

Gaussian Déjà-vu: Creating Controllable 3D Gaussian Head-Avatars with Enhanced Generalization and Personalization Abilities

Peizhi Yan¹, Rabab Ward¹, Qiang Tang², Shan Du^{3*}

¹University of British Columbia {yanpz, rababw}@ece.ubc.ca

²Huawei Canada qiang.tang@huawei.com

³University of British Columbia (Okanagan) shan.du@ubc.ca

Abstract

Recent advancements in 3D Gaussian Splatting (3DGS) have unlocked significant potential for modeling 3D head avatars, providing greater flexibility than mesh-based methods and more efficient rendering compared to NeRF-based approaches. Despite these advancements, the creation of controllable 3DGS-based head avatars remains time-intensive, often requiring tens of minutes to hours. To expedite this process, we here introduce the “Gaussian Déjà-vu” framework, which first obtains a generalized model of the head avatar and then personalizes the result. The generalized model is trained on large 2D (synthetic and real) image datasets. This model provides a well-initialized 3D Gaussian head that is further refined using a monocular video to achieve the personalized head avatar. For personalizing, we propose learnable expression-aware rectification blendmaps to correct the initial 3D Gaussians, ensuring rapid convergence without the reliance on neural networks. Experiments demonstrate that the proposed method meets its objectives. It outperforms state-of-the-art 3D Gaussian head avatars in terms of photorealistic quality as well as reduces training time consumption to at least a quarter of the existing methods, producing the avatar in minutes. Project homepage: <https://peizhiyan.github.io/docs/dejavu>

1. Introduction

The creation of photorealistic and controllable 3D head avatars has become a popular topic due to their employment in video gaming, VR and AR, filmmaking, telepresence, and many other fields. Three important factors must be considered when creating a 3D head avatar: **efficiency**, **quality**, and **controllability**. Efficiency is related to both the training and the rendering processes. Quality refers to achieving

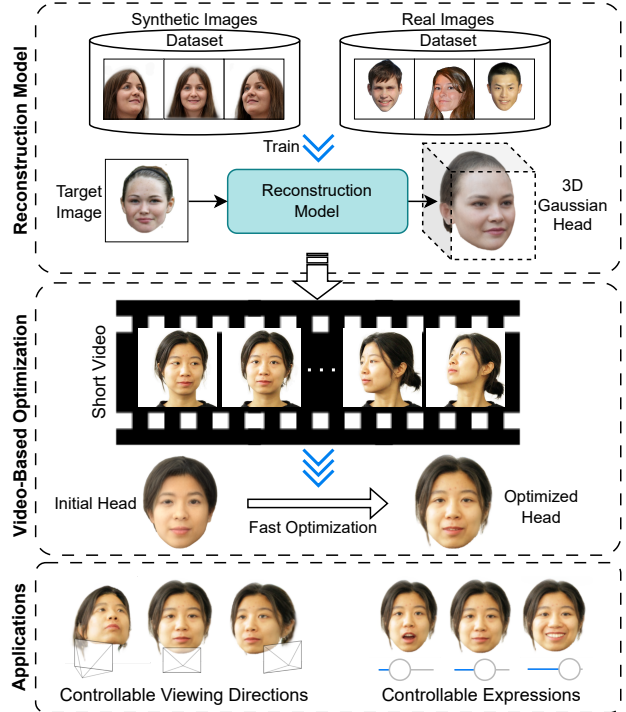


Figure 1. Gaussian Déjà-vu first trains a reconstruction model on large face image datasets and serves as a generalized base. This model initializes the 3D Gaussian head, which is then optimized to personalize the avatar to match the person in the video.

photorealistic results that resemble real human appearances. Controllability involves the ability to easily manipulate facial expressions, head poses, and camera views. Meeting these requirements ensures that the avatars can be effortlessly integrated into various applications. While existing methods may address some of these factors, they often fail to address all. To overcome this, we propose Gaussian Déjà-vu (Déjà-vu), which enables efficient, high-quality, and controllable 3D Gaussian head avatar creation. Table 1 summarizes the features supported by different 3D head

*Corresponding author: Shan Du.

Supported by funding: GR017752 and GR025942.

avatar creation methods.

	Mesh-based	NeRF-based	Existing 3DGSs	Ours
Single-Image Reconstruction	✓	✓	✗	✓
Easy Animation	✓	✗	✓	✓
Flexible Shape	✗	✓	✓	✓
Fast Training on Videos	N/A	✗	✗	✓
Efficient Rendering	✓	✗	✓	✓

Table 1. Comparison of supported features.

Mesh-based 3D Morphable Models (3DMMs) have been the foundation of existing 3D head avatars due to their efficiency in rendering and simplicity in animation [5, 6, 15, 43]. 3DMMs also support single-image-based reconstruction via either analysis-by-synthesis (fitting-based) [5, 45] or learning-based [11, 16, 31, 36, 57] methods to allow easy creation of a personal 3D head avatar. Despite their wide acceptance in industries, 3DMMs have a major drawback: the rigidity of mesh topology, making it challenging to model complex parts like hair.

Neural implicit-based methods [40, 41, 42], particularly the Neural Radiance Fields (NeRF) [41] model the 3D scene as a continuous volumetric representation field learned by a neural network, allowing for rendering complex scenes without the need for explicit 3D shapes. NeRF-based 3D face/head methods can be further categorized into conventional NeRF-based methods [8, 20, 21, 24, 49, 63] and StyleGAN+NeRF-based methods [1, 3, 7, 10, 25, 26, 39, 48, 55, 61]. Conventional NeRF-based methods mainly rely on multi-view images/videos as training data, while StyleGAN+NeRF-based methods leverage StyleGAN’s network design and learning strategy to enable the training on in-the-wild images. Although NeRF-based methods can provide photorealistic rendering results, the rendering efficiency of NeRF-based methods is still far from practical use, even with some techniques to improve the rendering efficiency [2, 7, 24, 48, 49, 55]. In addition, StyleGAN+NeRF-based methods often encounter flickering issues when used for generating video sequences.

