion-based works, only FOMM [29] and [30] models the occlusion using an occlusion map learned in a self-supervised way. In this paper, two categories of occluded area are perceived separately and recovered with the help of the 3D geometry embedding and the widely used inpainting technique.

Landmark-based methods refer to those works that utilize facial landmarks as conditions to synthesis animated source with generative models. Few-shot Talking Heads [42] uses landmarks as the conditional input and injects source appearance features to the generator through an adaptive instance normalization layer (AdaIn) [12]. Fast Bi-layer [41] adopts spatially-adaptive denormalization (SPADE) layers [21] to build a neural rendering system based on the landmark skeleton. Few-shot Vid2Vid [36] uses SPADE residual blocks to inject landmark features into the generator, where appearance features are adapted to driving landmarks. Even landmarks are possible to represent face poses and relative motions, it still lacks the ability to preserve identity information and to handle occlusion.

3D-based methods take advantages of the geometric prior of 3D faces, which are reconstructed by fitting 3DMM parameters to the source and driving frames. The traditional 3D method Face2face [35] transfers driving expressions to the source through deformation transfer [32] and the reenacted 3D face is rendered to synthesize videos. In the deep learning era, 3D-based methods are combined with convolutional neural networks to generate more photo-realistic videos. HeadGAN [7] uses the rendered 3D faces of the driving frames as conditional inputs to build a neural rendering generator, where both deformed source appearance features and driving audio features are injected to generate realistic animation videos. StyleRig [34] and GIF [8] manipulate face images through a pre-trained StyleGAN conditioned on 3DMM parameters. Basically, these methods can generate realistic images, but are intrinsically disadvan-

tageous in spatial and temporal consistency on videos.

In this paper, we propose a method that combines the advantage of motion-based and 3D-based methods which can generate photo-realistic and identity-preserving reenactments even in case of large pose deviation and severe occlusion.

3. Preliminaries

3.1. 3D face prior

We use FLAME [16] as the 3D face prior model in our framework. As a statistical 3D morphable model, FLAME acts as a decoder, which takes shape (or identity) $\alpha \in \mathbb{R}^{|\alpha|}$, expression $\beta \in \mathbb{R}^{|\beta|}$ and pose $\theta \in \mathbb{R}^{3k+3}$ (k joint rotations and one global rotation) parameters as inputs and outputs a 3D head mesh with $v = 5023$ vertices and $f = 9976$ faces. The model is mathematically defined as follows,

$$M(\alpha, \beta, \theta) = W(T_P(\alpha, \beta, \theta), \mathbf{J}(\alpha), \mathbf{w}), \quad (1)$$

where $W(\cdot)$ is the blend skinning function, which takes the offseted template mesh $T_P \in \mathbb{R}^{v \times 3}$, joints position $\mathbf{J}(\alpha) \in \mathbb{R}^{3k}$ and blendweights $\mathbf{w} \in \mathbb{R}^{k \times v}$ as inputs.

3.2. Differentiable rendering

We adopt the differentiable renderer in PyTorch3D [23] to render the reconstructed 3D face. The differentiable renderer in our case rasterizes 3D vertex attributes into the 2D image space. The rendering function \mathcal{R} can be mathematically formulated as

$$I = \mathcal{R}(M(\alpha, \beta, \theta), \mathbf{c}, A), \quad (2)$$

where $I \in \mathbb{R}^{H \times W \times a}$ is the rendered 2D image. $\mathbf{c} = (s, \mathbf{t})$ represents the scale and translation parameters of the weak perspective camera model that we adopted in this paper. $A \in \mathbb{R}^{v \times a}$ represents 3D vertex attributes to fill on the image and a is the attribute dimension. In the following of the paper, we use attributes of RGB texture, vertex normal and vertex offset (motion) between source and driving faces. We only keep the face region of the rendered image through a face mask in the UV space (eyes and mouth regions are not included). With the differentiable renderer, the whole framework can be trained end-to-end.

4. Method

Our method takes a source face image $S \in \mathbb{R}^{H \times W \times 3}$ and a driving face image $D \in \mathbb{R}^{H \times W \times 3}$ as inputs and outputs the reenacted image D_{rec} that mimics the face pose and expression of the driving face. As shown in Figure 1, our model contains four modules: (1) 3D motion module, (2) first order motion module, (3) dense motion module, (4) occlusion aware generator. In (1), we estimate the 3DMM

parameters $\Phi = (\alpha, \beta, \theta, \mathbf{c})$ of the source and the driving images and render a 3D face motion field $\mathcal{T}_{S \leftarrow D, 3D}$ from driving to source. In (2), we predict K abstract keypoints and Jacobians from source and driving images respectively and approximate K local motion fields $\{\mathcal{T}_{S \leftarrow D, k}\}_{k=1}^K$ from driving to source. (3) integrates all the estimated local motions into a dense motion field $\hat{\mathcal{T}}$ and estimates two occlusion maps, O_g and O_{ca} . (4) takes S , $\hat{\mathcal{T}}$, O_g , O_{ca} and the driving 3DMM parameters Φ_D as inputs and generates the reenacted image. In the following subsections, we will discuss the four modules in detail.

4.1. 3D motion module

The 3D motion module contains a 3DMM encoder, a FLAME decoder and a differentiable renderer. The 3DMM encoder is pre-trained to predict the 3DMM parameters Φ given face images.

