

# NPGA: Neural Parametric Gaussian Avatars

SIMON GIEBENHAIN, Technical University of Munich, Germany  
 TOBIAS KIRSCHSTEIN, Technical University of Munich, Germany  
 MARTIN RÜNZ, Synthesia, Germany  
 LOURDES AGAPITO, University College London, United Kingdom  
 MATTHIAS NIESSNER, Technical University of Munich, Germany



Fig. 1. **NPGA**: We utilize the rich expression space of Neural Parametric Head Models to create high-fidelity avatars with fine-grained expression control. Our avatars consist of a dynamics module and a canonical Gaussian point cloud, which is augmented with per-primitive features that encode valuable semantic information, as indicated on the left. On the right, we demonstrate a highly detailed cross-reenactment using the inset image as a driving expression.

The creation of high-fidelity, digital versions of human heads is an important stepping stone in the process of further integrating virtual components into our everyday lives. Constructing such avatars is a challenging research problem, due to a high demand for photo-realism and real-time rendering performance. In this work, we propose Neural Parametric Gaussian Avatars (NPGA), a data-driven approach to create high-fidelity, controllable avatars from multi-view video recordings. We build our method around 3D Gaussian splatting for its highly efficient rendering and to inherit the topological flexibility of point clouds. In contrast to previous work, we condition our avatars' dynamics on the rich expression space of neural parametric head models (NPHM), instead of mesh-based 3DMMs. To this end, we distill the backward deformation field of our underlying NPHM into forward deformations which are compatible with rasterization-based rendering. All remaining fine-scale, expression-dependent details are learned from the multi-view videos. For increased representational capacity of our avatars, we propose per-Gaussian latent features that condition each primitives dynamic behavior. To regularize this increased dynamic expressivity, we propose Laplacian terms on the latent features and predicted dynamics. We evaluate our method on the public NeRSemble dataset, demonstrating that NPGA significantly outperforms the previous state-of-the-art avatars on the self-reenactment task by  $\approx 2.6$  PSNR. Furthermore, we demonstrate accurate animation capabilities from real-world monocular videos.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

SA Conference Papers '24, December 3–6, 2024, Tokyo, Japan

© 2024 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-1131-2/24/12.

<https://doi.org/10.1145/3680528.3687689>

CCS Concepts: • **Computing methodologies** → **Tracking**; **Motion capture**.

Additional Key Words and Phrases: Virtual avatars, 3D Gaussian splatting, Data-driven animation, 3d morphable models

## ACM Reference Format:

Simon Giebenhain, Tobias Kirschstein, Martin Rünz, Lourdes Agapito, and Matthias Nießner. 2024. NPGA: Neural Parametric Gaussian Avatars. In *SIGGRAPH Asia 2024 Conference Papers (SA Conference Papers '24)*, December 3–6, 2024, Tokyo, Japan. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3680528.3687689>

## 1 INTRODUCTION

Creating photo-realistic 3D avatars is one of the core challenges in computer graphics and includes a wide range of applications such as movies, games, AR/VR teleconferencing, and the metaverse. In particular, there is a strong motivation to reconstruct digital avatars from real-world captures, such as multi-view recordings, to obtain a digital copy of a specific real person. The resulting digital avatars can then be animated and rendered from arbitrary viewpoints while expecting high visual fidelity; e.g., with respect to photo-realistic colors and details, preservation of identity, and the adoption of person-specific mannerisms. At the same time, many avatar applications demand real-time rendering capabilities without compromising visual quality.