3D Gaussian Splatting (3DGS) is the recent advancement in 3D representation and rendering [30]. It offers a promising way to model 3D head avatars, providing more flexibility than mesh-based methods and also more efficient rendering than NeRF-based methods. The core idea of 3DGS is to use 3D geometry primitives, namely 3D Gaussians, to represent 3D scenes, allowing for smooth and efficient rendering of complex scenes by optimizing the density and attributes of the 3D Gaussians. One important aspect of learning a 3DGS-based head avatar is the initialization of 3D Gaussians. Some works randomly initialize the locations of 3D Gaussians [9, 51], while most initialize the locations based on the 3D head mesh to make early-stage learning more efficient [13, 22, 32, 35, 37, 38, 44, 46, 47, 54, 56]. Most of these methods focus on the controllability in the an-

imation via training a personalized 3D Gaussian head avatar using the video of a person [9, 13, 22, 35, 37, 38, 44, 46, 47, 51, 54, 56]. However, a 3D Gaussian contains more information beyond its location, and only initializing the locations of the 3D Gaussians still leads to slow convergence. Most methods take tens of minutes or even hours to train an animatable 3D Gaussian head avatar for a single person, significantly limiting their wide application. We believe that a single-image-based reconstruction model can provide a good initialization of the 3D Gaussians. Nevertheless, existing 3DGS-based animatable head avatar creation methods do not support using a single 2D image.

In this work, we propose the “Gaussian Déjà-vu” framework (Déjà-vu) to efficiently create a controllable 3DGS-based head avatar. We parameterize the 3D Gaussians into a specialized UV map representation (a 2D representation of a 3D object’s surface), which we call the UV Gaussian map, aligned with the mesh-based FLAME model [34] to enable shape initialization and controllable animation. We train a reconstruction network to generate UV offsets that will be used to correct the FLAME-initialized 3D Gaussians [54] related to the single-image-based reconstruction task. The training primarily uses a synthetic dataset with a large number of identities, each captured from multiple viewpoints. We propose a training scheme that randomly pairs the input and target views of the same identity (person), along with a corresponding multi-view consistency regularization loss to enhance the 3D consistency. To further improve the generalization ability, we also fine-tune the network on real images. We only need to train the network once.

We first use our trained single-image-based reconstruction model to provide an initial 3D Gaussian head and then propose a rectification approach to enable further optimization of the head using a monocular video of the intended person. This rectification approach eliminates the need for a neural network by introducing learnable UV rectification maps, which are additional offsets to the initial 3D Gaussians. We propose expression-aware UV rectification blendmaps (blendmaps) to be learned using the given video. Different from 3D Gaussian blendshapes methods [13, 38], our proposed blendmaps only serve to rectify the initial 3D Gaussians, thus the optimization is faster than training a set of blendshapes from scratch. Due to our efficient design, Déjà-vu achieves a real-time animation speed of 220 frames-per-second (fps) at a resolution of 512×512 and a faster training speed than existing methods, reducing the training time to at least a quarter. Figure 1 provides a brief overview of the proposed framework.

In summary, our contributions are as follows:

- We propose Déjà-vu, a framework that first develops a single-image-based 3D Gaussian head reconstruction model to initialize a 3D Gaussian head, and further optimize the head using monocular videos to appear real and

gain personalization. This approach significantly reduces the training time on videos and produces state-of-the-art quality;

- We propose expression-aware UV rectification blendmaps, which allow easy control of facial animation without the need for a neural network;
- Our training strategy for the single-image-based reconstruction model involves the use of both synthetic and real images as training data, and also view-consistency regularization, which enhances generalization. The reconstruction model provides a robust initialization for subsequent video-based optimization. To our knowledge, this is the first single-image-based reconstruction model for the 3D Gaussian head.

2. Related Works

2.1. 3D Morphable Models in 3D Face Modeling

Over the past half-century, mesh-based methods have been the mainstream in 3D face modeling and animation [12]. A straightforward way to synthesize varying 3D faces is by linearly blending multiple template 3D faces, a method known as blendshapes [50]. Building on the concept of blendshapes, 3D Morphable Models (3DMMs) compress the 3D face shapes and textures from scanned data to principal components and allow the synthesis of a new 3D face from those principal components with given blending coefficients [5]. Earlier 3DMMs have only covered the facial region (including neck and ears) [5, 6, 43], while more advanced versions such as FLAME [34] model have addressed the entire head except for hair. Many 3DMM-based single-image 3D face reconstruction methods estimate the 3DMM coefficients to reconstruct the 3D face [11, 16, 31, 36]. To leverage the spatial capabilities of convolutional neural networks, some reconstruction methods utilize UV maps as a middle representation [16, 17, 62]. The UV maps are 2D representations of the 3D model with predefined 3D-2D correspondences that facilitate the application of textures and detailed surface information onto the 3D geometry [18].

2.2. NeRF-Based 3D Face and Head Models

The original Neural Radiance Fields method (NeRF) employs a neural network to represent the static 3D scene learned from multi-view images [41]. MoFaNeRF adapts NeRF to model the human face with parametric control by incorporating facial codes akin to 3DMM coefficients as additional network inputs [63]. HeadNeRF enhances rendering speed through a coarse-to-fine strategy, utilizing NeRF to render low-resolution feature maps and a 2D neural rendering network for producing high-resolution head images [24]. 3DMM-RF leverages 3DMM-based synthetic face images for training, achieving more accurate fitting over the facial region [21]. Inspired by the success of 2D

StyleGANs [27, 29] in generating photorealistic face images, StyleNeRF integrates a style-based network architecture into NeRF and trains the network through adversarial learning as in StyleGAN [23]. EG3D further refines the network architecture by forming a tri-plane representation for volume rendering of the raw image [7]. Due to efficiency considerations, the tri-plane generally has low resolution, prompting the use of a 2D super-resolution module to enhance the final image resolution. This tri-plane NeRF scheme is now widely adopted by numerous 3D-aware head models [1, 10, 26, 39, 48, 55]. Although StyleGAN+NeRF methods generate photorealistic rendered head images, they can introduce artifacts, such as noticeable flickering during animations. Moreover, these models do not support the customization of the reconstructed head based on video data of a person, limiting their applicability in more personalized head avatar scenarios. In contrast, the Déjà-vu framework we propose here produces an explicit 3D Gaussian-based head avatar that can be rendered directly through 3DGS [30] and optimized using video data to enhance personalization.