Given Φ_S and Φ_D , we obtain the 3D face vertices in the source and driving images from the FLAME decoder and then compute RGB texture, normal and motion attributes of these vertices. We interpret the source image $S(\ast)$ and the camera projection function $\mathbf{Proj}(\ast, \mathbf{c})$ as functions supporting arrayed inputs and outputs for convenience. The texture attributes $A_{S, tex} \in \mathbb{R}^{v \times 3}$ of the source image are obtained as

$$A_{S, tex} = S(\mathbf{Proj}(M(\alpha_S, \beta_S, \theta_S), \mathbf{c}_S)), \quad (3)$$

where bi-linear interpolation is used to obtain colors from S . The normal attributes $A_{S, n} \in \mathbb{R}^{v \times 3}$ and $A_{D, n} \in \mathbb{R}^{v \times 3}$ of both images (Omitted in the Figure 1) are computed from the neighbor faces of each vertex. The vertex motion attributes $A_m \in \mathbb{R}^{v \times 2}$ are computed as

$$A_m = \mathbf{Proj}(M(\alpha_S, \beta_S, \theta_S), \mathbf{c}_S) - \mathbf{Proj}(M(\alpha_D, \beta_D, \theta_D), \mathbf{c}_D). \quad (4)$$

Then by using Equation 2, we rasterize the following maps:

$$I_r = \mathcal{R}(M(\alpha_D, \beta_D, \theta_D), \mathbf{c}_D, A_{S, tex}), \quad (5)$$

$$I_{n, S} = \mathcal{R}(M(\alpha_S, \beta_S, \theta_S), \mathbf{c}_S, A_{S, n}), \quad (6)$$

$$I_{n, D} = \mathcal{R}(M(\alpha_D, \beta_D, \theta_D), \mathbf{c}_D, A_{D, n}), \quad (7)$$

$$\mathcal{T}_{S \leftarrow D, 3D} = \mathcal{R}(M(\alpha_D, \beta_D, \theta_D), \mathbf{c}_D, A_m). \quad (8)$$

3D reenactment $I_r \in \mathbb{R}^{H \times W \times 3}$ refers to the rendered driving 3D face with texture transferred from source 3D face. $I_{n, S} \in \mathbb{R}^{H \times W \times 3}$ and $I_{n, D} \in \mathbb{R}^{H \times W \times 3}$ denote the normal maps of source and driving images, respectively. $\mathcal{T}_{S \leftarrow D, 3D} \in \mathbb{R}^{H \times W \times 2}$ denotes the sparse motion field of the face area from the driving image to the source image. Exemplar visualization of I_r and $\mathcal{T}_{S \leftarrow D, 3D}$ is shown in Figure 1.

4.2. First order motion module

We adopt the same sparse motion estimator as proposed in FOMM [29], where an Hourglass Network