Project Website: <https://simongiebenhain.github.io/NPGA/>

Recent advances at the intersection of computer graphics and vision research have steadily improved methods to digitally reconstruct 3D objects with photo-realistic rendering quality [Kerbl et al. 2023a; Mildenhall et al. 2021; Müller et al. 2022]. In particular, 3D Gaussian Splatting (3DGS) [Kerbl et al. 2023a] has been quickly adopted in recent work on digital humans, e.g., [Li et al. 2024; Zielonka et al. 2023] and virtual head avatars, e.g., [Qian et al. 2023; Xu et al. 2024], due to its efficient rendering and photo-realistic reconstructions. At the same time recent publicly available multi-view datasets [Işık et al. 2023; Kirschstein et al. 2023b; Pan et al. 2024; Wu et al. 2022], offer an ever more exciting basis for avatar research. A central question is how controllability can be achieved. The most prominent approach for heads is to utilize a 3D morphable model (3DMM), which offers compact descriptions of faces using disentangled parametric spaces for identity and facial expressions. When utilizing the expressions of an underlying 3DMM, e.g., [Gafni et al. 2021; Grassal et al. 2022; Qian et al. 2023; Xu et al. 2024; Zielonka et al. 2022], the avatar is optimized to follow a generalized expression space which enables expression transfer or animation through tracking in monocular videos [Thies et al. 2016]. While 3DMMs offer a compact parameterization, their linear nature inherently limits the fidelity of represented expressions. At the same time, we argue that the underlying expression space plays a crucial role in determining the quality of the created avatars. It not only influences the controllability of the resulting avatars but it also limits the sharpness of details. If the stream of input expression codes is insufficiently correlated with the observed images, the optimization problem can become fundamentally ill-posed and lead to overfitting.

To this end, we propose NPGA, a new avatar representation that leverages a learned deformation representation while ensuring that the predicted facial dynamics stay close to the prior of an underlying neural parametric head model (NPHM) [Giebenhain et al. 2023, 2024]. NPHM provides our avatars with more fine-grained expression control compared to classical, public 3DMMs [Li et al. 2017; Paysan et al. 2009], which were previously used for avatar creation. In another line of work Cao et al. [2022] explore a learned neural prior over faces, that helps to map baked uv-space information from a monocular video into a volumetric avatar representation.

NPGA consists of a canonical Gaussian point cloud, that can be forward-deformed using an expression code and rendered using 3DGS, similar to previous work [Qian et al. 2023; Xu et al. 2024; Zielonka et al. 2023]. As a rasterization-based approach, 3DGS cannot be efficiently combined with the *backward* deformation field of MonoNPHM. Therefore, we propose a distillation strategy to invert the deformation direction of MonoNPHM’s expression prior using a cycle consistency loss, similar to SCANimate [Saito et al. 2021]. The resulting *forward* deformation field becomes compatible with the rasterization-based rendering of 3DGS. To increase the overall dynamic expressivity of our avatars, we further propose to augment our canonical Gaussian with per-Gaussian latent features similar to GaussianHeadAvatar [Xu et al. 2024] and HeadGAS [Dhamo et al. 2023]. Compared to these works, we allow these features to influence the movement of Gaussians instead of only influencing dynamic changes in appearance. This enables our deformation module to operate in a higher dimensional space, that can describe

facial dynamics more effectively. We show that this added expressivity results in higher-fidelity image synthesis, but is required to be appropriately regularized to achieve artifact-free renderings. To this end, we formulate Laplacian smoothness terms on the latent features and predicted dynamics, based on the k-nearest neighbor graph in canonical space. Furthermore, we modify the adaptive density control (ADC) strategy of 3DGS for more detailed avatar reconstructions.

To summarize, our contributions are the following:

- To leverage MonoNPHM’s rich expression prior for fine-grained animation control and effective optimization, we utilize a cycle-consistency-based distillation strategy.
- We condition our deformation network on per-Gaussian latent features for increased dynamic capacity, use Laplacian regularization, and adjust the ADC for a dynamic setting.
- We outperform the previous state-of-the-art by 2.6 PSNR and 0.021 SSIM. Furthermore, we demonstrate accurate avatar re-animation from monocular RGB sequences.

## 2 RELATED WORK

### 2.1 Dynamic Scene Representations

Since humans are inherently dynamic, and subject to topological variation when for example opening and closing the mouth, general dynamic scene representations aim to solve a set of similar problems with the exemption of controllability. Implicit approaches based on Neural Radiance Fields [Mildenhall et al. 2021] provide such flexibility and can be extended to deformable scenes. This extension is either achieved by modeling temporal changes via a deformation field that accompanies a canonical frame [Park et al. 2021a,b], by adding time as a conditioning variable [Li et al. 2022], or directly decomposing the 4D scene volume into a computationally manageable representation [Attal et al. 2023; Song et al. 2023]. The extensive use of MLPs in neural radiance fields entails a high computational burden, however, and still proves costly after improvements such as hash encodings [Müller et al. 2022], triplanes [Chan et al. 2022] or other low-rank approximations [Shao et al. 2023]. MVP [Lombardi et al. 2021] uses a CNN for amortized decoding of small primitives that can be compared to the Gaussians in 3DGS [Kerbl et al. 2023a], but typically at lower sharpness. One approach to extend 3DGS to dynamic scenes is to optimize parameters like positions over time [Luiten et al. 2024], resulting in a representation that is hard to control. This limitation can be overcome by animating a canonical 3DGS representation with a deformation field [Yang et al. 2023], similar to the Nerfies [Park et al. 2021a] approach. NPGA adopts this paradigm too.

