# HeadNeRF: A Real-time NeRF-based Parametric Head Model

Yang Hong    Bo Peng    Haiyao Xiao    Ligang Liu    Juyong Zhang*

University of Science and Technology of China

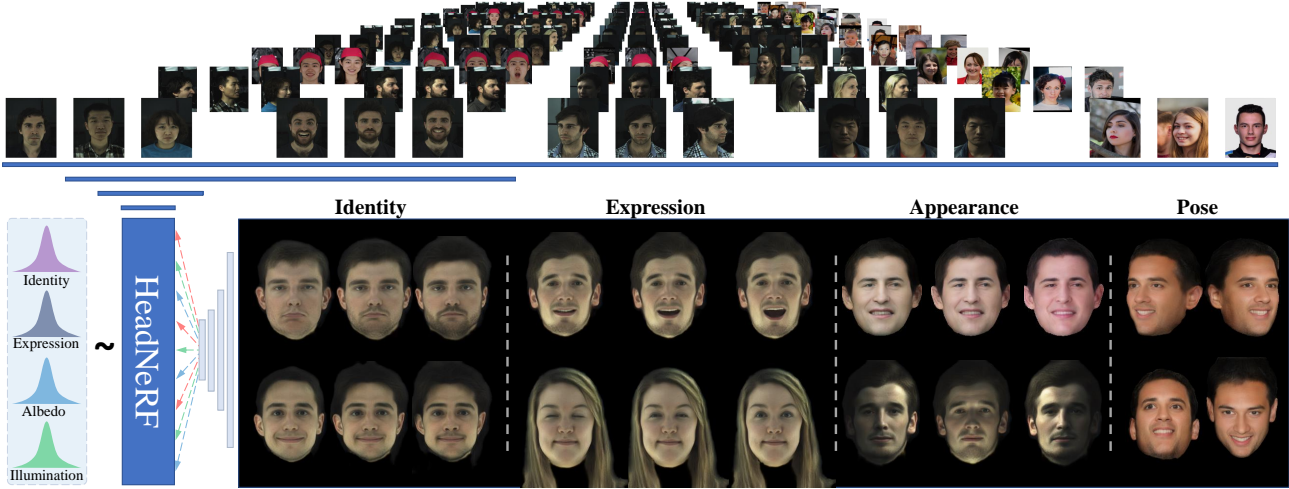{hymath@mail., pb15881461858@mail., xhy1999512@mail., lgliu@, juyong@}ustc.edu.cn

Figure 1. HeadNeRF, a NeRF-based parametric head model, is able to render high fidelity head images in real-time, and supports directly controlling the generated images' rendering pose, and various semantic attributes. The images in the black area are generated by exploring the latent space of HeadNeRF.

## Abstract

*In this paper, we propose HeadNeRF, a novel NeRF-based parametric head model that integrates the neural radiance field to the parametric representation of the human head. It can render high fidelity head images in real-time, and supports directly controlling the generated images' rendering pose and various semantic attributes. Different from existing related parametric models, we use the neural radiance fields as a novel 3D proxy instead of the traditional 3D textured mesh, which makes that HeadNeRF is able to generate high fidelity images. However, the computationally expensive rendering process of the original NeRF hinders the construction of the parametric NeRF model. To address this issue, we adopt the strategy of integrating 2D neural rendering to the rendering process of NeRF and design novel loss terms. As a result, the rendering speed of HeadNeRF can be significantly accelerated, and the rendering time of one frame is reduced from 5s to 25ms. The novel-designed loss terms also improve the rendering accuracy, and the fine-*
*level details of the human head, such as the gaps between teeth, wrinkles, and beards, can be represented and synthesized by HeadNeRF. Extensive experimental results and several applications demonstrate its effectiveness. We will release the code and trained model to the public.*

## 1. Introduction

The parametric face/head model, which encodes the human face/head in low-dimensional space, is a hot research topic in computer vision and computer graphics and widely used in many applications like identity recognition [34, 52], face analysis [11, 64] and film/game production [47], etc. Early works of the parametric face/head model [3, 6, 11, 30, 41, 57] mainly model 3D faces with the topologically uniformed face template mesh and usually ignore to represent the non-face parts, such as hair and teeth. With the development of deep learning, 2D generative adversarial networks (GAN) [21, 22] are able to directly render photo-realistic face images without the help of 3D modeling. Some methods [5, 9, 14, 26, 36] further introduce seman-

---

*Corresponding Author

tically disentangled constraints to render the face images in a user-controlled way. However, their rendered results from different views often tend to be inconsistent as they do not explicitly encode or model 3D geometry.

Recently, Mildenhall *et al*. [35] propose to represent 3D scenes using neural radiance fields (NeRF). This strategy can synthesis photorealistic images and has emerged as a compelling technique. Compared with the above-mentioned generative methods, the density field from NeRF actually implicitly encodes the 3D geometry shape of the scene. Therefore, the results from NeRF have excellent multi-view consistency. In another opinion, NeRF itself can be regarded as a novel 3D representation equivalent to a textured mesh and naturally supports differentiable rendering.

Based on this observation, we apply the NeRF structure to the representation of human heads and propose a novel NeRF-based parametric head model, HeadNeRF. HeadNeRF inherits the excellent properties of NeRF, which can generate high fidelity head images and maintain remarkable multi-view consistency. Moreover, NeRF itself supports freely changing the camera perspective used for rendering, and thus HeadNeRF naturally supports the pose editing of rendered images. It is challenging for the above-mentioned 2D generative methods. In addition, HeadNeRF intrinsically supports differentiable rendering. In contrast, traditional 3D representations, such as mesh, point cloud, voxel, etc., need to design various approximation methods [19, 23, 27, 28, 32] to alleviate the non-differentiable problems of their rendering process. Therefore, compared with the previous methods, which need to capture and process a large amount of high-quality 3D scan data, the construction of HeadNeRF only needs 2D images as input. Specifically, we collect and process three large-scale face datasets and design novel loss terms to disentangle this parametric representation. With the well designed network structure and loss function, HeadNeRF can semantically disentangle the identity, expression, and appearance of the rendered images. Fig. 1 shows some results by freely exploring within the space of HeadNeRF.