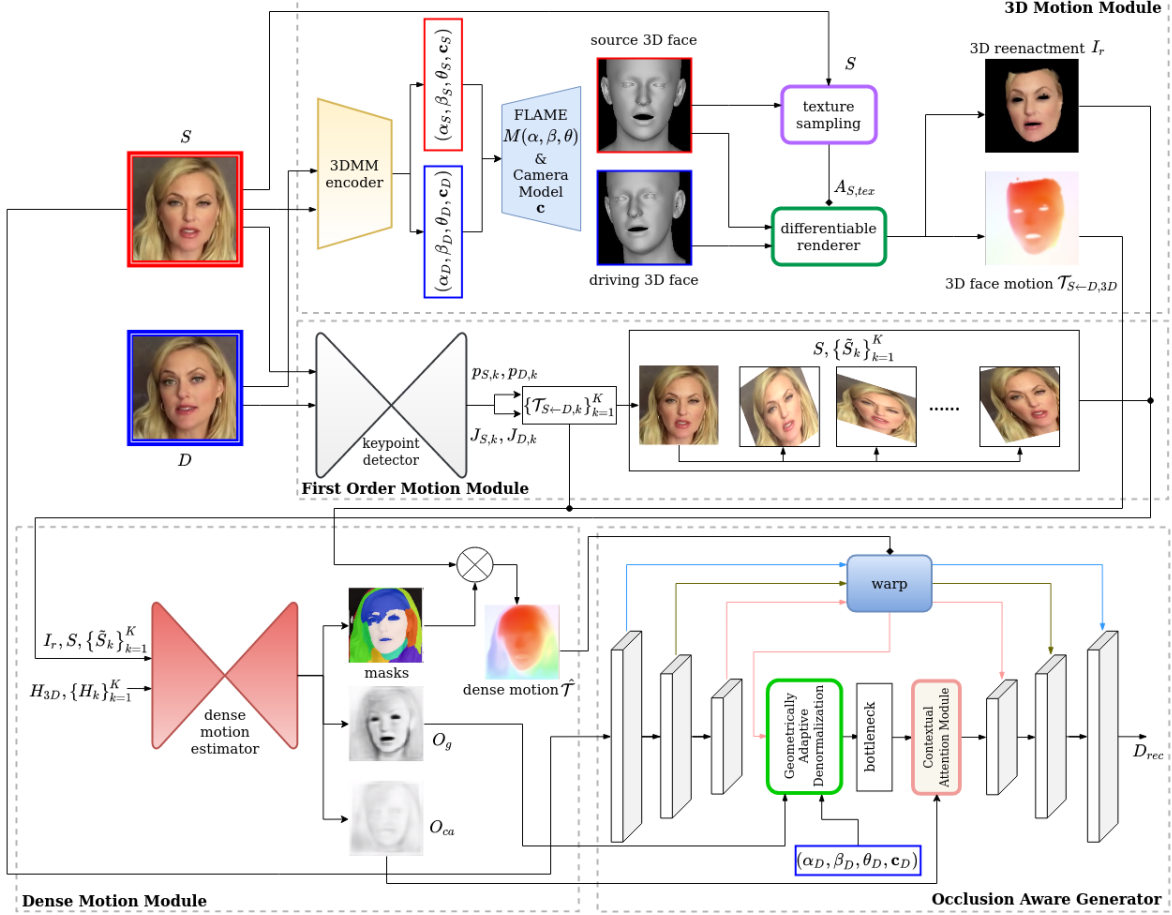


Figure 1. Overview of our framework. Our method takes a source image and a driving image as inputs. The 3D motion module reconstructs source and driving 3D faces and models 3D face motion using a differentiable renderer. The first order motion module dedicates to model the rest motions in terms of 2D local affine transformations. The dense motion module integrates both 3D and 2D motions and perceives where occlusion might happen. Finally, the U-Net shaped occlusion aware generator aligns the source image to the driving image.

[20] is used. The sparse motion estimator estimates $K = 10$ keypoints $\{p_{S,k}, p_{D,k} \in \mathbb{R}^2\}_{k=1}^K$ and Jacobians $\{J_{S,k}, J_{D,k} \in \mathbb{R}^{2 \times 2}\}_{k=1}^K$ from the source and driving images respectively. The keypoints are represented in terms of heatmaps and Jacobians are calculated by average pooling of the four output channels. The local affine transformations are approximated by the first order Taylor expansion,

$$\mathcal{T}_{S \leftarrow D,k}(z) \approx p_{S,k} + J_{S,k} J_{D,k}^{-1} (z - p_{D,k}), \quad (9)$$

where $z \in \mathbb{R}^2$ is an arbitrary point on the driving frame. K sparse affine motion fields $\{\mathcal{T}_{S \leftarrow D,k}(z)\}_{k=1}^K$ are used to warp the source image separately. Warping is simply implemented using a differentiable bi-linear sampling operator [13]. Besides the K sparse motions, an identity motion field is appended to model the static background area.

4.3. Dense motion module

The dense motion estimator is an encoder-decoder network with similar design to [29]. It takes deformed source

images and a set of heatmaps as inputs, which are downsampled to 1/4 of the original image resolution. The deformed source images include I_r , S , and $\{\tilde{S}_k\}_{k=1}^K$. \tilde{S}_k is obtained by deforming S using $\mathcal{T}_{S \leftarrow D,k}$. I_r corresponds to the face motion computed from the 3D face model. S corresponds to staticity of the background. \tilde{S}_k corresponds to the 2D affine motion. The heatmaps indicate where motion might happen and also have $K + 2$ instances. K heatmaps corresponding to the affine motion are formulated in the same way as FOMM, which are calculated by the difference of the heatmaps centered in $p_{D,k}$ and $p_{S,k}$.

$$H_k(z) = \exp\left(\frac{-(p_{D,k} - z)^2}{2\sigma}\right) - \exp\left(\frac{-(p_{S,k} - z)^2}{2\sigma}\right), \quad (10)$$

where $\sigma = 0.01$. For the static background, the heatmap is represented by a zero map. For the heatmap of 3D face motion, we make use of the z channel of the normal maps $I_{n,S}$ and $I_{n,D}$. The z component indicates whether the 3D face is visible in the camera frame, which helps estimate not

only the 3D motions but also the occlusion maps.

$$H_{3D}(z) = I_{n,D}^Z(z) - I_{n,S}^Z(z). \quad (11)$$

The dense motion module estimates $K + 2$ soft masks M_{3D} , M_b , and $\{M_k\}_{k=1}^K$ (as shown in the bottom-left block of Figure 1). M_{3D} corresponds to the rigid 3D face (in white). M_b corresponds to the static background (in black) and $\{M_k\}_{k=1}^K$ correspond to 2D affine motion areas (in other colors). All the masks are softmaxed to guarantee their pixel-wise sum to be 1. Sparse motion fields are fused into a dense motion field through masked averaging:

$$\hat{\mathcal{T}}(z) = M_b z + M_{3D} \mathcal{T}_{S \leftarrow D, 3D}(z) + \sum_{k=1}^K M_k \mathcal{T}_{S \leftarrow D, k}(z). \quad (12)$$

Besides, the dense motion module also predicts two occlusion maps $O_g, O_{ca} \in [0, 1]^{\frac{H}{4} \times \frac{W}{4}}$ indicating the invisible region in the source image but exposed in the generated image, which are used separately in the two elaborately designed modules in the occlusion aware generator to generate geometrically realistic face appearance and perceive areas need to be inpainted.

4.4. Occlusion aware generator

The occlusion aware generator takes S , $\hat{\mathcal{T}}$, O_g , O_{ca} and the driving 3DMM parameters Φ_D as inputs and generates the reenacted image D_{rec} . We adopt a U-Net [24] shaped network, which contains an encoder, a geometrically-adaptive denormalization layer (GADE), a bottleneck block, a contextual attention module and a decoder. The encoder takes the S as input and outputs source feature maps of different scales. Then those feature maps are warped by $\hat{\mathcal{T}}$. The highest level warped features are ported to the GADE layer, the bottleneck block and the contextual attention module to further recover the occluded area. The GADE layer and the contextual attention module are guided by O_g and O_{ca} separately to avoid potential conflicts. Finally, the recovered feature map is fed to the decoder. Lower level warped features are also concatenated to their corresponding intermediate features in the decoder in a U-Net fashion. Figure 1 shows this procedure in the bottom-right block.