### 2.2 3D Morphable Models

Traditional morphable models [Blaž and Vetter 1999; Paysan et al. 2009] learn a representation of body geometry via PCA. They are one of the primary tools to drive human animation and are heavily used in industry applications, making them a core building block for work on virtual avatars. While some models are dedicated to specific regions like the face [Paysan et al. 2009] and head [Li et al. 2017], some variants include the neck [Zhang et al. 2023] or even the entire body [Pavlakos et al. 2019]. More recently, neural equivalents of

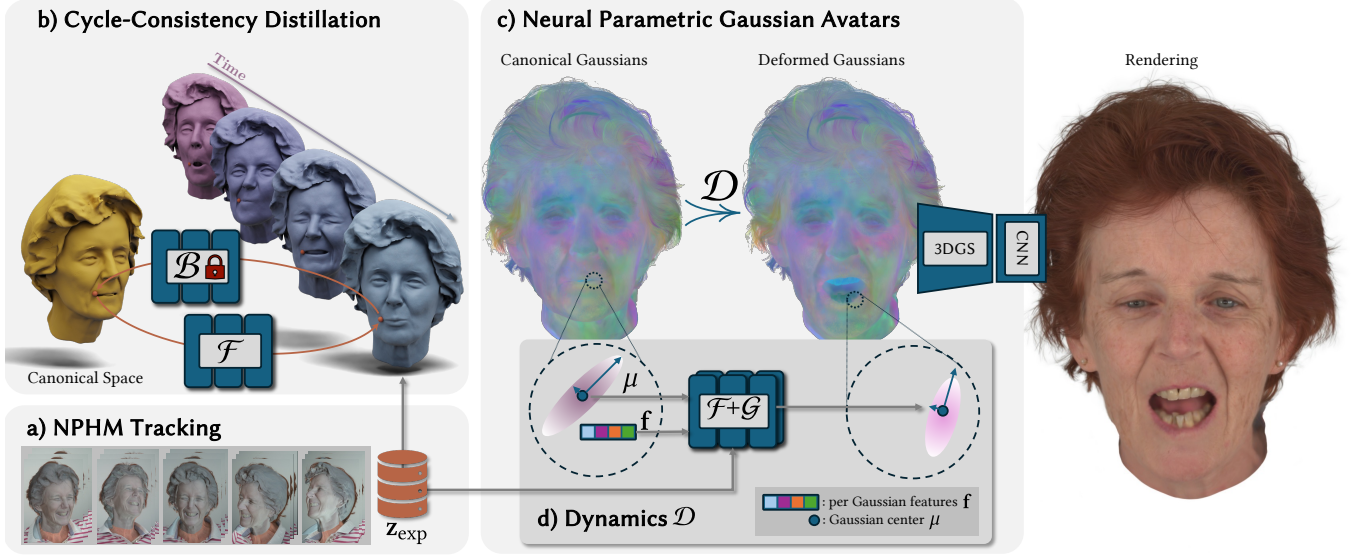


Fig. 2. **Method Overview:** The basis of our avatar optimization are multi-view video recordings alongside a MonoNPHM tracking thereof, see (a). Next, we extract a forward-deformation prior  $\mathcal{F}$  from MonoNPHM’s backward deformation field  $\mathcal{B}$  using a cycle-consistency loss, see (b). Our avatars consist of a canonical Gaussian point cloud (c), which is warped into posed space using our dynamics module  $\mathcal{D}$ , consisting of the coarse pre-trained component  $\mathcal{F}$  and a detail network  $\mathcal{G}$ . We condition both networks on per Gaussian features, which dictate each primitive’s behavior. After rendering the avatar with 3DGS, we employ a screen-space CNN to suppress small-scale artifacts.

3DMMs, such as i3DMM [Yenamandra et al. 2021], ImFace [Zheng et al. 2022a] or NPHM [Giebenhain et al. 2023, 2024], have improved upon the expression fidelity compared to classical PCA-based models. They encode a canonical representation of the geometry via a signed distance functions (SDF), which can be mapped to arbitrary expressions via a deformation field. NPGA utilizes NPHM as it offers several beneficial characteristics: It models the face densely including eyes, hair, and teeth, it captures local details well and it disentangles shape from expressions.