We further integrate the volume rendering of NeRF with 2D neural rendering to achieve real-time rendering. Similar with GIRAFFE [38] and StyleNeRF [16], this coarse to fine strategy significantly accelerates the rendering speed of HeadNeRF, and it can exceed 40fps without sacrificing the rendering quality. Benefiting from its nice disentangled representation, real-time rendering of inference, and high fidelity generated results, we apply HeadNeRF to various applications, including novel view synthesis from a single face image, semantically editing face attributes, and even facial reenactment where the expressions of one person are transferred to another person. In summary, the main contributions of this paper include:

- We propose the first NeRF-based parametric human head model, which can directly control the rendering pose, identity, expression, and appearance in real-time.

- We propose an effective training strategy to train the model from general 2D image datasets, and the trained model can generate high fidelity rendered images.

- We design and implement several novel applications with HeadNeRF, and the results verify its effectiveness. We believe that more interesting applications can be explored with our HeadNeRF.

## 2. Related Work

**Parametric Face/Head Model.** As different people share similar face shapes and appearances, human face can be embedded into a low-dimensional parametric space. Based on this observation, Blanz and Vetter [3] propose to build 3D morphable model (3DMM), which has been further improved and widely used for 3D face representation [6, 11, 41, 51]. 3DMM decomposes the intrinsic attributes of the human face into identity, expression, and reflectance. They are encoded as low-dimensional vectors and can be used to restore 3D textured face mesh using corresponding blendshapes. The scene illumination is often simultaneously modeled via the spherical harmonics function [15]. Recently, some generative adversarial [21, 22] methods can directly generate photorealistic face images without the aid of 3D modeling. However, these methods pay more attention to generating images that meet the specified distribution and lack semantic and interpretable control over the image synthesis. Some strategies, such as adding disentangled loss [9, 26, 49], embedding the above-mentioned parametric face model into generative adversarial network (GAN) [5, 14], etc., are proposed to alleviate this problem. However, these 2D GAN methods essentially still lack an explicit 3D geometric structure. Therefore, the rendering results tend to be inconsistent as the camera pose changes.

**Neural Radiance Field.** NeRF [35] represents 3D scenes using an implicit MLP-based function and shows convincing rendering quality for the novel view synthesis task. Meanwhile, the rendering process of NeRF is inherently differentiable. Therefore, the training of NeRF can be completed only using multi-view images with camera parameters. Benefiting from these advantages, NeRF has been widely used in many fields, such as 3D modeling [53, 58], human face/body digitization [12, 17, 42, 43, 48, 55], generating 4D free-view video [31, 39], etc. Besides, many works are proposed to further improve NeRF, including speeding up training [1, 2] and inference [13, 18, 44, 56, 59], improving rendering quality [55], and reducing the number of required inputs [8, 60]. Please refer to [50] for a complete summary.
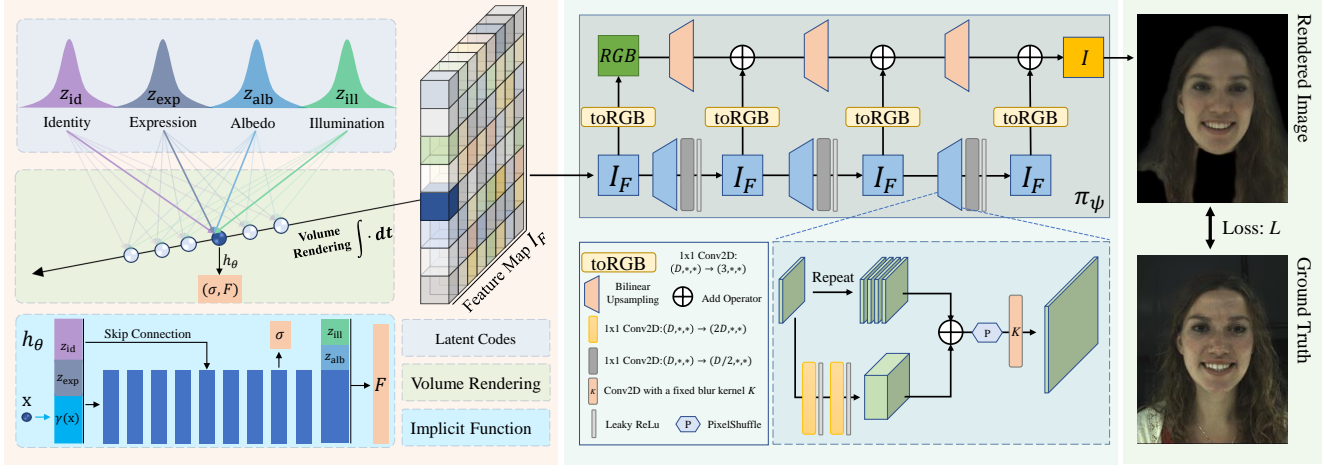
Figure 2. Overview of HeadNeRF. Given semantic latent codes and camera parameters, the MLP-based implicit function $h_\theta$ is utilized to predict the density $\sigma(\mathbf{x})$ and feature vector $F(\mathbf{x})$ of the 3D point $\mathbf{x}$ sampled from one ray. Then we perform volume rendering to generate a low-resolution feature map $I_F$, which is further used to render the final result $I$ by our well-designed 2D neural rendering module $\pi_\psi$. The whole process is differentiable, and thus the construction of HeadNeRF can be completed using only 2D images.

**NeRF-Based GANs.** Some works [7, 16, 38, 45, 63] integrate NeRF with GANs to design 3D-aware generators. Thanks to the introduction of NeRF structure, these methods generally can directly control the pose of synthesized results and effectively improve the multi-view consistency of generated images, which is challenging for 2D generative models [21, 22]. Voxel-based GANs [36, 37] can alleviate this problem, but their generated results tend to lack fine details due to the voxel resolution restriction. GRAF [45] is the first to introduce the neural radiance field into GAN. Although the quality and 3D consistency of the rendering results can be improved by using this method, it still struggles to render high-resolution and high-fidelity images due to the expensive rendering process of the neural radiance field. GIRAFFE [38] further improved the training and rendering efficiency by combining NeRF with a 2D CNN-based neural renderer and can significantly improve the computational speed of NeRF with a slight loss of accuracy. In addition, some works [7, 16, 63] attempt to design novel ways of applying conditions to generate more fine details.

## 3. Method

In this work, we propose HeadNeRF, a novel parametric model integrating neural radiance field to human head representation. Unlike the previous 3D mesh-based topologically uniformed face parametric model [3,4,6,30,41,51,57], HeadNeRF takes the accelerated variant of neural radiance field as a unified 3D proxy, which allows it to directly control the viewing pose of the rendered result and generate high fidelity head images in real-time. To train HeadNeRF, we collected and processed three large-scale face image datasets. Meanwhile, the novel network structure and loss

terms are well designed such that the trained parametric model can semantically control and edit the rendering result's identity, expression, and appearance. In the following, we will describe the algorithm details.