4.4.1 Geometrically-adaptive denormalization

The geometrically-adaptive denormalization (GADE) layer takes normalized features from the encoder as input and denormalizes the features on the condition of the driving 3D face geometry. Figure 2 depicts the detail of GADE. We first normalize α, β, θ, c , separately for numerical stability. α, β, s are divided by their average norms calculated from the training set and θ, t are unchanged. Then we adopt two fully connected layers to convert these 3DMM parameters

to a scale vector γ and a shift vector δ . The warped source feature F_w is denormalized as

$$\hat{F}_w = \gamma \otimes F_w + \delta, \quad (13)$$

where \otimes denotes the broadcastable element-wise multiplication. The occlusion map O_g is used to adaptively select between F_w and \hat{F}_w ,

$$F_{out} = O_g \otimes F_w + (1 - O_g) \otimes \hat{F}_w. \quad (14)$$

In the original FOMM, F_{out} only contains the first term and F_w is simply diminished where occlusion occurs, i.e. O_g is small. In our approach, we combine \hat{F}_w in this case to exploit information from the driving image to recover the occlusion area. The injected face geometry code depicts 3D face shape, pose and expression of the driving image, which makes it easier to generate realistic features and to follow the driving face pose accurately. Note that we use Φ_D instead of its precedent features from the 3DMM encoder. This is because the precedent features might contain appearance information of the driving image and might cause the reenacted image to collapse to reconstruct the driving image without using the source image features. Denormalized features from GADE are fed to a bottleneck block that contains several residual convolution blocks, then used as inputs to the contextual attention module.

4.4.2 Contextual attention module

The original Contextual Attention Module (CAM) proposed in [40] was used for image inpainting, which takes a feature map and a hard binary mask indicating the regions to be inpainted as input. In our case, we adopt this module to realistically recover the perceived occluded areas indicated by the soft occlusion map O_{ca} from the dense motion module. Note that we do not use face inpainting methods such as [17] since we need to recover both the foreground and the background areas. The contextual attention module contains two parallel branches: (1) contextual attention layer, (2) dilated convolution block. The contextual attention layer is a differentiable layer. Patches are extracted from output features from the bottleneck block and each patch is used as a kernel to compute a matching score with the occluded area through convolution. Then softmax function is applied to the convoluted features to get attention scores. Finally occluded areas are reconstructed through deconvolution on the scores with the related feature patches. The parallel branch, dilated convolution block contains pyramid dilated convolution layers of dilation ratios 1, 2, 4, 8. It not only increases the reception field but also enables the network to hallucinate according to neighbor regions. The output from contextual attention layer and the dilated convolution block is then concatenated and fed to the decoder of the generator to obtain the final reenactment D_{rec} . Figure 3 depicts the detail of the contextual attention module.

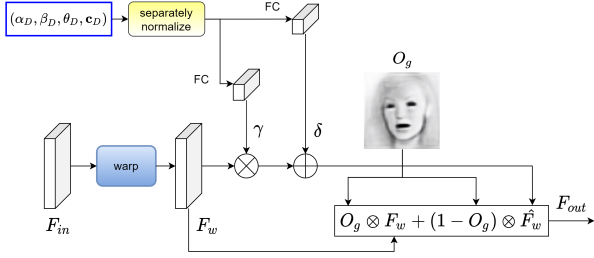


Figure 2. Illustration of the Geometrically-Adaptive Denormalization Layer

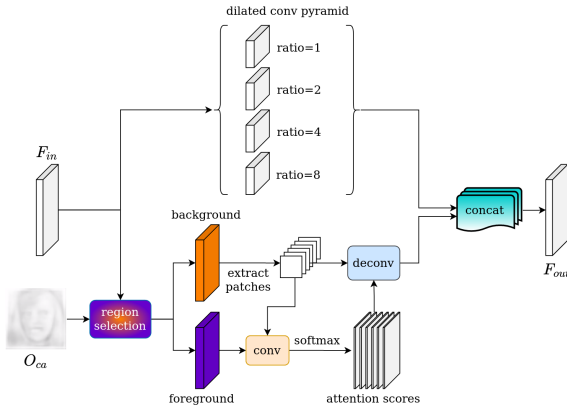


Figure 3. Illustration of the Contextual Attention Module

4.5. Training

The 3DMM encoder is pre-trained before we train the whole model end-to-end. We project 68 3D facial landmarks on the surface of the reconstructed 3D face $l_i \in \mathbb{R}^3$ into 2D image space and use the 2D landmark prediction $k_i \in \mathbb{R}^3$ from FAN [3] to formulate the $L1$ landmark re-projection loss. Besides, we use $L2$ parameter loss to regulate the 3DMM shape parameters α and expression parameters β to obtain reasonable reconstruction results.

Then, our framework is trained in a self-supervised strategy. We perform self-reenactment, where a pair of source and driving images are sampled from each video. The driving images are used as the ground truth for supervision. We adopt loss functions of two main objectives: (1) to improve the faithfulness and quality of output images; (2) to regulate the 3D motion module and the first order motion module to predict reasonable motion fields. For the first objective, we follow the same loss functions in FOMM, including perceptual loss [14] and GAN-related loss. To regulate 3D motion, we preserve the 3DMM loss function in the pre-train stage. To stabilize 2D sparse keypoints and jacobians, we adopt the equivariance constraints from FOMM. Please refer to the supplementary material for the detailed expressions of loss functions.

5. Experiments

5.1. Dataset

We use Voxceleb1 [19] dataset to train and evaluate our approach. The Voxceleb1 dataset contains 22496 videos of 1251 celebrities extracted from YouTube videos. We adopt the same pre-processing strategy and train-test split as that of the FOMM. Videos of size 256×256 are cropped from the original videos covering the heads. In total, we obtain 17941 videos for training and 506 videos for testing after cropping. Besides, we also extract images from the pre-processed videos and use FAN [3] to predict 68 2D facial landmarks of each image.