2.3. 3D Gaussian-Based Head Avatars

Like the original NeRF, the 3D Gaussian Splatting (3DGS) is designed to represent and render static 3D scenes [30]. To achieve controllability, 3DGS-based head avatar methods typically employ 3DMM. GaussianAvatars introduces a technique to rig the 3D Gaussians to the FLAME model, where the Gaussians are anchored to the triangle facets of the mesh [44]. Similarly, SplattingAvatar proposes a method to embed 3D Gaussians into any head mesh model [47]. SplatFace aligns the 3D Gaussians with the underlying head mesh through a non-rigid process and jointly optimizes both the Gaussians and the mesh [37]. Numerous 3DGS head avatar methods leverage neural rendering to achieve photorealism by incorporating high-dimensional latent features into each 3D Gaussian and utilizing neural networks to decode these features to produce the final images [4, 13, 22, 32, 35, 46, 51, 56]. For instance, GaussianHead and Gaussian Splatting Decoder employ a tri-plane representation to learn the 3D Gaussian latent features [4, 51]. HeadGaS utilizes a Multi-layer Perceptron (MLP) network to decode the latent features of each 3D Gaussian into color and alpha blending weights [13]. Additionally, NPGA and Gaussian Head Avatar initially render a 2D feature map using a modified 3DGS, which is then decoded into an image through 2D neural rendering [22, 56]. However, these neural-rendering-based methods deviate from the simplicity of traditional 3DGS, resulting in reduced compatibility with existing 3D engines. In addition, some methods employ a “canonical + deformation” strategy which leverages a neural network, typically an MLP, to estimate the deformation for each 3D Gaussian [9, 35, 54]. Approaches like HeadGaS and 3D Gaussian Blendshapes integrate the con-

cept of blendshapes by learning a set of 3D Gaussians as expression bases [13, 38].

Despite the success of existing 3DGS-based head avatar methods, this work takes us one step further, introducing the first single-image-based reconstruction of a 3D Gaussian head avatar, which provides robust initialization for our video-based optimization. Our 3D Gaussian head avatar aligns with the original 3DGS and does not require a neural network for animation, ensuring compatibility and ease of integration with standard 3D engines.

3. Method

3.1. Background: 3D Gaussian Splatting

The 3D Gaussian Splatting (3DGS) presents a 3D Gaussians-based 3D representation and its corresponding splatting-based rendering process, allowing for differentiable and efficient rendering [30]. 3DGS modifies the Probability Density Function (PDF) of a multivariate Gaussian distribution to define the impact of a 3D Gaussian \mathbf{g}_i on a given 3D location $\mathbf{x} \in \mathbb{R}^3$:

$$G_i(\mathbf{x}) = e^{-\frac{1}{2}(\mathbf{x}-\mu_i)^\top \Sigma_i^{-1}(\mathbf{x}-\mu_i)}, \quad (1)$$

where $\mu_i \in \mathbb{R}^3$ is the mean (center of the 3D Gaussian) and $\Sigma_i \in \mathbb{R}^{3 \times 3}$ is the covariance matrix. To ensure the covariance matrix Σ_i to be positive semi-definite to have physical meaning, 3DGS computes the Σ_i based on a scaling matrix $S_i \in \mathbb{R}^{3 \times 3}$ and a rotation matrix $R_i \in \mathbb{R}^{3 \times 3}$:

$$\Sigma_i = R_i S_i S_i^\top R_i^\top. \quad (2)$$

This representation is analogous to defining the scaling and rotation of an ellipsoid [30]. Further, the scaling and rotation matrices can be derived by a scaling vector $\mathbf{s}_i \in \mathbb{R}^3$ and a quaternion $\mathbf{q}_i = (1, \mathbf{r}_i)$, where $\mathbf{r}_i \in \mathbb{R}^3$ defines the three rotation angles. The 3D Gaussian \mathbf{g}_i is then parameterized by $\mathbf{g}_i = \{\mu_i, \mathbf{s}_i, \mathbf{r}_i, \alpha_i, \mathbf{a}_i\}$, where α_i is the alpha blending weight, and \mathbf{a}_i is the vector of Spherical Harmonic (SH) coefficients for computing the view-dependent RGB color [19]. The 3D Gaussians are projected to 2D for rendering. The image formation model of 3DGS mirrors the volume rendering techniques utilized in NeRF [14, 41]:

$$\mathbf{C}(p) = \sum_{i \in \mathcal{N}_p} \mathbf{c}_i^p \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (3)$$

where p is the 2D-pixel location on the plane onto which 3D Gaussians are projected. Here, $\mathbf{C}(p)$ represents the final color of the pixel, \mathcal{N}_p is an ordered set of the 2D projections of the 3D Gaussians that cover the pixel, and \mathbf{c}_i^p is the computed color of the i^{th} 2D projection at location p [58].

3.2. Our UV Gaussian Map Representation

Our 3D Gaussian head avatar is based on FLAME [34] to gain controllability. We denote F as the FLAME shape model. The 3D head mesh can be constructed through: $m = F(\beta_{id}, \beta_{exp}, \beta_{jaw})$, where $\beta_{id} \in \mathbb{R}^{100}$, $\beta_{exp} \in \mathbb{R}^{50}$, and $\beta_{jaw} \in \mathbb{R}^3$ are the coefficients of identity, expression, and jaw pose respectively. We modify the original UV correspondences of FLAME to remove the neck and add the mouth interior shape (see supplementary material).

We define a set of UV maps with height and width of H and W to represent the parameters of 3D Gaussians. These UV maps include U_μ for the 3D Gaussian positions, U_s for the scales, U_r for the rotations, U_α for the alpha blending weights, and U_a for the SH coefficients. Given the relatively simple surface materials of the human head, we do not consider complex anisotropic reflections and refractions. Therefore, we use SHs with a degree of zero to represent RGB colors uniformly across all directions. The number of channels in U_a is thus equals to three. Our UV Gaussian map $U \in \mathbb{R}^{H \times W \times 13}$ is created by channel-wise concatenation of the above set of UV maps. We mask out the invalid pixels (undefined locations in the UV map) in the UV Gaussian map, and each valid pixel corresponds to a 3D Gaussian. Therefore the total number of 3D Gaussians is less than $H * W$.

To prepare an initial UV Gaussian map U_{init} , we rasterize the mesh shape m to initialize the Gaussian positions U_μ . We empirically initialize the values of alpha weights in U_α to 0.1 and the values of scales in U_s to -8.3533 (note that 3DGS applies exponential activation on the scales, to ensure the final scales are positive). All the other UV maps are initialized with zeros.