### 3.1. Model Representation

HeadNeRF is a NeRF-based parametric model $\mathcal{R}$, which can render an image $I$ with specified attributes for a given camera parameter and some semantic codes. It is formulated as:

$$I = \mathcal{R}(\mathbf{z}_{\text{id}}, \mathbf{z}_{\text{exp}}, \mathbf{z}_{\text{alb}}, \mathbf{z}_{\text{ill}}, P), \quad (1)$$

where $P$ is the camera parameter used for rendering, including the extrinsic matrix and the intrinsic matrix. $\mathbf{z}_*$ represent the latent codes for four independent factors: identity $\mathbf{z}_{\text{id}}$, expression $\mathbf{z}_{\text{exp}}$, the albedo $\mathbf{z}_{\text{alb}}$ of the face, and the illumination $\mathbf{z}_{\text{ill}}$ of the scene. To obtain a disentangled and interpretable control to the rendered results, we integrate the neural radiance field to the parametric head model and describe it as follows.

### 3.2. Network Architecture

Similar to 3DMM, we consider that the underlying geometric shape of the head image is mainly controlled by latent codes related to identity and expression, and the latent codes of albedo and illumination are responsible for the appearance of rendered heads. Therefore, the MLP-based implicit function $h_\theta$ of NeRF is adjusted as:

$$h_\theta : (\gamma(\mathbf{x}), \mathbf{z}_{\text{id}}, \mathbf{z}_{\text{exp}}, \mathbf{z}_{\text{alb}}, \mathbf{z}_{\text{ill}}) \mapsto (\sigma, F), \quad (2)$$

where $\theta$ represents the network parameters, and the network architecture of the implicit function is shown in Fig. 2. $\gamma(*)$ is the pre-defined positional encoding from NeRF [35].

$\mathbf{x} \in \mathbb{R}^3$ is a 3D point sampled from one ray. Like previous works [16, 38, 63], instead of directly predicting $\mathbf{x}$'s RGB, we predict a high-dimensional feature vector $F(\mathbf{x}) \in \mathbb{R}^{256}$ for the 3D sampling point $\mathbf{x}$. Specifically, $h_\theta$ takes as input the concatenation of $\gamma(\mathbf{x})$, $\mathbf{z}_{id}$, $\mathbf{z}_{exp}$ and output the density $\sigma$ of $\mathbf{x}$ and an intermediate feature, the latter and $\mathbf{z}_{alb}$, $\mathbf{z}_{ill}$ are used to further predict $F(\mathbf{x})$. Thus, the prediction of the density field is mainly affected by the identity and expression code. The albedo and illumination codes only affect the feature vector prediction, which is consistent with our previous description. Similarly, we remove the viewing direction to avoid capturing undesired inconsistencies caused by dataset bias [16, 63].

Then a low-resolution 2D feature map $I_F \in \mathbb{R}^{256 \times 32 \times 32}$ can be obtained by the following volume rendering strategy:

$$I_F(r) = \int_0^\infty w(t) \cdot F(r(t)) dt,$$

$$\text{where} \quad w(t) = exp(-\int_0^t \sigma(r(s))ds) \cdot \sigma(r(t)). \quad (3)$$

$r(t)$ represents a ray emitted from the camera center. Finally, we map $I_F$ to the final predicted image $I \in \mathbb{R}^{3 \times 256 \times 256}$ with a 2D neural rendering module $\pi_\psi$, which is mainly composed of 1x1 Conv2D and leaky ReLU [33] activation layer to alleviate possible multi-view inconsistent artifacts [16]. $\psi$ denotes the learnable parameters. Similar with the strategy used in StyleNeRF [16], the resolution of $I_F$ is gradually increased through a series of upsampling layers. The upsampling process can be formulated as:

$$\text{Upsample}(X) = \text{Conv2D}(Y, K)$$
$$Y = \text{Pixelshuffle}(\text{repeat}(X, 4) + \beta_\theta(X), 2), \quad (4)$$

where $\beta_\theta : \mathbb{R}^D \to \mathbb{R}^{4D}$ is a learnable 2-layer MLP, and $K$ is a fixed blur kernel [62]. Like GIRAFFE [38], we map each feature tensor to an RGB image and take the sum of all RGB as the final predicted image. The difference is that we use 1x1 convolution instead of 3x3 convolution to avoid possible multi-view inconsistencies [16]. The network architecture of 2D neural rendering module is shown in Fig. 2.

### 3.3. Latent Codes and Canonical Coordinate

To efficiently train HeadNeRF, we utilize the 3DMM to initialize the latent codes of each image of our training dataset. Specifically, we set the dimensionality of the latent codes of HeadNeRF to be the same with the dimensionality of the corresponding codes from 3DMM [51] and initialize them by solving inverse rendering optimization [10, 54] based on the 3DMM model. Although the initial identity code of 3DMM only describes the geometry of the face area (without hair, teeth, etc.), it will be adaptively adjusted through the backpropagation gradient of training.



Figure 3. Some human head images used for training HeadNeRF.

On the other hand, the images used for building our model come from different channels (See sec. 3.4). To stabilize the training of HeadNeRF, we need to align each image's underlying geometry to a similar center before training. To this end, for each image, we can solve the above-mentioned 3DMM parameter optimization to obtain its corresponding global rigid transformation $T \in \mathbb{R}^{4 \times 4}$, which transforms the 3DMM geometry of the image from 3DMM canonical coordinate to camera coordinate. We take this transformation as the camera extrinsic parameter of the image. This strategy actually implicitly aligns the underlying geometry of each image to the center of the 3DMM template mesh.

### 3.4. Datasets and Preprocessing

We collected and processed three datasets to train HeadNeRF, and the details are described in the following.

**FaceSEIP Dataset.** This dataset includes 51 subjects with different genders, ages, races, and illumination conditions. These subjects are photographed in their daily dress-up and asked to perform 25 specific expressions with 13 cameras and 4 lighting conditions. This dataset contains 66300 face images, and some instances are shown in Fig. 3. Besides, the face mask is generated by the off-the-shelf segmentation methods [24, 61] and the images that fail to segment foreground are manually removed.

**FaceScape Dataset.** This dataset is from FaceScape [57] and contains 359 valid subjects. Each subject wears a hood and is asked to perform 20 specific expressions. Since the camera number of each subject is not fixed, for the convenience of training, we select subjects with 10 common perspectives from all subjects and filter out the rest. We also adopt the above-mentioned strategy to generate the face mask for each image. In the final, it contains 124 valid subjects, and some instances are shown in Fig. 3.