5.2. Implementation details

We first pre-train the 3DMM encoder using the images extracted from pre-processed videos for every 10 frame. We adopt the lightweight MobilenetV2 [26] as our 3DMM encoder and modify the output fully connected layer to output 156 dimensions of FLAME parameters and 3 dimensions of camera parameters. We train the 3DMM estimator for 50 epochs with batch size of 128. Adam optimizer is used with learning rate $2e - 4$ and $\beta_1 = 0.9, \beta_2 = 0.999$

For the end-to-end model training, we perform self-reenactment on the pre-processed videos. We randomly sample a pair of source and driving images from each video and use the driving image as ground truth. The 3DMM encoder is jointly fine-tuned with the animation model during end-to-end training. We train the whole model for 50 epochs and for each epoch we repeat the video list for 150 times. We adopt synchronized Batch Normalization to further improve the performance. Adam optimizer is used with learning rate $2e - 4$ and $\beta_1 = 0.5, \beta_2 = 0.9$ for all network modules. We use 8 24GB NVIDIA P40 GPUs to train our model.

5.3. Evaluation Metrics

To quantitatively evaluate our approach and the previous state-of-the-art methods, we perform video reconstruction experiments, where the first frame of each test video is used as the source image and the rest frames are used as driving images. The following metrics are computed:

(1) **L1 distance**: Averaged L1 distance between the reconstructed images and the ground truth images.

(2) **Average Keypoint Distance (AKD)**: We extract facial landmarks from the reconstructed and ground truth images using FAN [3] and compute the average distance between the two set of landmarks. AKD measures how accurately the generated face video follows the face poses of the driving video.

(3) **Average Euclidean Identity Distance (AEID)**: We adopt the widely used face recognition network, ArcFace [6] to extract identity features from both generated



Figure 4. Qualitative comparison with SOTA methods on motion transfer.

and ground truth images. Then we calculate the averaged euclidean distance between the two. AEID measures how much identity information is preserved.

(4) **Fréchet Inception Distance (FID)**: FID was proposed in [11], which measures how realistic the visual quality of the generated images is, compared with the real images. We adopt the implementation of [27], which uses a pre-trained InceptionV3 [33] network to extract features from generated and ground truth images and calculates distribution distance between the two sets of features.

(5) **Peak Signal-to-Noise Ratio (PSNR)**: PSNR measures the image quality between generated and real images in terms of the absolute mean squared error (MSE).

(6) **Structural Similarity Index (SSIM)**: SSIM measures the structural similarity between the reconstructed and real images, which is more robust to global illumination

Table 1. Quantitative comparison with SOTA methods on Vox-celeb1 test set.

	L1↓	AKD↓	AEID↓	FID↓	PSNR↑	SSIM↑
FS-Vid2vid [36]	0.0869	6.214	0.3538	18.490	29.349	0.559
FS T. Heads [42]	0.1053	10.247	0.3243	27.629	28.906	0.507
Fast Bi-layer [41]	0.3775	12.486	0.3544	109.187	28.118	0.319
FOMM[29]	0.0448	1.386	0.1654	14.219	30.567	0.764
Ours	0.0401	1.222	0.1532	8.546	31.079	0.783

changes compared with PSNR.

The units of L1, AKD, and PSNR are RGB intensity, pixels, and dB, respectively. FID, AEID and SSIM have no units since they are defined perceptually.

5.4. Comparison with state of the art

We compare our approach with the state-of-the-art face animation methods FS-Vid2vid [36], FS Talking

Table 2. Quantitative results of ablation study on Voxceleb1 test set.

	L1↓	AKD↓	AEID↓	FID↓	PSNR↑	SSIM↑
FOMM baseline	0.0448	1.386	0.1654	14.219	30.567	0.764
FOMM+3D	0.0437	1.315	0.1676	13.894	30.644	0.776
FOMM+3D+CAM	0.0418	1.306	0.1642	11.273	30.965	0.780
Full model	0.0401	1.222	0.1532	8.546	31.079	0.783

Heads [42], Fast Bi-layer [41], FOMM [29], and One-Shot Talking Heads [38]. For FS-Vid2vid and Fast Bi-layer, we use the pre-trained models provided by the authors. But note that their training sets are much larger than ours (FS-Vid2vid: FaceForensics [25], Fast Bi-layer: Voxceleb2 [5]). For FS Talking Heads and FOMM, we train their models on the same training set as our method. For the latest proposed method One-Shot Talking Heads [38], due to the unavailability of the official source code, we only provide qualitative comparisons using cases reported in their paper.

We provide quantitative evaluation results for video reconstruction on the test set of Voxceleb1 dataset, as listed in Table 1. Our method performs best in all metrics. To visually illustrate the improvement besides the above abstract metrics, we provide qualitative visualizations on motion transfer, where source and driving faces are of different identities, as shown in Figure 4. We intentionally choose image pairs with large pose differences, translational motions and challenging expressions to show our superiority. The qualitative comparison with One-Shot Talking Heads is shown in Figure 5. We generate face images using a similar relative motion transfer strategy as in FOMM, please refer to the supplementary material for the details of the inference process of the motion transfer. Our method shows superiority in generating realistic, pose-preserving and identity-preserving face animation results. Moreover, we improve the image quality for the occluded failure cases by [38] with fewer parameters and training data.