Instead of replacing mesh-based 3DMMs altogether, several approaches have been proposed to model details on top of an underlying PCA-driven mesh [Cao et al. 2022; Wang et al. 2022; Yang et al. 2020]. Specifically, Cao et al. [2022] leverage high-frequency information from tracked UV texture maps to create and animate avatars, using a high-detailed, learned expression space.

### 2.3 Human Head Reconstruction and Animation

Existing approaches for animating avatars mainly differ in two fundamental aspects: first, the utilized 3D representation in combination with its rendering mechanism, and second, how the expression codes are transferred into scene dynamics. Existing work has explored, among others, meshes in combination with deferred neural rendering [Grassal et al. 2022; Kim et al. 2018], neural radiance fields [Gafni et al. 2021; Zhao et al. 2023; Zielonka et al. 2022], and CNNs in combination with primitive-based volume rendering [Cao et al. 2022; Lombardi et al. 2019, 2021]. Recently, there has been a lot of interest in point-based representations and rendering [Dhamo et al. 2023; Qian et al. 2023; Xu et al. 2024; Zheng et al. 2023], especially since Kerbl et al. [2023a] proposed 3D Gaussian Splatting. While

some approaches choose to explain the face’s movement explicitly through the underlying mesh of a 3DMM, e.g., [Athar et al. 2022; Qian et al. 2023; Zielonka et al. 2022], others choose the opposing extreme of a more data-driven approach by freely learning the face movement using a neural component, which is directly conditioned on the expression codes [Gafni et al. 2021; Lombardi et al. 2021; Xu et al. 2024]. For NPGA we adopt the latter idea but replace the 3DMM with a neural parametric model, anticipating that the improved tracking quality translates to higher fidelity avatar reconstructions and animations.

## 3 PRELIMINARIES

### 3.1 3D Gaussian Splatting (3DGS)

3DGS uses a point-based scene representation, where each point represents a Gaussian primitive that is described by a position  $\mu$ , rotation  $\mathbf{q}$ , scale  $\mathbf{S}$ , opacity  $\alpha$  and spherical harmonics coefficients  $\mathbf{SH}$ . In the following, we let the following notation

$$\mathcal{A} = \{\mu, \mathbf{q}, \mathbf{S}, \alpha, \mathbf{SH}\}, \quad I = \text{3DGS}(\mathcal{A}, \pi_{K,E}) \quad (1)$$

denote the set of attributes  $\mathcal{A}$  composing the Gaussian point cloud, and its tile-based differentiable rasterization into an image  $I$  under the camera projection described by intrinsic and extrinsic parameters  $K$  and  $E$  respectively.

### 3.2 Neural Parametric Head Models

3DMMs describe the geometry (and appearance) of faces (or heads) using disentangled parametric spaces for identity and expression variations. NPHM is a special case of a 3DMM, which represents a person’s head geometry using a neural SDF and deformation field

conditioned on identity latent codes  $\mathbf{z}_{\text{id}}$  and expression codes  $\mathbf{z}_{\text{exp}}$ , respectively. In particular, our work builds on top of MonoNPHM formulation, which describes expressions using a neural backward deformation field

$$\mathbf{x}_c = \mathcal{B}(\mathbf{x}_p; \mathbf{z}_{\text{id}}, \mathbf{z}_{\text{exp}}) \quad (2)$$

that warps points  $\mathbf{x}_p$  in posed space into canonical space  $\mathbf{x}_c$ .

## 4 METHOD

Fig. 2 shows an overview of our proposed representation and methodology to build our Neural Parametric Gaussian Avatars (NPGA). As described in Section 4.1, our avatars are composed of two key components: a canonical Gaussian point cloud  $\mathcal{A}_c$  and a dynamics module  $\mathcal{D}$  which deforms the Gaussians when provided with an expression code, similar to recent work [Dhamo et al. 2023; Qian et al. 2023; Xu et al. 2024]. In Section 4.2 we describe our distillation strategy that allows NPGA to leverage the rich latent expression space and detailed motion prior of MonoNPHM [Giebenhain et al. 2024]. Given multi-view video recordings alongside tracked MonoNPHM expression codes, we jointly optimize for our canonical Gaussians and dynamics module, as described in Section 4.3.