**FFHQ Dataset.** This dataset is from StyleGAN [22] and

4

contains 70000 high-resolution face images. The purpose of using this dataset is to utilize various in-the-wild face images to enhance the generalization ability of HeadNeRF. We also apply the above segmentation strategy to generate the segmentation mask and then manually select 4133 images with good masks. As shown in Fig. 8, this single-view in-the-wild image dataset effectively improves the generalization ability of HeadNeRF.

### 3.5. Loss Function

All the above-mentioned images are used to train Head-NeRF. The learnable variables include the latent codes of each image and the shared network parameters of volume rendering and neural rendering. The loss terms used to train the model include the following three terms.

**Photometric Loss.** For each image, it is required that the rendered result of the head area to be consistent with the corresponding real image, this loss term is formulated as:

$$L_{\text{data}} = \|M_{\text{h}} \odot (\mathcal{R}(\mathbf{z}_{\text{id}}, \mathbf{z}_{\text{exp}}, \mathbf{z}_{\text{alb}}, \mathbf{z}_{\text{ill}}, P) - I_{\text{GT}})\|^2, \quad (5)$$

where $\mathcal{R}(\mathbf{z}_{\text{id}}, \mathbf{z}_{\text{exp}}, \mathbf{z}_{\text{alb}}, \mathbf{z}_{\text{ill}}, P)$ is the rendered image, $I_{\text{GT}}$ is the corresponding real image, $M_{\text{h}}$ is the head mask and $\odot$ indicates Hadamard product operator.

**Perceptual Loss.** Compared with the vanilla NeRF, Head-NeRF can directly predict the color of all pixels in the rendered image via one inference. Therefore, we adopt the perceptual loss [20] in Eq. (6) to further improve the image details of the rendered results.

$$L_{\text{per}} = \sum_i \|\phi_i(\mathcal{R}(\mathbf{z}_{\text{id}}, \mathbf{z}_{\text{exp}}, \mathbf{z}_{\text{alb}}, \mathbf{z}_{\text{ill}}, P)) - \phi_i(I_{\text{GT}})\|^2, \quad (6)$$

where $\phi_i(*)$ denotes the activations of the $i$-th layer in VGG16 [46] network. As shown in Fig. 6, the perceptual loss can significantly improve the details of rendered results.

**Disentangled Loss.** In order to achieve semantically disentangled control to the rendered results, we let all images of one subject share the same identity latent code, and the images of the same expression with different lighting conditions and different captured cameras from the same subject share the same expression latent code.

For the purpose of disentangled representation, different subjects with similar expression should have similar expression codes. The initial expression codes generated by the initialization method of Sec.3.3 satisfy this requirement thanks to the introduction of 3DMM. Therefore, we require that the learnable expression code cannot be far away from the initial expression code. Other attributes are also constrained similarly. As 3DMM mainly models the face area without hair, teeth, etc., we relax this constraint for these non-expression attributes. Accordingly, the loss term for
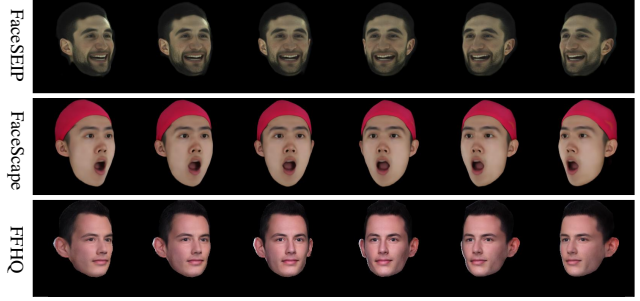


Figure 4. Disentangled control on camera pose. HeadNeRF can directly control the generated images' rendering pose and synthesize high fidelity rendered results with excellent multi-view consistency.

disentanglement is designed as:

$$L_{\text{dis}} = w_{\text{id}}\|\mathbf{z}_{\text{id}} - \mathbf{z}_{\text{id}}^0\|^2 + w_{\text{exp}}\|\mathbf{z}_{\text{exp}} - \mathbf{z}_{\text{exp}}^0\|^2 +$$
$$w_{\text{alb}}\|\mathbf{z}_{\text{alb}} - \mathbf{z}_{\text{alb}}^0\|^2 + w_{\text{ill}}\|\mathbf{z}_{\text{ill}} - \mathbf{z}_{\text{ill}}^0\|^2, \quad (7)$$

where $\mathbf{z}_*$ is the learnable latent code and $\mathbf{z}_*^0$ is the initial latent code from the 3DMM model.

In summary, the overall loss of HeadNeRF is defined as:

$$L = L_{\text{data}} + L_{\text{per}} + L_{\text{dis}}. \quad (8)$$

## 4. Experiments

### 4.1. Implementation Details

We implement HeadNeRF with Pytorch [40] and the learnable parameters are updated using Adam solver [25] on 3 NVIDIA 3090 GPUs. The sizes of different latent codes are $\mathbf{z}_{\text{id}} \in \mathbb{R}^{100}$, $\mathbf{z}_{\text{exp}} \in \mathbb{R}^{79}$, $\mathbf{z}_{\text{alb}} \in \mathbb{R}^{100}$, and $\mathbf{z}_{\text{ill}} \in \mathbb{R}^{27}$ respectively. In the volume rendering process, 64 points are sampled for each ray. In addition, we remove the hierarchical volume sampling of NeRF to further speed up inference. As a result, our HeadNeRF is able to render one head image more than 40fps without any other specific acceleration or optimization. The loss weights in Eq. (7) are set to $w_{\text{id}} = w_{\text{alb}} = w_{\text{ill}} = 0.001, w_{\text{exp}} = 0.1$.

### 4.2. Evaluations

**Disentangled Control.** In this part, we test HeadNeRF's ability to independently control various semantic attributes of rendered results. First, we train HeadNeRF with the datasets mentioned in Sec. 3.4. As shown in Fig. 4, for a given combination of latent codes $(\mathbf{z}_{\text{id}}, \mathbf{z}_{\text{exp}}, \mathbf{z}_{\text{alb}}, \mathbf{z}_{\text{ill}})$, we can directly adjust the camera parameters to continuously change the rendering view of the face image. These rendering results have excellent multi-view consistency, and the identity, expression, and appearance attributes of the rendered object are also well maintained.