5.5. Ablation study

We conduct both quantitative and qualitative experiments on video reconstruction to show the contributions of the 3D motion module, the Contextual Attention Module (CAM) and the Geometrically-Adaptive Denormalization layer (GADE). Table 2 shows the quantitative results on the test set of Voxceleb1. With the integration of 3D face prior knowledge, the AKD score decrease significantly, which means the face pose of the reenactment becomes more accurate. The additional contextual attention module further improves the FID scores with context-coherent inpainting. Finally, injection 3D face geometry information into the generation process benefits both image quality and face appearance of the generated image. This shows the effectiveness of 3D geometry embedding in learning the scene structures. We also provide qualitative results as shown in Figure 6. Without 3D modelling, the faces might be improperly



Figure 5. Qualitative comparison with [38] on two successful (row 1-2) and two failed (row 3-4) cases.



Figure 6. Qualitative results of ablation study on video reconstruction.

distorted. With the contextual attention module, the visual quality of the generated images improves. With the GADE layer, facial details and expressions become more realistic and self-occluded face area is properly recovered. The full model outperforms all the model variations in terms of visual sharpness of both foreground and background.

6. Conclusions

In this paper, we propose SAFA, a structure aware face animation method. It enhances the animation quality by exploiting the scene structure knowledge to model each component in the image. We show the effectiveness of combining 3DMM with the 2D affine motion model in modelling the foreground. Besides, the contextual attention module for image inpainting further helps generate realistic images by recovering the perceived occluded area. Moreover, the proposed geometrically-adaptive denormalization layer boosts the quality of face appearance by injection of 3D geometry information. Both quantitative and qualitative experiments show that our method outperforms the state-of-the-art methods.

References

- [1] Volker Blanz and Thomas 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
- [2] James Booth, Anastasios Roussos, Stefanos Zafeiriou, Allan Ponniah, and David Dunaway. A 3d morphable model learnt from 10,000 faces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5543–5552, 2016. 2
- [3] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In *International Conference on Computer Vision*, 2017. 6
- [4] Chen Cao, Yanlin Weng, Shun Zhou, Yiyong Tong, and Kun Zhou. Facewarehouse: A 3d facial expression database for visual computing. *IEEE Transactions on Visualization and Computer Graphics*, 20(3):413–425, 2013. 2
- [5] J. S. Chung, A. Nagrani, and A. Zisserman. Voxceleb2: Deep speaker recognition. In *INTERSPEECH*, 2018. 8
- [6] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2019. 6
- [7] Michail Christos Doukas, Stefanos Zafeiriou, and Viktoriia Sharmanska. Headgan: Video-and-audio-driven talking head synthesis. *arXiv preprint arXiv:2012.08261*, 2020. 1, 2
- [8] Partha Ghosh, Pravir Singh Gupta, Roy Uziel, Anurag Ranjan, Michael J. Black, and Timo Bolkart. GIF: Generative interpretable faces. In *International Conference on 3D Vision (3DV)*, pages 868–878, 2020. 2
- [9] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*, 2014. 1
- [10] Sungjoo Ha, Martin Kersner, Beomsu Kim, Seokjun Seo, and Dongyoung Kim. Marionette: Few-shot face reenactment preserving identity of unseen targets. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020. 1
- [11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Günter Klambauer, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a nash equilibrium. *CoRR*, abs/1706.08500, 2017. 7
- [12] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1501–1510, 2017. 2
- [13] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer networks. *arXiv preprint arXiv:1506.02025*, 2015. 4
- [14] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 6
- [15] Hyeonwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Niessner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. Deep video portraits. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018. 1
- [16] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4d scans. *ACM Trans. Graph.*, 36(6):194–1, 2017. 1, 2, 3
- [17] Xiaoming Li, Guosheng Hu, Jieru Zhu, Wangmeng Zuo, Meng Wang, and Lei Zhang. Learning symmetry consistent deep cnns for face completion. *IEEE Transactions on Image Processing*, 29:7641–7655, 2020. 5
- [18] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2794–2802, 2017. 11
- [19] A. Nagrani, J. S. Chung, and A. Zisserman. Voxceleb: a large-scale speaker identification dataset. In *INTER-SPEECH*, 2017. 6
- [20] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *European conference on computer vision*, pages 483–499. Springer, 2016. 4
- [21] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2337–2346, 2019. 2
- [22] Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas 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
- [23] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501*, 2020. 3
- [24] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 5
- [25] Andreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. FaceForensics++: Learning to detect manipulated facial images. In *International Conference on Computer Vision (ICCV)*, 2019. 8
- [26] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018. 6
- [27] Maximilian Seitzer. pytorch-fid: FID Score for PyTorch. <https://github.com/mseitzer/pytorch-fid>, August 2020. Version 0.1.1. 7
- [28] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. Animating arbitrary objects via