3.3. Single-Image-Based Reconstruction

Our single-image-based 3D Gaussian head reconstruction utilizes a convolutional encoder-decoder network to estimate the offsets to the initialized 3D Gaussians, see Figure 2 (a). The initial 3D Gaussian head is defined by using the initial UV Gaussian map, as detailed in Section 3.2. We implemented a FLAME tracking method to derive the FLAME coefficients and the camera pose given an image (see our homepage). Our goal is to estimate a set of offsets, referred to as UV offsets ΔU , which adjust the initial UV Gaussian map U_{init} . The encoder processes a head image I to produce low-resolution feature maps, which are subsequently input into the decoder. Additionally, the decoder incorporates the initial UV position map U_μ to enhance adaptation to the UV map layout. The corrected UV Gaussian map for the single-image-reconstruction phase is computed as $U_{SIR} = U_{init} + \Delta U$. Using the estimated camera pose for I , we render U_{SIR} through 3DGS, with the resulting image denoted by I'_{SIR} . The reconstruction loss, which quantifies

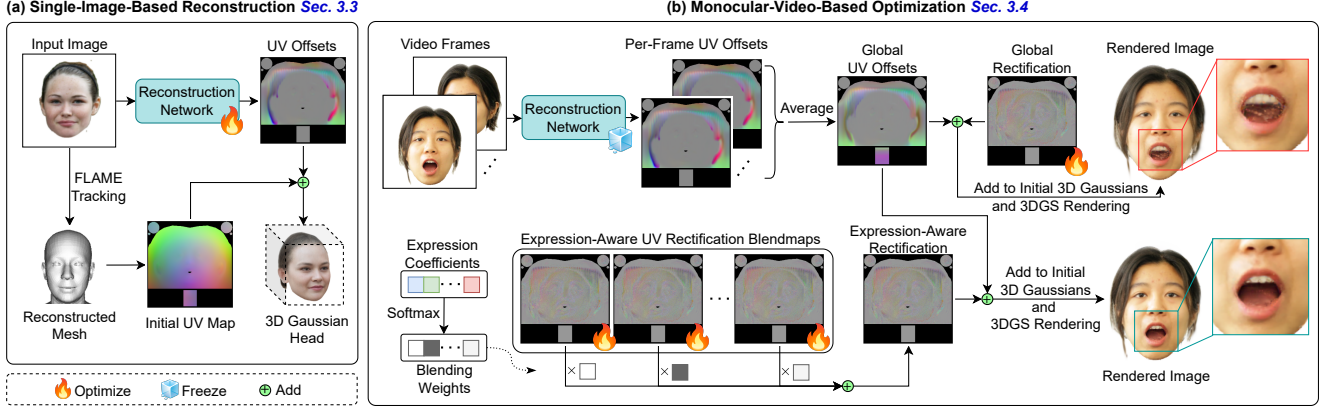


Figure 2. Detailed flowcharts for our single-image-based reconstruction (a) and monocular-video-based further optimization (b) processes.

the discrepancy between I and I'_{SIR} , is computed as follow:

$$\mathcal{L}_{\text{photo}} = \lambda_{\text{RGB}} \mathcal{L}_{\text{RGB}} + \lambda_{\text{LPIPS}} \mathcal{L}_{\text{LPIPS}} + \lambda_{\text{SSIM}} \mathcal{L}_{\text{SSIM}}, \quad (4)$$

where \mathcal{L}_{RGB} is the pixel-level L1 loss, $\mathcal{L}_{\text{LPIPS}}$ is the perceptual loss [59], $\mathcal{L}_{\text{SSIM}}$ is the differentiable SSIM-based loss [53], and the λ s are corresponding loss term weights.

To ensure that the positions of the 3D Gaussians do not deviate significantly from the original reconstructed FLAME shape, we introduce a position regularization term, denoted as \mathcal{L}_{μ} :

$$\mathcal{L}_{\mu} = \|\Delta U_{\mu} \circ M_{\mu}\|_2, \quad (5)$$

where ΔU_{μ} is the network estimated UV position offsets map, \circ denotes the element-wise Hadamard product operation, and M_{μ} is a predefined weights map. This map allocates more weight to the facial region and less to the scalp, ensuring that the facial area more closely conforms to the reconstructed FLAME shape, while permitting the hair, represented by the 3D Gaussians originally on the scalp, to maintain some positional variability. Similarly, we define the scale regularization term, denoted as \mathcal{L}_s to regularize the size of the Gaussians from growing excessively:

$$\mathcal{L}_s = \|\Delta U_s\|_2. \quad (6)$$

We introduce a view consistency regularization loss, $\mathcal{L}_{\text{view}}$, to ensure stability across different views of the same identity:

$$\mathcal{L}_{\text{view}} = \sum_{\Delta U \in \mathcal{B}} \|\Delta U - \Delta \bar{U}_B\|_2 / |\mathcal{B}|, \quad (7)$$

where \mathcal{B} represents a set of UV offsets from multi-view input images for the same identity, and $\Delta \bar{U}_B$ denotes the mean UV offsets within the set. To further enhance view consistency, we also randomly pair the input and target images from the same batch of multi-view images belonging to the same identity during training. The overall loss is:

$$\mathcal{L} = \mathcal{L}_{\text{photo}} + \lambda_{\mu} \mathcal{L}_{\mu} + \lambda_s \mathcal{L}_s + \lambda_{\text{view}} \mathcal{L}_{\text{view}}. \quad (8)$$

3.4. Monocular-Video-Based Optimization

Using a monocular video of a person, we apply our single-image-based reconstruction to each frame and then average the estimated UV offsets to obtain the mean UV offsets for the video, denoted as $\Delta \bar{U}$. Merely using $U_{\text{init}} + \Delta \bar{U}$ to form the reconstructed 3D Gaussian head may result in a lack of personalization. To address this, we propose a two-stage optimization method aimed at further enhancing the personalization and quality of the reconstructed 3D Gaussian head. See Figure 2 (b).