Furthermore, we can utilize HeadNeRF to directly edit semantic attributes of rendered results. As shown in Fig. 5,
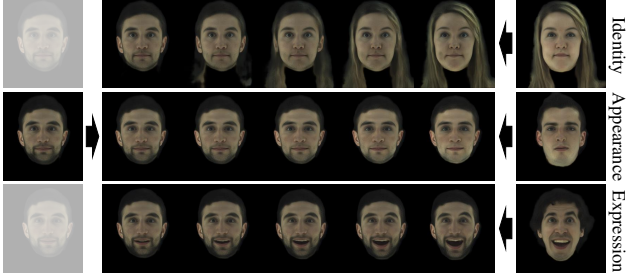
Figure 5. Disentangle facial attributes. HeadNeRF can independently edit the specific attributes of the rendered result by performing linear interpolating on the latent codes of the specified attributes.

we first sample several combinations of latent codes and render their corresponding frontal face images. Then two latent codes describing the same attribute are selected and the intermediate latent code can be obtained by performing linear interpolation on them. Finally, we update the interpolation results to the relevant attributes and use HeadNeRF to re-render the head image. As shown in this figure, HeadNeRF can maintain the remaining semantic attributes when editing a specific attribute, which verifies that HeadNeRF effectively disentangles different facial attributes.

**Ablation Study on the Perceptual Loss.** We attempt to remove the perceptual loss from the total loss, and record the corresponding trained model as HeadNeRF-noPerc. Then, for a given image in the training dataset, we use the trained HeadNeRF and HeadNeRF-noPerc to generate their prediction results respectively. As shown in Fig. 6, it can be found that the perceptual loss does effectively enhance the fine-level details of the generated results.

**Ablation Study on 2D Neural Rendering.** To verify the effectiveness of 2D neural rendering module in HeadNeRF, we design the following baseline version. All the latent codes, viewing direction and hierarchical sampling are kept except the 2D neural rendering module, and we denote this baseline as HeadNeRF-vanilla. Noting that the baseline can only be trained using the way of sampling batch rays due to the expensive rendering process of NeRF. Thus, the perceptual loss is not available to the HeadNeRF-vanilla. For a fair comparison, our HeadNeRF also removes the perceptual loss. Then, we use FaceSEIP dataset (See sec. 3.4) to train HeadNeRF-vanilla and HeadNeRF, respectively. HeadNeRF-vanilla is trained about 7 days with 3 NVIDIA 3090 GPUs, and HeadNeRF is trained about 3 days with a single NVIDIA 3090 GPU. As shown in Fig. 7, the results of HeadNeRF-vanilla tend to be blurred, which may be caused by its inefficient training. In contrast, thanks to the effectiveness and efficiency endowed by the neural rendering module, HeadNeRF can use much less training time while achieving better rendered result. Meanwhile, HeadNeRF-vanilla takes ∼5s to render a frame image, while HeadNeRF can render the result in real-time.
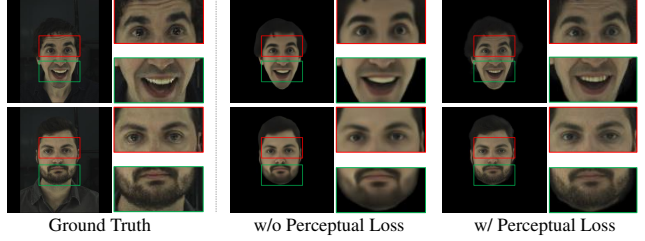


Figure 6. Ablation study on perceptual loss. The perceptual loss effectively enhances the fine-level details(wrinkle in red area and beard in green area) of the generated results.
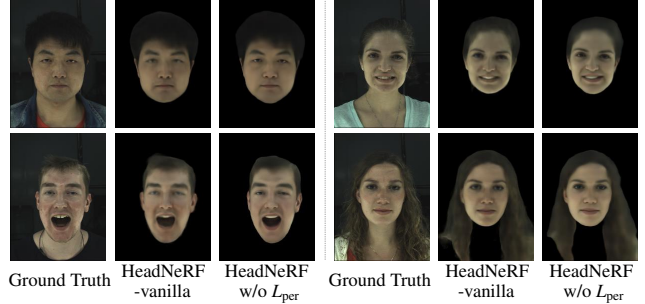


Figure 7. Ablation study on 2D neural rendering. Due to the effectiveness and efficiency endowed by the 2D neural rendering module, HeadNeRF can use much less training time while achieving better rendered results.

**Ablation Study on Using FFHQ Dataset.** To verify the importance of FFHQ dataset, we train two parametric models for comparison, where one is only trained with FaceSEIP dataset and the other is trained with Face-SEIP and FFHQ dataset. These two models are denoted as HeadNeRF-SEIP and HeadNeRF-SEIP+FFHQ, respectively. Images which are not used in the training are used for test. The following objective function Eq. (9) is optimized to obtain the latent codes embedding of the image in the above two models. The network parameters are fixed during optimization.

$$L_{\text{fitting}} = L_{\text{data}} + L_{\text{per}}. \tag{9}$$

The fitting results are shown in Fig. 8. It can be found that the introduction of FFHQ dataset significantly promotes the generalization ability of HeadNeRF. It is worth noting that we can further semantically modify the specified attributes of the optimization result based on our HeadNeRF, such as adjusting the rendering pose and changing the expression, etc. Meanwhile, the driven results are plausible and maintain excellent multi-view consistency. It shows that our HeadNeRF has successfully learned the statistical priors of human head from FaceSEIP dataset so that reasonable semantic editing results can be generated. In summary, the datasets used in our model training are complementary and effectively improve the representation ability of HeadNeRF.