- deep motion transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2377–2386, 2019. [2](#)
- [29] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. *arXiv preprint arXiv:2003.00196*, 2020. [1](#), [2](#), [3](#), [4](#), [7](#), [8](#)
- [30] Aliaksandr Siarohin, Oliver Woodford, Jian Ren, Menglei Chai, and Sergey Tulyakov. Motion representations for articulated animation. In *CVPR*, 2021. [2](#)
- [31] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [11](#)
- [32] Robert W Sumner and Jovan Popović. Deformation transfer for triangle meshes. *ACM Transactions on graphics (TOG)*, 23(3):399–405, 2004. [2](#)
- [33] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. [7](#)
- [34] Ayush Tewari, Mohamed Elgharib, Gaurav Bharaj, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhofer, and Christian Theobalt. Stylerig: Rigging stylegan for 3d control over portrait images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6142–6151, 2020. [2](#)
- [35] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2387–2395, 2016. [2](#)
- [36] Ting-Chun Wang, Ming-Yu Liu, Andrew Tao, Guilin Liu, Jan Kautz, and Bryan Catanzaro. Few-shot video-to-video synthesis. *arXiv preprint arXiv:1910.12713*, 2019. [1](#), [2](#), [7](#)
- [37] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018. [11](#)
- [38] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. *arXiv preprint arXiv:2011.15126*, 2020. [1](#), [2](#), [8](#)
- [39] Olivia Wiles, A Koepke, and Andrew Zisserman. X2face: A network for controlling face generation using images, audio, and pose codes. In *Proceedings of the European conference on computer vision (ECCV)*, pages 670–686, 2018. [1](#), [2](#)
- [40] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5505–5514, 2018. [2](#), [5](#)
- [41] Egor Zakharov, Aleksei Ivakhnenko, Aliaksandra Shysheya, and Victor Lempitsky. Fast bi-layer neural synthesis of one-shot realistic head avatars. In *European Conference on Computer Vision*, pages 524–540. Springer, 2020. [1](#), [2](#), [7](#), [8](#)
- [42] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. Few-shot adversarial learning of realistic neural talking head models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9459–9468, 2019. [1](#), [2](#), [7](#), [8](#)

A. Details of loss function

Perceptual loss We adopt a pre-trained VGG-19 network [31] to extract features of l layers from the driving image D and the reconstructed image D_{rec} . Original images are downsampled 4 times to extract a feature pyramid. Then $L1$ loss is applied to the two sets of feature maps.

$$\mathcal{L}_p = \sum_{i=1}^4 \sum_{j=1}^l \left\| f_{VGG}^{i,j}(D) - f_{VGG}^{i,j}(D_{rec}) \right\|_1. \quad (15)$$

GAN loss We adopt a multi-scale discriminator \mathbf{D}_m as used in [37]. The driving image and the reconstructed image of multiple resolutions are fed into the discriminator. We use the least square loss proposed in [18] to enable high-quality image generation and stable training.

$$\mathcal{L}_G = \sum_{i=1}^s \mathbb{E} \left[(1 - \mathbf{D}_m^i(D_{rec}))^2 \right], \quad (16)$$

$$\mathcal{L}_D = \sum_{i=1}^s \mathbb{E} \left[(\mathbf{D}_m^i(D_{rec}))^2 \right] + \mathbb{E} \left[(1 - \mathbf{D}_m^i(D))^2 \right]. \quad (17)$$

3DMM loss The 3DMM loss contains a $L1$ landmark re-projection term and a $L2$ parameter regularization term.

$$L_{3DMM} = \sum_{i=1}^{68} \|k_i - \mathbf{c}(l_i)\|_1 + \lambda_\alpha \|\alpha\|_2^2 + \lambda_\beta \|\beta\|_2^2, \quad (18)$$

where $\lambda_\alpha, \lambda_\beta$ are weights of the regularization on shape and expression parameters. The loss is used in both pre-training and the end-to-end stage.

Equivariance constraints We adopt this loss to ensure consistency of the sparse keypoints. Specifically, we randomly sample an affine transformation $\mathcal{T}(z)$ and apply it to the original driving image. If we transform the predicted keypoints from the deformed driving image back, the result should be sufficiently close to the keypoints estimated from the original driving image. We minimize the $L1$ loss between the two value, similar constrains can also be applied to the Jacobians.

$$\mathcal{L}_e = \sum_{k=1}^K \|p_{D,k} - T^{-1}(p_{D_T,k})\|_1. \quad (19)$$

In total, the final loss function is given by

$$\mathcal{L} = \lambda_p \mathcal{L}_p + \lambda_G \mathcal{L}_G + \lambda_D \mathcal{L}_D + \lambda_{3DMM} \mathcal{L}_{3DMM} + \lambda_e \mathcal{L}_e, \quad (20)$$

where λ_* are weights that balance contributions of each loss term. The loss weights are set as $\lambda_p = 10, \lambda_G = 1, \lambda_D = 1, \lambda_{3DMM} = 1, \lambda_\alpha = 1e - 2, \lambda_\beta = 8e - 3, \lambda_e = 10$.

B. Inference

Our inference procedure is similar to that of FOMM. We first select a reference frame D_r from the driving video,

which has the nearest pose with respect to S . Then we impose the relative motion between the D_r and each frame D_t of the driving video on S to generate the reenactment S_t . For the 2D affine motion, the relative motion $\mathcal{T}_{S \leftarrow S_t, k}$ is represented as ¹

$$\mathcal{T}_{S \leftarrow S_t, k}(z) \approx p_{S,k} + J_{D_r, k} J_{D_t, k}^{-1} (z - p_{S,k} + p_{D_r, k} - p_{D_t, k}). \quad (21)$$

For the 3D face motion, we calculate the relative motion in terms of relative 3DMM parameters. The 3DMM parameters of the relatively reenacted image are

$$\Phi_{S_t} = (\alpha_S, \beta_S + \beta_{D_t} - \beta_{D_r}, \theta_S + \theta_{D_t} - \theta_{D_r}, (s_S \frac{s_{D_t}}{s_{D_r}}, \mathbf{t}_S + \mathbf{t}_{D_t} - \mathbf{t}_{D_r})). \quad (22)$$

Then the relative vertex motion attributes are

$$A_{m, S_t} = \mathbf{Proj}(M(\alpha_S, \beta_S, \theta_S), \mathbf{c}_S) - \mathbf{Proj}(M(\alpha_{S_t}, \beta_{S_t}, \theta_{S_t}), \mathbf{c}_{S_t}). \quad (23)$$