### 4.1 Neural Parametric Gaussian Avatars

**4.1.1 Canonical Representation.** Compared to the default scene representation of 3DGS, outlined in Eq. (1), we augment our canonical Gaussian point cloud

$$\mathcal{A}_c = \{\mu, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{SH}, \} \cup \{\mathbf{f}\} \quad (3)$$

with per-Gaussian features  $\mathbf{f} \in \mathbb{R}^{N \times 8}$ . While these features are static themselves, they provide crucial semantic information to describe the dynamic behavior of the respective primitives. In some sense, our per-Gaussian features serve a similar purpose as positional encodings [Mildenhall et al. 2021; Müller et al. 2022], which are uncorrelated with spatial coordinates, have infinite spatial resolution and do not require additional data structures. The idea of per-Gaussian features has been previously proposed in [Dhamo et al. 2023; Xu et al. 2024]. Crucially, however, both works do not utilize these features to better predict the movement of Gaussians, which we ablate in Section 5.5. [Dhamo et al. 2023] goes beyond static per-primitive features by dynamically blending them with expression codes. In the formulation of [Xu et al. 2024] and ours, this happens implicitly inside the MLP introduced in the next paragraph.

**4.1.2 Dynamics Module.** We model facial expressions using a dynamics model  $\mathcal{D}$  which is decomposed into two Multi-Layer Perceptrons (MLPs), a coarse prior-based network  $\mathcal{F}$  and a network  $\mathcal{G}$  responsible for modeling all remaining details. Our prior-guided forward deformation field

$$\delta_\mu^\mathcal{F} = \mathcal{F}(\mu, \mathbf{f}; \mathbf{z}_{\text{exp}}) \in \mathbb{R}^3 \quad (4)$$

is optimized to act as the inverse of MonoNPHM’s backward deformations  $\mathcal{B}$ , as we describe later in Section 4.2.  $\mathcal{F}$  is a coordinate-based network which predicts offsets  $\delta_\mu^\mathcal{F}$  to the Gaussian centers  $\mu$  and is conditioned on spatial coordinates  $\mu$ , features  $\mathbf{f}$  and expression code  $\mathbf{z}_{\text{exp}}$ . Note, that  $\mathcal{F}$  acts independently on each primitive, which we omit in Eq. (4) and below for clarity.

To represent dynamics beyond NPHM’s prior, such as fine-scaled expression-dependent wrinkles, and appearance changes, e.g. due to ambient occlusions and changes in blood flow concentration, we rely on a second MLP

$$\delta_a^\mathcal{G} = \mathcal{G}_a(\mu, \mathbf{f}; \mathbf{z}_{\text{exp}}) \quad (\forall a \in \mathcal{A}_c), \quad (5)$$

which predicts offsets for all canonical Gaussian attributes  $a$ . We use the same architecture for  $\mathcal{G}$  as for  $\mathcal{F}$ , besides having more output channels due to the increased number of attribute offsets. In total, given an expression code  $\mathbf{z}_{\text{exp}} \in \mathbb{R}^{100}$  we obtain the Gaussian point cloud in posed space  $\mathcal{A}_p$  by adding  $\delta^\mathcal{F}$  and  $\delta^\mathcal{G}$  to their respective canonical attributes, which we denote as

$$\mathcal{A}_p = \mathcal{D}(\mathcal{A}_c; \mathbf{z}_{\text{exp}}). \quad (6)$$

**4.1.3 Screen-Space Refinement.** Finally, after rendering the posed Gaussians  $\mathcal{A}_p$  using the differentiable rasterizer from Kerbl et al. [2023b], we apply a screen-space CNN network [Xu et al. 2023]:

$$[\hat{I}_{\text{rgb}}, \hat{I}_h] = 3\text{DGS}(\mathcal{A}_p; \pi_{K,E}), \quad \hat{I}_{\text{cnn}} = \text{CNN}([\hat{I}_{\text{rgb}}, \hat{I}_h]), \quad (7)$$

where the  $\hat{I}_{\text{rgb}}$  denotes rendered RGB colors.  $\hat{I}_h$  is a latent image used for the CNN refinement module, which is obtained by rendering  $\mathbf{h} + \delta_h^\mathcal{G}$ , where  $\mathbf{h}$  are latent CNN features which we additionally include in  $\mathcal{A}_c$ . We took inspiration from Xu et al. [2024], to include a screen-space CNN, but note that we decided against performing super-resolution.