Fitting | Editing



Input     w/o FFHQ     w/ FFHQ       Yaw     Pitch      Identity    Expression    Albedo    Illumination
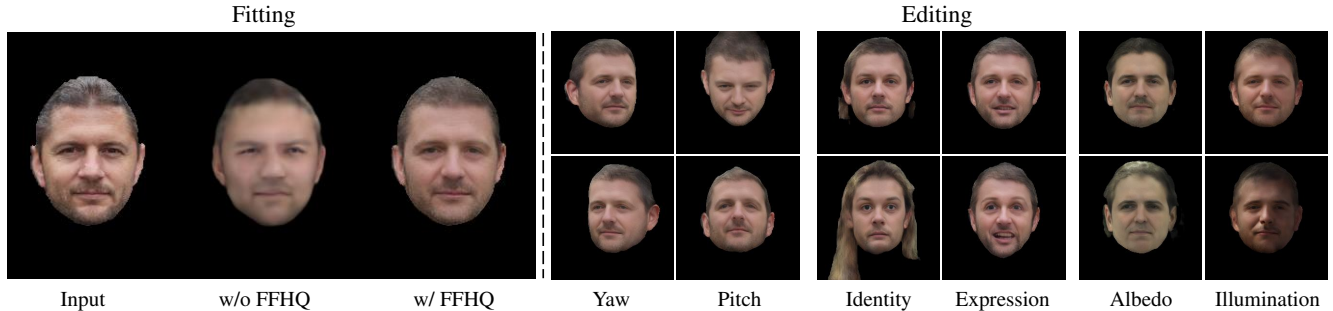
Figure 8. Ablation study on using FFHQ dataset. The introduction of FFHQ dataset significantly improves the generalization ability of HeadNeRF. Based on HeadNeRF, we can modify the specified attributes of the optimization result only with one single image as input, such as adjusting the rendering pose and changing the identity, expression and appearance of the rendered result.
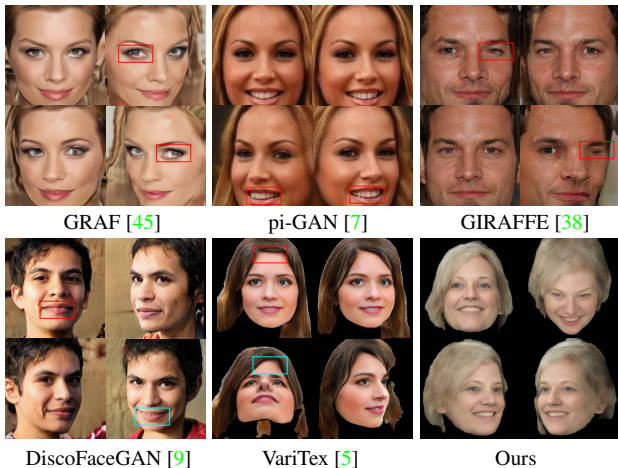


GRAF [45]       pi-GAN [7]      GIRAFFE [38]

DiscoFaceGAN [9]     VariTex [5]       Ours

Figure 9. Qualitative comparison with state-of-the-art methods. The undesired or inconsistent parts of other methods are marked with rectangles.

## 4.3. Comparisons

**Qualitative Comparison.** Fig. 9 shows the qualitative comparison between HeadNeRF and some related state-of-the-art methods. Each method generates several face images with different camera poses according to their sampled noise code. Among them, pi-GAN [7], GRAF [45] and GIRAFFE [38] are the generative models based on NeRF structure. It can be observed that the generated results of GRAF have obviously inconsistent artifacts, which may be caused by GRAF's inability to globally enforce the rendering results to meet the specified data distribution. Although the images generated by pi-GAN improve visual consistency, they still lack details. Because of the use of discriminant loss and the pure NeRF structure, pi-GAN cannot be trained at high resolution. GIRAFFE also designs a 2D neural rendering module to accelerate the rendering process. However, as shown in this figure, GIRAFFE doesn't get rid of the problem of multi-view inconsistency, which may be caused by the 3x3 CNN in GIRAFFE damaging the 3D information encoded by the NeRF.

It needs to be pointed out that these methods are designed

|  | CelebAMask-HQ | | | FFHQ | | |
|---|---|---|---|---|---|---|
| Method | L1 ↓ | PSNR ↑ | SSIM ↑ | L1 ↓ | PSNR ↑ | SSIM ↑ |
| pi-GAN [7] | 0.543 | **24.8** | **0.799** | 0.483 | **24.9** | **0.810** |
| GIRAFFE [38] | 0.420 | 20.6 | 0.628 | 0.428 | 20.3 | 0.635 |
| Ours | **0.306** | 19.6 | 0.702 | **0.344** | 21.6 | 0.755 |

Table 1. Quantitative comparison with state-of-the-art generative adversarial models based on NeRF structure. Some metrics about fitting a single image are calculated.



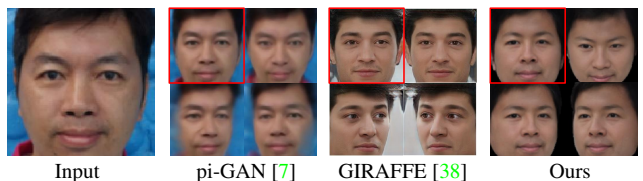Input      pi-GAN [7]     GIRAFFE [38]      Ours

Figure 10. Fitting results of different methods. The figure with red box is the fitting result, and others are generated by changing the rendering pose of the fitting result.

for general objects, and they cannot be used to semantically edit and control various attributes of the generated results. Some other methods like VariTex [5] and DiscoFaceGAN [9] attempt to integrate the priors from 3DMM into 2D generative models, but their rendering results still suffer from multi-view inconsistencies due to the limited representation ability of 3DMM. In contrast, HeadNeRF can generate high fidelity images in real-time while maintaining excellent multi-view consistency.

**Quantitative Comparison.** To quantitatively evaluate the effectiveness of HeadNeRF, we randomly select 400 in-the-wild face images from the CeleAMask-HQ dataset [29] and the test dataset of FFHQ [21], respectively. These images are used as the evaluation dataset and did not participate in the training of HeadNeRF at all. For each image in the evaluation dataset, the fitting results of pi-GAN [7], GIRAFFE [38] and our HeadNeRF are obtained by solving the inverse rendering optimization. For pi-GAN, the official fitting code and pre-training model are used. For GIRAFFE, we use the official pre-trained model, and the fitting code is implemented by ourselves. Then, the mean $L_1$ distance, Peak Signal-to-Noise Ratio (PSNR), and the Structure Similarity Index (SSIM) is calculated between the input image
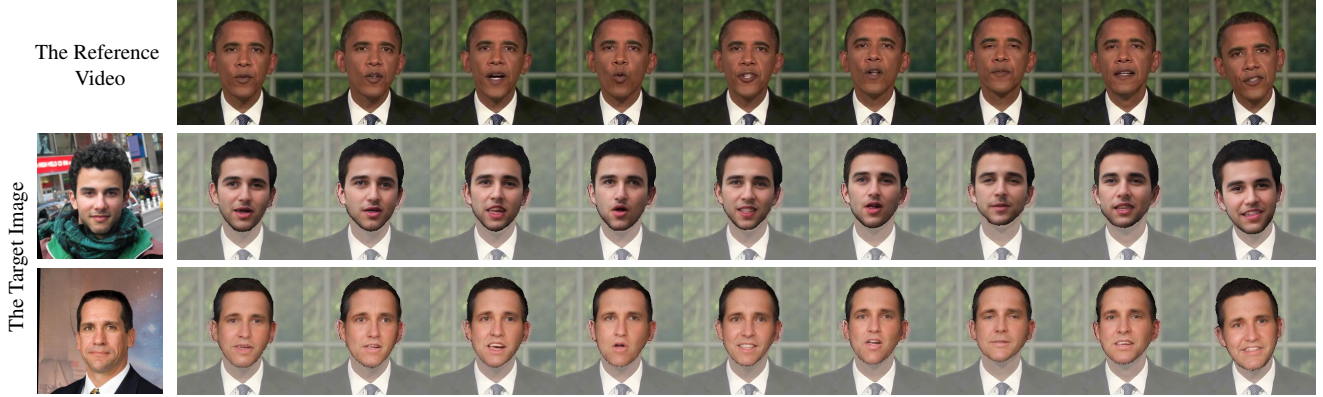
Figure 11. Expression Transfer. HeadNeRF is used to transfer facial expressions from the reference video to the persons in target images.