In the first stage (Figure 2 (b) top), we define the global UV rectification, ΔU_{global} as a tensor in $\mathbb{R}^{H \times W \times 13}$, which is initially set to zeros and made learnable. The UV Gaussian map, after applying this global UV rectification, is expressed as:

$$U_{\text{global}} = U_{\text{init}} + \Delta \bar{U} + \Delta U_{\text{global}}. \quad (9)$$

The objective of the first stage is to optimize ΔU_{global} to minimize \mathcal{L} in Equation 8 (excluding $\mathcal{L}_{\text{view}}$). After the first stage of optimization, the resulting 3D Gaussian head already closely resembles the person in the video. However, relying solely on ΔU_{global} can still lead to artifacts during animation, as this approach does not account for variations in facial expressions.

To address this limitation, we introduce a second stage (Figure 2 (b) bottom) that optimizes the UV rectification blendmaps $\Delta U_{\text{blend}} \in \mathbb{R}^{D \times H \times W \times 13}$, where D denotes the number of blending weights and $\Delta U_{\text{blend}}^i$ denotes the i^{th} blendmap. Our blending weights $b \in \mathbb{R}^D$ are derived from the estimated FLAME expression coefficients β_{exp} :

$$b = \text{softmax}(\beta_{\text{exp}}), \quad (10)$$

where b_i indicates the i^{th} blending weight. Utilizing the softmax function ensures the blending weights are positive values and have a sum of exactly 1. We initialize the

blendmaps with the ΔU_{global} learned in the first stage. The final UV Gaussian map after rectification is expressed as:

$$U_{\text{blend}} = U_{\text{init}} + \Delta \bar{U} + \sum_{i=1}^D b_i \Delta U_{\text{blend}}^i. \quad (11)$$

We optimize ΔU_{blend} with the same objective as in the first stage.

4. Experiments

4.1. Datasets

We prepare two datasets for training our single-image-based reconstruction model. The first is a synthetic dataset created using the 3D-aware face image generation model, PanoHead [1], which synthesized images for 18,000 identities across 25 pre-defined camera views (detailed in the supplementary material). The second dataset comprises real-face images from the FFHQ dataset [28] to improve the robustness of our network on real images. We filtered out images where faces are obscured by hats or clothing, or of low quality, resulting in 38,000 usable images. Of these, the first 1,000 images are set aside for evaluation, while the remaining 37,000 images are used for training. For video-based optimization, we utilize the dataset from PointAvatar [60], which includes monocular video clips of three subjects [60]. For each subject, we allocate the first 90% of the frames from each video clip for training purposes, reserving the remaining 10% for evaluation.

4.2. Implementation Details

We set the resolution of our UV Gaussian map to 320×320 in our experiment. The actual number of Gaussians is 74,083. The loss term weights are set as follows: $\lambda_{\text{RGB}} = 1.0$, $\lambda_{\text{LPIPS}} = 0.5$, $\lambda_{\text{SSIM}} = 0.05$, $\lambda_{\mu} = 0.5$, $\lambda_s = 5 \times 10^{-5}$, and $\lambda_{\text{view}} = 0.1$. The rendering resolution is 512×512 .

During the training of the single-image-based reconstruction network, each training step starts with training on a random batch from the synthetic dataset at a higher learning rate. This is followed by fine-tuning on a random batch from the real dataset at a reduced learning rate, specifically one-tenth of that used for the synthetic dataset. It's important to note that the $\mathcal{L}_{\text{view}}$ term is omitted during fine-tuning on the in-the-wild dataset, as multi-view images are not available. The training batch size is set to 16.

For video-based optimization, the distribution of steps between the first and second stages is split into a 3:7 ratio. The batch size for each step is 8. We utilize only the first ten expression coefficients to compute the blending weights, where $D = 10$.

All the experiments are performed on one Nvidia A6000 GPU (48GB VRAM). For additional details, please refer to the supplementary material.

4.3. Single-Image-Based Reconstruction Results

We compare the performance of our single-image-based reconstruction model with leading NeRF-based parametric face/head models, MoFaNeRF [63], HeadNeRF [24], and 3DMM-RF [21], due to their aligned goal with us on training a general model for face and head reconstruction. We exclude comparisons with StyleGAN-based or diffusion-based methods, as these primarily focus on image synthesis rather than creating a general head model. Since MoFaNeRF and 3DMM-RF only represent the facial region, we compare only the facial region with these methods. For a fair comparison, we employ the same set of 550 testing samples used by 3DMM-RF, extracted from the CelebAMask-HQ dataset [33]. We assess reconstruction performance using several metrics: pixel-level L1 distance, LPIPS [59], Structural Similarity Index (SSIM) [52], and Peak Signal-to-Noise Ratio (PSNR). Detailed results of these comparisons are presented in Table 2. Additionally, we conduct tests on a separate set of 1,000 images from the FFHQ dataset [28], with outcomes detailed in Table 3. Our method demonstrates competitive performance, frequently outperforming these benchmarks in most metrics. Figure 3 presents qualitative comparisons, showing that our method achieves a visual quality closely resembling the source image.

Method	Region	L1 ↓	LPIPS ↓	SSIM ↑	PSNR ↑
MoFaNeRF [63]	Facial	0.273	0.442	0.910	14.713
3DMM-RF* [21]		0.216	-	0.956	-
HeadNeRF [24]		0.072	0.113	0.940	24.937
Ours		0.033	0.061	0.935	28.755
HeadNeRF [24]	Head	0.354	0.431	0.693	15.429
Ours		0.125	0.211	0.740	20.525

* indicates results reported from the original paper (not open-source).
- indicates the value is missing or not comparable.

Table 2. Reconstruction performances (CelebAMask-HQ).

Method	Region	L1 ↓	LPIPS ↓	SSIM ↑	PSNR ↑
MoFaNeRF [63]	Facial	0.251	0.416	0.720	15.084
HeadNeRF [24]		0.087	0.151	0.920	23.858
Ours		0.040	0.071	0.922	28.030
HeadNeRF [24]	Head	0.211	0.310	0.798	18.794
Ours		0.088	0.155	0.813	23.184

Table 3. Reconstruction performances (FFHQ).

We compare our method with HeadNeRF to showcase the superior 3D view consistency of our 3D Gaussian-based head representation. Figure 4 presents comparisons, illustrating the rendering of the reconstructed head from various viewing angles. As shown, our method consistently outperforms HeadNeRF in maintaining view consistency, attributable to the explicit 3D structure represented by the 3D Gaussians. Notably, at extreme viewing angles, where



Figure 3. Qualitative comparison results on single-image-based reconstruction.