### 4.2 Cycle-Consistency Distillation

One of our core ideas is to leverage MonoNPHM’s motion prior and expression space. However, since it utilizes a *backward* deformation field  $\mathcal{B}$ , which warps points into canonical space, we cannot directly incorporate this deformation prior in our pipeline. Instead, we need *forward* deformations, which warp points into posed space, such that they can be directly rasterized. While it is possible to numerically approximate the inverse of  $\mathcal{B}$  using iterative root-finding [Chen et al. 2023, 2021], we search for a more computationally efficient method. To this end we resort to a cycle consistency loss, similar to [Saito et al. 2021] and follow-ups [Dhamo et al. 2023; Yang et al. 2022] which optimize for volumetric skinning weights used for animation via linear blend skinning.

We propose to directly distill a forward deformation network  $\mathcal{F}$  as the inverse of  $\mathcal{B}$  (see Fig. 2, using a cycle consistency loss

$$\mathcal{L}_{\text{cyc}}(\mathbf{x}_c) = \|\mathcal{B}(\mathcal{F}(\mathbf{x}_c, \mathbf{f}(\mathbf{x}_c))) - \mathbf{x}_c\|_2^2. \quad (8)$$

Using Eq. (8) we can directly supervise  $\mathcal{F}$  with the knowledge of  $\mathcal{B}$  for arbitrarily sampled points  $\mathbf{x}_c \in \mathbb{R}^3$  in canonical space and expression codes  $\mathbf{z}_{\text{exp}}$ . Note, that in Eq. (8) we omit the dependence on expression codes  $\mathbf{z}_{\text{exp}}$  for clarity. During this stage, the canonical space is not yet discretized into a set of Gaussian primitives. Hence, we utilize a feature *field*

$$\mathbf{f}(\mathbf{x}_c) = \text{TRIPLANE}(\mathbf{x}_c) \quad (9)$$

represented as low-resolution 64x64 triplanes [Chan et al. 2022; Peng et al. 2020], which can be evaluated at arbitrary points  $\mathbf{x}_c$  in canonical space. For convenience, we train  $\mathcal{F}$  once, using all six identities from the NeRSemble dataset [Kirschstein et al. 2023b] that we perform experiments with.



During training, each person has their own **TriPLANE**, which we regularize using a total variation loss. Furthermore, we regularize the norm of predicted offsets  $\|\mathcal{F}(x_c, \mathbf{f}(x_c))\|_2$  to be small. Note that distilling  $\mathcal{F}$  separately results in insignificant performance changes.

### 4.3 Avatar Optimization Strategy

After obtaining a forward deformation field  $\mathcal{F}$  using our cycle-consistency distillation strategy, we aim to jointly optimize for the canonical parameters  $\mathcal{A}_c$  and MLP  $\mathcal{G}$  to minimize a photometric energy term. To initialize the canonical Gaussians centers  $\mu$  we sample 50,000 points uniformly on the iso-surface of the tracked MonoNPHM model. The per Gaussian features are initialized by querying **TriPLANE**( $\mu$ ) at the sampled Gaussian centers. All remaining attributes are initialized using the default 3DGS procedure. In practice, we observed that keeping  $\mathcal{F}$  frozen results in sub-optimal performance, which is likely caused through topological issues during distillation in the mouth region. Hence, we decide to further optimize  $\mathcal{F}$  alongside  $\mathcal{G}$  and  $\mathcal{A}_c$ , however, using a significantly smaller learning rate and a warm-up schedule that encourages the preservation of the distilled prior.

Our optimization strives to minimize the photometric data term

$$\mathcal{L} = \|I - \hat{I}_{\text{rgb}}\|_1 + \lambda \left(1 - \text{SSIM}(I, \hat{I}_{\text{rgb}})\right) + \lambda \left(1 - \text{SSIM}(I, \hat{I}_{\text{cnn}})\right), \quad (10)$$

where  $I$  denotes a randomly sample ground truth image from the NeRSemble dataset with corresponding expression codes  $\mathbf{z}_{\text{exp}}$ .