and fitting results. The statistical results are shown in Tab. 1.

Although the PSNR and SSIM of pi-GAN are optimal, we find that the fitting results of pi-GAN are often prone to overfitting. As shown in Fig. 10, the fitting result of pi-GAN is indeed visually optimal. However, if we edit and change the rendering pose of the fitting result, the rendered results often become blurred and damaged. On the other hand, GI-RAFFE can generate plausible rendered results when we change the rendering pose of GIRAFFE's fitting result, but the multi-view consistency of the rendered results is undesired. Compared with GIRAFFE, HeadNeRF has higher fitting accuracy and can maintain better multi-view consistency. Please note that our HeadNeRF can further modify and edit other semantic attributes, including identity, expression and appearance. As shown in Fig. 8. In addition, we want to point out that the insufficient training data may prevent our HeadNeRF from thoroughly exhibiting its advances. We will discuss this issue further in Sec. 5.

## 4.4. Application: Expression Transfer

As HeadNeRF has strong representation ability and can disentangle various attributes of the rendered result, it can be used for many applications such as novel-view synthesis, style mixing, etc,. In this part, we utilize HeadNeRF to perform expression transfer, i,e., transfer the facial expressions from the reference video to the people in target picture. To this end, we only need to extract the latent codes of all images from the reference video and the target image, and replace the expression latent code of the target image with the expression latent codes from the reference video. Finally, the trained HeadNeRF is employed to generate the desired face image sequence where the characters in the target image are driven to make expressions from the reference video. The qualitative results are shown in Fig. 11.

## 5. Limitation and Future Work

There still exist some limitations in HeadNeRF. Although images from FFHQ dataset are added to enhance



(a) Input   Fitting   (b) Continuously change the illumination latent code

Figure 12. Failure results of fitting and re-illumination.

the representation ability of HeadNeRF, the current training dataset is still not enough to cover various cases. For images that are quite different from our training data, HeadNeRF can only return similar results for the fitting task. As the examples shown in Fig. 12, because our training data rarely involves images with headgear, it is difficult to render the content of the headgear in the fitting results of HeadNeRF. In the future, we consider using a large amount of in-the-wild face image data to further enhance the representation ability of HeadNeRF in a self-supervised manner.

The current training dataset only contains four types of illuminations in our multi-view dataset (FaceSEIP), which is insufficient for the coverage of illumination types. Therefore, when we edit the illumination attribute, the change of image shading is not like a continuous movement of the light source position. In the future, we can consider adding data captured by the light stage to alleviate this problem.

## 6. Conclusion

In this paper, we proposed HeadNeRF, a novel NeRF-based parametric head model that integrates neural radiance field to parametric human head representation. Thanks to our well-designed network module and loss terms, Head-NeRF can render high fidelity head images in real-time and support directly controlling the pose of rendered images and independently editing the identity, expression, and appearance of generated images. Extensive experimental results have demonstrated that HeadNeRF outperforms state-of-the-art related models. We believe that HeadNeRF has taken a significant step toward the realistic digital human.

# References

[1] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *IEEE/CVF International Conference on Computer Vision*, 2021. 2

[2] Alexander W Bergman, Petr Kellnhofer, and Gordon Wetzstein. Fast training of neural lumigraph representations using meta learning. In *Advances in Neural Information Processing Systems*, 2021. 2

[3] 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. 1, 2, 3

[4] James Booth, Anastasios Roussos, Allan Ponniah, David Dunaway, and Stefanos Zafeiriou. Large scale 3d morphable models. *International Journal of Computer Vision*, 126(2):233–254, 2018. 3

[5] Marcel C Bühler, Abhimitra Meka, Gengyan Li, Thabo Beeler, and Otmar Hilliges. Varitex: Variational neural face textures. In *IEEE/CVF International Conference on Computer Vision*, 2021. 1, 2, 7

[6] Chen Cao, Yanlin Weng, Shun Zhou, Yiying 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. 1, 2, 3

[7] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5799–5809, 2021. 3, 7

[8] Julian Chibane, Aayush Bansal, Verica Lazova, and Gerard Pons-Moll. Stereo radiance fields (srf): Learning view synthesis for sparse views of novel scenes. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7911–7920, 2021. 2

[9] Yu Deng, Jiaolong Yang, Dong Chen, Fang Wen, and Xin Tong. Disentangled and controllable face image generation via 3d imitative-contrastive learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5154–5163, 2020. 1, 2, 7

[10] Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. 4

[11] Bernhard Egger, William AP Smith, Ayush Tewari, Stefanie Wuhrer, Michael Zollhoefer, Thabo Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, et al. 3d morphable face models—past, present, and future. *ACM Transactions on Graphics (TOG)*, 39(5):1–38, 2020. 1, 2

[12] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8649–8658, 2021. 2

[13] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *IEEE/CVF International Conference on Computer Vision*, 2021. 2

[14] 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. 1, 2

[15] Robin Green. Spherical harmonic lighting: The gritty details. In *Archives of the game developers conference*, volume 56, page 4, 2003. 2

[16] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis. *arXiv preprint arXiv:2110.08985*, 2021. 2, 3, 4

[17] Yudong Guo, Keyu Chen, Sen Liang, Yongjin Liu, Hujun Bao, and Juyong Zhang. Ad-nerf: Audio driven neural radiance fields for talking head synthesis. In *IEEE/CVF IEEE Conference on Computer Vision*, 2021. 2