HeadNeRF’s renderings show significant deterioration, our 3D Gaussian-based head continues to maintain high-quality visual outputs.

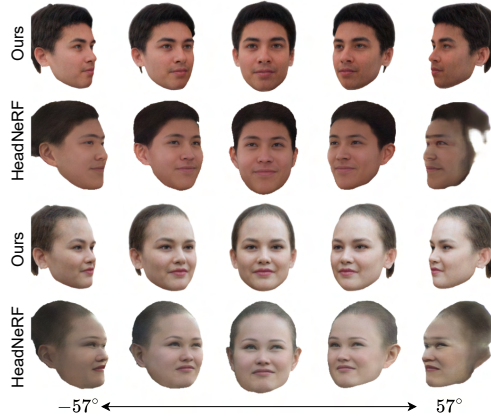


Figure 4. Qualitative comparison with HeadNeRF across varying viewing angles. Our method works even at extreme angles.

4.4. Monocular-Video-Based Optimization Results

In the monocular-video-based optimization, we compare our method with the state-of-the-art 3D Gaussian head avatar methods, GaussianAvatars [44] and FlashAvatar [54], both of which are based on the FLAME model. To ensure fairness, we use the same FLAME tracking results to obtain expression coefficients, the underlying mesh, and the 6DOF camera pose. Additionally, as FlashAvatar also utilizes a fixed number of 3D Gaussians, we use the same amount of 74,083 3D Gaussians as in our method. GaussianAvatars

incorporates dynamic density control from 3DGS, but each learned head model has around 60,000 to 100,000 3D Gaussians, which is similar to ours. Table 4 shows the quantitative comparison results. Our method outperforms both GaussianAvatars and FlashAvatar in terms of LPIPS, SSIM, and PSNR metrics across all three subjects, achieving the best results in the shortest training time of only 8 minutes (on a single Nvidia A6000 GPU). We demonstrate qualitative expression reconstruction results in Figure 5, with the first column displaying source expressions that are subsequently applied to the learned head avatars.

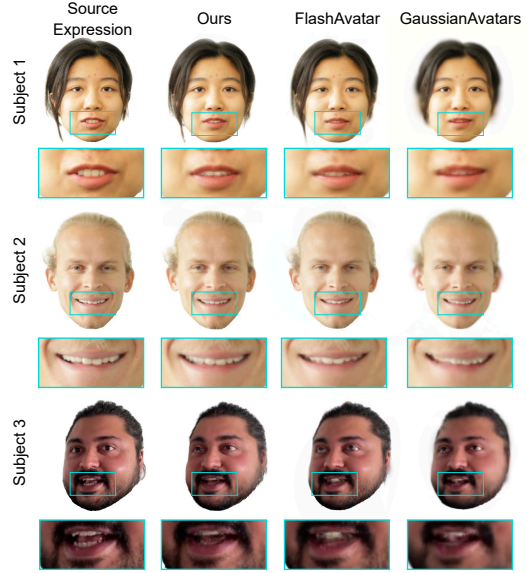


Figure 5. Expression reconstruction results.

We compare the convergence speed with FlashAvatar, which is known for its rapid training. Figure 6 displays the comparison. Our approach, utilizing a better initialization and a non-neural network strategy, achieves faster convergence than FlashAvatar. Ultimately, our method enables the efficient optimization of a high-density 3D Gaussian head avatar within minutes.

4.5. Ablation Studies

We first use an example to demonstrate the importance of our view consistency regularization loss $\mathcal{L}_{\text{view}}$, see Figure 7. The first row displays input images of the same identity from different camera views fed into our single-image-based reconstruction model. The second (w/o $\mathcal{L}_{\text{view}}$) and third (w/ $\mathcal{L}_{\text{view}}$) rows show the reconstructed heads rendered from a frontal view. As evident, without $\mathcal{L}_{\text{view}}$, the reconstruction quality is biased towards the side visible in the input image. However, with $\mathcal{L}_{\text{view}}$, the reconstructions appear nearly identical across different views of the same identity, showcasing the robustness of our approach in maintaining consistent visual outputs.

	Training Time	Subject 1			Subject 2			Subject 3		
		LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑
GaussianAvatars (CVPR'24) [44]	3 hrs	0.155	0.877	22.650	0.141	0.932	28.672	0.135	0.894	22.152
	6 hrs	0.155	0.879	22.999	0.141	0.933	29.094	0.134	0.895	22.101
FlashAvatar (CVPR'24) [54]	30 mins	0.104	0.867	23.565	0.091	0.917	28.447	0.096	0.872	22.028
	60 mins	0.097	0.875	23.767	0.091	0.917	28.561	0.091	0.874	21.973
Ours (global rectification only)	8 mins	0.082	0.883	23.157	0.073	0.931	29.083	0.077	0.883	21.487
Ours	8 mins	0.073	0.899	24.277	0.070	0.936	29.455	0.067	0.899	22.342

Table 4. Quantitative comparison with state-of-the-art 3D head avatar methods.



Figure 6. Convergence speed on video-based optimization.



Figure 7. Ablation results on view consistency regularization.

We demonstrate the critical role of our 3D Gaussian position and scale regularization losses in Figure 8. In this experiment, we trained two additional single-image-based reconstruction models, each omitting either the \mathcal{L}_μ (position loss) or \mathcal{L}_s (scale loss). These models were used to initialize 3D Gaussian heads, which were then further optimized using video data. The results clearly show that without adequate regularization of the position and scale of the 3D Gaussians, noticeable artifacts persist even after subsequent video-based optimization.

In Figure 9, we showcase the importance of our expression-aware blendmaps. We apply the expression estimated from the source image to our learned head avatars. The second image displays the head enhanced with our

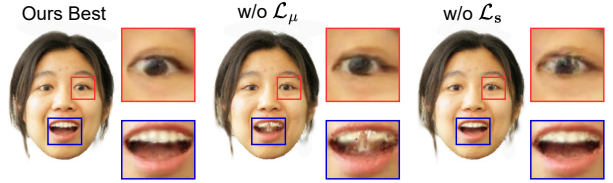


Figure 8. Ablation results on position and scale regularizations.

expression-aware UV blendmaps and the third illustrates the face with only global UV rectification. Both approaches use the same amount of time for optimization in this ablation experiment. It becomes clear that global UV rectification introduces artifacts when dealing with exaggerated expressions. Quantitative results, which further substantiate the limitations of using only global rectification, are presented in Table 4.