Finally, the relative 3D face motion is rendered as

$$\mathcal{T}_{S \leftarrow S_t, 3D}(z) = \mathcal{R}(M(\alpha_{S_t}, \beta_{S_t}, \theta_{S_t}), \mathbf{c}_{S_t}, A_{m, S_t}). \quad (24)$$

C. Intermediate components visualization

In Figure 7, we provide additional qualitative visualizations of the intermediate components produced by each module in our framework, including the 2D affine keypoints $\{p_{S,k}, p_{D,k}\}_{k=1}^K$, 3D reenactment I_r , 3D face motion $\mathcal{T}_{S \leftarrow D, 3D}$, motion masks M_{3D}, M_b , and $\{M_k\}_{k=1}^K$, dense motion $\hat{\mathcal{T}}$, and two occlusion maps O_g and O_{ca} for the GADE layer and the contextual attention module respectively. The excellent accuracy of the 3DMM encoder can be demonstrated from the 3D reenactment I_r (column 6).

¹ $J_{S,k}$ is canceled out here.

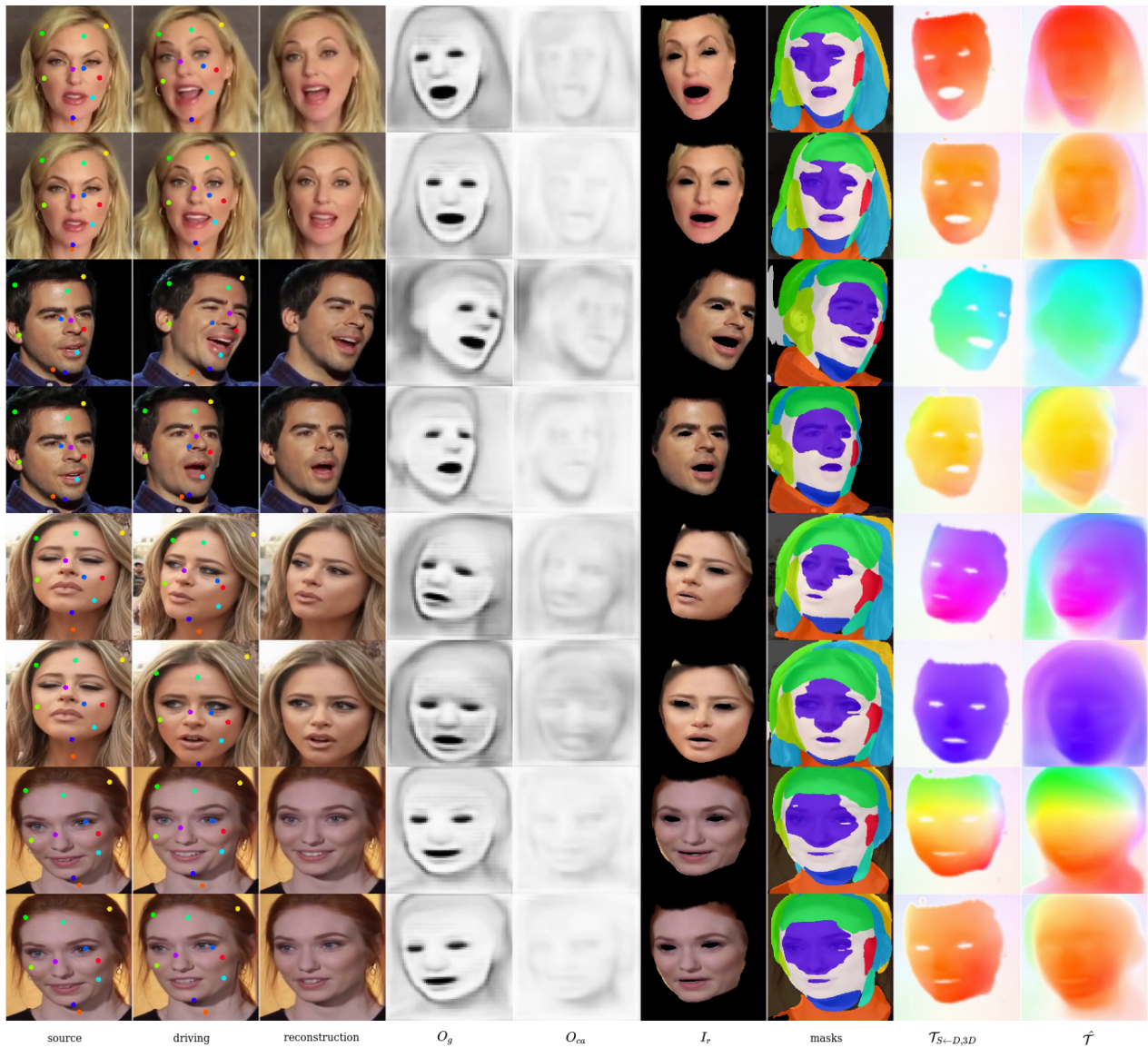


Figure 7. Visualization of intermediate components. From left to right are: source image with 2D keypoints, driving image with 2D keypoints, reconstructed image, occlusion map for GADE O_g , occlusion map for contextual attention module O_{ca} , 3D reenactment I_r , dense motion mask (the white part corresponds to M_{3D} , the black part corresponds to M_b and others correspond to $\sum_{k=1}^K M_k$), 3D motion field $\mathcal{T}_{S \leftarrow D, 3D}$ and dense motion field $\hat{\mathcal{T}}$