[18] Peter Hedman, Pratul P Srinivasan, Ben Mildenhall, Jonathan T Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. In *IEEE/CVF International Conference on Computer Vision*, 2021. 2

[19] Yue Jiang, Dantong Ji, Zhizhong Han, and Matthias Zwicker. Sdfdiff: Differentiable rendering of signed distance fields for 3d shape optimization. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1251–1261, 2020. 2

[20] 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. 5

[21] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. 1, 2, 3, 7

[22] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8110–8119, 2020. 1, 2, 3, 4

[23] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3907–3916, 2018. 2

[24] Zhanghan Ke, Kaican Li, Yurou Zhou, Qiuhua Wu, Xiangyu Mao, Qiong Yan, and Rynson W.H. Lau. Is a green screen really necessary for real-time portrait matting? *ArXiv*, abs/2011.11961, 2020. 4

[25] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *International Conference on Learning Representations (ICLR)*, 2015. 5

[26] Marek Kowalski, Stephan J. Garbin, Virginia Estellers, Tadas Baltrušaitis, Matthew Johnson, and Jamie Shotton. Config: Controllable neural face image generation. In *European Conference on Computer Vision (ECCV)*, 2020. 1, 2

[27] Samuli Laine, Janne Hellsten, Tero Karras, Yeongho Seol, Jaakko Lehtinen, and Timo Aila. Modular primitives for high-performance differentiable rendering. *ACM Transactions on Graphics*, 39(6), 2020. 2

[28] Christoph Lassner and Michael Zollhofer. Pulsar: Efficient sphere-based neural rendering. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1440–1449, 2021. 2

[29] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 7

[30] 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, 3

[31] Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, and Zhaoyang Lv. Neural 3d video synthesis. *arXiv preprint arXiv:2103.02597*, 2021. 2

[32] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In *IEEE/CVF International Conference on Computer Vision*, pages 7708–7717, 2019. 2

[33] Andrew L. Maas, Awni Y. Hannun, and Andrew Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *in ICML Workshop on Deep Learning for Audio, Speech and Language Processing*, 2013. 4

[34] Iacopo Masi, Yue Wu, Tal Hassner, and Prem Natarajan. Deep face recognition: A survey. In *SIBGRAPI conference on graphics, patterns and images (SIBGRAPI)*, pages 471–478. IEEE, 2018. 1

[35] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision*, pages 405–421. Springer, 2020. 2, 3

[36] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. Hologan: Unsupervised learning of 3d representations from natural images. In *IEEE/CVF International Conference on Computer Vision*, pages 7588–7597, 2019. 1, 3

[37] Thu Nguyen-Phuoc, Christian Richardt, Long Mai, Yong-Liang Yang, and Niloy J. Mitra. Blockgan: Learning 3d object-aware scene representations from unlabelled images. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 3

[38] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021. 2, 3, 4, 7

[39] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *ACM Transactions on Graphics*, 2021. 2

[40] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, 2019. 5

[41] Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas Vetter. A 3d face model for pose and illumination invariant face recognition. In *IEEE international conference on advanced video and signal based surveillance*, pages 296–301. Ieee, 2009. 1, 2, 3

[42] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9054–9063, 2021. 2

[43] Amit Raj, Michael Zollhöfer, Tomas Simon, Jason M. Saragih, Shunsuke Saito, James Hays, and Stephen Lombardi. Pixel-aligned volumetric avatars. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11733–11742, 2021. 2

[44] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *IEEE/CVF International Conference on Computer Vision*, 2021. 2

[45] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: generative radiance fields for 3d-aware image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 3, 7

[46] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations (ICLR)*, 2015. 5

[47] Steven L Song, Weiqi Shi, and Michael Reed. Accurate face rig approximation with deep differential subspace reconstruction. *ACM Transactions on Graphics (TOG)*, 39(4):34–1, 2020. 1

[48] Shih-Yang Su, Frank Yu, Michael Zollhoefer, and Helge Rhodin. A-nerf: Surface-free human 3d pose refinement via neural rendering. In *Advances in Neural Information Processing Systems*, 2021. 2

[49] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6142–6151, 2020. 2

[50] Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Tretschk, Yifan Wang, Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, Tomas Simon, Christian Theobalt, Matthias Niessner, Jonathan T. Barron, Gordon Wetzstein anad Michael Zollhoefer, and Vladislav Golyanik. Advances in neural rendering, 2021. 2

[51] Luan Tran and Xiaoming Liu. Nonlinear 3d face morphable model. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7346–7355, 2018. 2, 3, 4

[52] Mei Wang and Weihong Deng. Deep face recognition: A survey. *Neurocomputing*, 429:215–244, 2021. 1

[53] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2

[54] Xueying Wang, Yudong Guo, Zhongqi Yang, and Juyong Zhang. Prior-guided multi-view 3d head reconstruction. *IEEE Transactions on Multimedia*, 2021. 4

[55] Ziyan Wang, Timur Bagautdinov, Stephen Lombardi, Tomas Simon, Jason Saragih, Jessica Hodgins, and Michael Zollhofer. Learning compositional radiance fields of dynamic human heads. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5704–5713, 2021. 2

[56] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. Nex: Real-time view synthesis with neural basis expansion. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8534–8543, 2021. 2

[57] Haotian Yang, Hao Zhu, Yanru Wang, Mingkai Huang, Qiu Shen, Ruigang Yang, and Xun 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*, pages 601–610, 2020. 1, 3, 4

[58] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *arXiv preprint arXiv:2106.12052*, 2021. 2

[59] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenoctrees for real-time rendering of neural radiance fields. In *IEEE/CVF International Conference on Computer Vision*, 2021. 2

[60] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021. 2

[61] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *European conference on computer vision (ECCV)*, pages 325–341, 2018. 4

[62] Richard Zhang. Making convolutional networks shift-invariant again. In *International Conference on Machine Learning*, pages 7324–7334, 2019. 4

[63] Peng Zhou, Lingxi Xie, Bingbing Ni, and Qi Tian. CIPS-3D: A 3D-Aware Generator of GANs Based on Conditionally-Independent Pixel Synthesis. *arXiv preprint arXiv:2110.09788*, 2021. 3, 4

[64] Michael Zollhöfer, Justus Thies, Pablo Garrido, Derek Bradley, Thabo Beeler, Patrick Pérez, Marc Stamminger, Matthias Nießner, and Christian Theobalt. State of the art on monocular 3d face reconstruction, tracking, and applications. In *Computer Graphics Forum*, volume 37, pages 523–550. Wiley Online Library, 2018. 1