For more ablation study results please refer to our supplementary material.

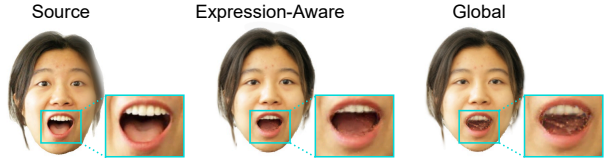


Figure 9. Ablation results on different levels of rectification.

5. Conclusion

In this work, we propose Déjà-vu, a novel framework designed to create controllable 3D Gaussian head avatars with fast training. Notably, Déjà-vu is the first method capable of reconstructing a 3D Gaussian head based solely on a single image as input and trained using 2D images. This reconstruction model provides robust initialization for our video-based optimization. We propose learnable UV blendmaps to adjust the head to have the desired expression, a solution that is both effective and quick to train. Our framework outperforms existing state-of-the-art methods in both rendering quality and training speed. In the future, we will enhance the adaptability of Déjà-vu to a wider range of facial expressions and explore more applications.

References

- [1] Sizhe An, Hongyi Xu, Yichun Shi, Guoxian Song, Umit Y Ogras, and Linjie Luo. Panohead: Geometry-aware 3d full-head synthesis in 360deg. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20950–20959, 2023. 2, 3, 6
- [2] Ziqian Bai, Feitong Tan, Sean Fanello, Rohit Pandey, Mingsong Dou, Shichen Liu, Ping Tan, and Yinda Zhang. Efficient 3d implicit head avatar with mesh-anchored hash table blendshapes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1975–1984, 2024. 2
- [3] Chong Bao, Yinda Zhang, Yuan Li, Xiyu Zhang, Bangbang Yang, Hujun Bao, Marc Pollefeys, Guofeng Zhang, and Zhaopeng Cui. Geneavatar: Generic expression-aware volumetric head avatar editing from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8952–8963, 2024. 2
- [4] Florian Barthel, Arian Beckmann, Wieland Morgenstern, Anna Hilsman, and Peter Eisert. Gaussian splatting decoder for 3d-aware generative adversarial networks. In *Workshop of IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 3
- [5] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 157–164. ACM New York, NY, USA, 2023. 2, 3
- [6] 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
- [7] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. 2, 3
- [8] Aggelina Chatziagapi, Grigorios G Chrysos, and Dimitris Samaras. Mi-nerf: Learning a single face nerf from multiple identities. *arXiv preprint arXiv:2403.19920*, 2024. 2
- [9] Yufan Chen, Lizhen Wang, Qijing Li, Hongjiang Xiao, Shengping Zhang, Hongxun Yao, and Yebin Liu. Monogaussianavatar: Monocular gaussian point-based head avatar. *arXiv preprint arXiv:2312.04558*, 2023. 2, 3
- [10] Xuangeng Chu, Yu Li, Ailing Zeng, Tianyu Yang, Lijian Lin, Yunfei Liu, and Tatsuya Harada. Gpavatar: Generalizable and precise head avatar from image (s). *ICLR*, 2024. 2, 3
- [11] 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 0–0, 2019. 2, 3
- [12] Zhigang Deng and Junyong Noh. Computer facial animation: A survey. In *Data-driven 3D facial animation*, pages 1–28. Springer, 2008. 3
- [13] Helisa Dharmo, Yinyu Nie, Arthur Moreau, Jifei Song, Richard Shaw, Yiren Zhou, and Eduardo Pérez-Pellitero. Headgas: Real-time animatable head avatars via 3d gaussian splatting. *arXiv preprint arXiv:2312.02902*, 2023. 2, 3, 4
- [14] Robert A Drebin, Loren Carpenter, and Pat Hanrahan. Volume rendering. *ACM Siggraph Computer Graphics*, 22(4):65–74, 1988. 4
- [15] 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. 2
- [16] Yao Feng, Haiwen Feng, Michael J Black, and Timo Bolkart. Learning an animatable detailed 3d face model from in-the-wild images. *ACM Transactions on Graphics (ToG)*, 40(4):1–13, 2021. 2, 3
- [17] Yao Feng, Fan Wu, Xiaohu Shao, Yanfeng Wang, and Xi Zhou. Joint 3d face reconstruction and dense alignment with position map regression network. In *Proceedings of the European conference on computer vision (ECCV)*, pages 534–551, 2018. 3
- [18] James D Foley. *Computer graphics: principles and practice*, volume 12110. Addison-Wesley Professional, 1996. 3
- [19] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022. 4
- [20] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8649–8658, 2021. 2
- [21] Stathis Galanakis, Baris Gecer, Alexandros Lattas, and Stefanos Zafeiriou. 3dmm-rf: Convolutional radiance fields for 3d face modeling. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3536–3547, 2023. 2, 3, 6
- [22] Simon Giebenhain, Tobias Kirschstein, Martin Rünz, Lourdes Agapito, and Matthias Nießner. Npga: Neural parametric gaussian avatars. *arXiv preprint arXiv:2405.19331*, 2024. 2, 3
- [23] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d aware generator for high-resolution image synthesis. In *International Conference on Learning Representations*, 2022. 3
- [24] Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. Headnerf: A real-time nerf-based parametric head model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20374–20384, 2022. 2, 3, 6
- [25] Zixiong Huang, Qi Chen, Libo Sun, Yifan Yang, Naizhou Wang, Qi Wu, and Minghui Tan. G-nerf: Geometry-enhanced novel view synthesis from single-view images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10117–10126, 2024. 2

- [26] Berna Kabadayi, Wojciech Zielonka, Bharat Lal Bhatnagar, Gerard Pons-Moll, and Justus Thies. Gan-avatar: Controlable personalized gan-based human head avatar. In *2024 International Conference on 3D Vision (3DV)*, pages 882–892. IEEE, 2024. 2, 3
- [27] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*, 34:852–863, 2021. 3
- [28] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 6
- [29] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 3
- [30] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023. 2, 3, 4
- [31] Chenyi Kuang, Jeffrey O Kephart, and Qiang Ji. Au-aware dynamic 3d face reconstruction from videos with transformer. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 6237–6247, 2024. 2, 3
- [32] Yushi Lan, Feitong Tan, Di Qiu, Qiangeng Xu, Kyle Genova, Zeng Huang, Sean Fanello, Rohit Pandey, Thomas Funkhouser, Chen Change Loy, et al. Gaussian3diff: 3d gaussian diffusion for 3d full head synthesis and editing. *arXiv preprint arXiv:2312.03763*, 2023. 2, 3
- [33] 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. 6
- [34] 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. 2, 3, 4
- [35] Xinyang Li, Jiaxin Wang, Yixin Xuan, Gongxin Yao, and Yu Pan. Ggavatar: Geometric adjustment of gaussian head avatar. *arXiv preprint arXiv:2405.11993*, 2024. 2, 3
- [36] Jiangke Lin, Yi Yuan, Tianjia Shao, and Kun Zhou. Towards high-fidelity 3d face reconstruction from in-the-wild images using graph convolutional networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5891–5900, 2020. 2, 3
- [37] Jiahao Luo, Jing Liu, and James Davis. Splatface: Gaussian splat face reconstruction leveraging an optimizable surface. *arXiv preprint arXiv:2403.18784*, 2024. 2, 3
- [38] Shengjie Ma, Yanlin Weng, Tianjia Shao, and Kun Zhou. 3d gaussian blendshapes for head avatar animation. *arXiv preprint arXiv:2404.19398*, 2024. 2, 4
- [39] Zhiyuan Ma, Xiangyu Zhu, Guo-Jun Qi, Zhen Lei, and Lei Zhang. Otavatar: One-shot talking face avatar with controllable tri-plane rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16901–16910, 2023. 2, 3
- [40] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4460–4470, 2019. 2
- [41] 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 *ECCV*, 2020. 2, 3, 4
- [42] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 165–174, 2019. 2
- [43] 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, 3
- [44] Shenhan Qian, Tobias Kirschstein, Liam Schoneveld, Davide Davoli, Simon Giebenhain, and Matthias Nießner. Gaussianavatars: Photorealistic head avatars with rigged 3d gaussians. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20299–20309, 2024. 2, 3, 7, 8
- [45] Romdhani and Vetter. Efficient, robust and accurate fitting of a 3d morphable model. In *Proceedings Ninth IEEE International Conference on Computer Vision*, pages 59–66. IEEE, 2003. 2
- [46] Shunsuke Saito, Gabriel Schwartz, Tomas Simon, Junxuan Li, and Giljoo Nam. Relightable gaussian codec avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 130–141, 2024. 2, 3
- [47] Zhijing Shao, Zhaolong Wang, Zhuang Li, Duotun Wang, Xiangru Lin, Yu Zhang, Mingming Fan, and Zeyu Wang. Splattingavatar: Realistic real-time human avatars with mesh-embedded gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1606–1616, 2024. 2, 3
- [48] Jingxiang Sun, Xuan Wang, Lizhen Wang, Xiaoyu Li, Yong Zhang, Hongwen Zhang, and Yebin Liu. Next3d: Generative neural texture rasterization for 3d-aware head avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20991–21002, 2023. 2, 3
- [49] Kartik Teotia, Xingang Pan, Hyeonwoo Kim, Pablo Garrido, Mohamed Elgharib, and Christian Theobalt. Hq3davatar: High quality implicit 3d head avatar. *ACM Transactions on Graphics*, 2024. 2
- [50] Thomas Vetter and Volker Blanz. Estimating coloured 3d face models from single images: An example based approach. In *Computer Vision—ECCV’98: 5th European Conference on Computer Vision Freiburg, Germany, June 2–6, 1998 Proceedings, Volume II 5*, pages 499–513. Springer, 1998. 3

- [51] Jie Wang, Xianyan Li, Jiucheng Xie, Feng Xu, and Hao Gao. Gaussianhead: Impressive 3d gaussian-based head avatars with dynamic hybrid neural field. *arXiv e-prints*, pages arXiv-2312, 2023. 2, 3
- [52] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 6
- [53] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multi-scale structural similarity for image quality assessment. In *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, volume 2, pages 1398–1402. Ieee, 2003. 5
- [54] Jun Xiang, Xuan Gao, Yudong Guo, and Juyong Zhang. Flashavatar: High-fidelity digital avatar rendering at 300fps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 2, 3, 7, 8
- [55] Hongyi Xu, Guoxian Song, Zihang Jiang, Jianfeng Zhang, Yichun Shi, Jing Liu, Wanchun Ma, Jiashi Feng, and Linjie Luo. Omniavatar: Geometry-guided controllable 3d head synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12814–12824, 2023. 2, 3
- [56] Yuelang Xu, Benwang Chen, Zhe Li, Hongwen Zhang, Lizhen Wang, Zerong Zheng, and Yebin Liu. Gaussian head avatar: Ultra high-fidelity head avatar via dynamic gaussians. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1931–1941, 2024. 2, 3
- [57] Peizhi Yan, James Gregson, Qiang Tang, Rabab Ward, Zhan Xu, and Shan Du. Neo-3df: Novel editing-oriented 3d face creation and reconstruction. In *Proceedings of the Asian Conference on Computer Vision*, pages 486–502, 2022. 2
- [58] Wang Yifan, Felice Serena, Shihao Wu, Cengiz Öztireli, and Olga Sorkine-Hornung. Differentiable surface splatting for point-based geometry processing. *ACM Transactions on Graphics (TOG)*, 38(6):1–14, 2019. 4
- [59] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 5, 6
- [60] Yufeng Zheng, Wang Yifan, Gordon Wetzstein, Michael J Black, and Otmar Hilliges. Pointavatar: Deformable point-based head avatars from videos. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 21057–21067, 2023. 6
- [61] Peng Zhou, Lingxi Xie, Bingbing Ni, and Qi Tian. Cips-3d++: End-to-end real-time high-resolution 3d-aware gans for gan inversion and stylization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 2
- [62] Wenbin Zhu, HsiangTao Wu, Zeyu Chen, Noranart Vespapunt, and Baoyuan Wang. Reda: reinforced differentiable attribute for 3d face reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4958–4967, 2020. 3
- [63] Yiyu Zhuang, Hao Zhu, Xusen Sun, and Xun Cao. Mofan-erf: Morphable facial neural radiance field. In *Computer*

Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part III, pages 268–285. Springer, 2022. 2, 3, 6