**4.3.1 Regularization.** We find that regularizing both our canonical representation  $\mathcal{A}_c$ , as well as our dynamics module  $\mathcal{D}$  is crucial to avoid overfitting to the training expressions. To regularize NPGA we utilize a Laplacian smoothness term based on the  $k$ NN graph of canonical Gaussian centers. To this end let

$$\mathcal{R}_{\text{lap}}(x) = \left\| \frac{1}{|\mathcal{N}_i|} \left( \sum_{j \in \mathcal{N}_i} x_j \right) - x_i \right\|_2^2 \quad \left( x \in \{\mathbf{f}\} \cup \bigcup_{a \in \mathcal{A}_c} \delta_a^{\mathcal{G}} \right), \quad (11)$$

denote a Laplacian smoothness term, which we use to regularize the per Gaussian features  $\mathbf{f}$ , as well as, the offset predictions  $\delta_a^{\mathcal{G}}$  for all attributes  $a \in \mathcal{A}_c$ . Note, that whenever the number of canonical Gaussians changes due to densification or pruning, we recompute the  $k$ NN graph.

In addition to these smoothness terms, we encourage  $\mathcal{F}$  and  $\mathcal{G}$  to predict small offsets

$$\mathcal{R}_{\delta} = \lambda_{\mu}^{\mathcal{F}} \|\delta_{\mu}^{\mathcal{F}}\|_2^2 + \sum_{a \in \mathcal{A}} \lambda_a^{\mathcal{G}} \|\delta_a^{\mathcal{G}} - \mathbf{e}_a\|_2^2, \quad (12)$$

where  $\mathbf{e}_a$  denotes the neutral element for the group operation acting on attribute  $a$ . Similarly, we impose regularization on the per Gaussian attributes  $\mathcal{R}_{\mathbf{f}} = \|\mathbf{f}\|_2^2$  to remain small, and utilize the scale regularization loss of [Saito et al. 2024], which punishes scales  $S$  lying outside of a well-behaved range.

**4.3.2 Adaptive Density Control (ADC).** A central ingredient to the success of 3DGS is its strategy to adaptively add and prune Gaussians in areas where they are needed or redundant, based on a set of simple yet effective heuristics that are periodically invoked. The rules of ADC have been designed with static scenes in mind and we find the default settings to be suboptimal for our avatar creation. In the

dynamic scenario, there can be areas that remain hidden for large parts of the training sequence, such as the mouth interior. Therefore, we adjust the ADC by employing a generalized mean

$$M_i^e = \left( \frac{1}{N} \sum_{t \in T} \tau_t^e \right)^{1/e} \quad (13)$$

to aggregate the view-space gradients  $\tau_t$  of the  $i$ -th primitive of all frames  $T$  between invocations of the ADC mechanism. Note, that  $e = 1$  results in the default 3DGS settings. By increasing the exponent  $e$  the aggregation  $M_i^e$  becomes closer to a maximum function. Therefore, a few visible frames can be sufficient for the ADC to trigger densification, which is especially important for regions like the teeth and mouth interior. We find that  $e = 2$  already results in an increased number of Gaussians, leading to more detailed reconstruction, especially in the mouth interior.

Furthermore, we replace the hard opacity reset mechanism of 3DGS, which we find to be harmful to our optimization, with a softer variant proposed in [Bulò et al. 2024]. Instead of infrequently setting the  $\alpha$  values to be almost transparent, the opacities get reduced frequently for a small amount only, i.e. in our experiments by 0.01.

### 4.4 Differences to GaussianHeadAvatar (GHA)

In general our avatar formulation shares large parts with that of GHA [Xu et al. 2024]. A major difference is our use of MonoNPHM as underlying 3DMM, which we enabled through our distillation strategy. Furthermore, we utilize the per-Gaussian features to also influence the prediction of  $\delta_{\mu}$ , we add a Laplacian regularization, do not rely on screen-space super resolution or perceptual losses, and modify the ADC strategy instead of resorting to a fixed number of Gaussians as Xu et al. [2024].

## 5 RESULTS

For our experiments, we use the state-of-the-art, public multi-view video NeRSemble dataset [Kirschstein et al. 2023b], from which we choose a diverse set of six subjects performing challenging facial expressions. After providing additional details on the performed experiments and baseline methods in Sections 5.1 and 5.2, we present our main results on the tasks of self- and cross-reenactment in Sections 5.3 and 5.4, respectively. Finally, in Section 5.5 we ablate a series of experiments validating our proposed model components. We highly encourage the reader to consult our supplemental video for complete results including temporal information.

### 5.1 Experimental Setup

The NeRSemble dataset provides 16 synchronized and calibrated videos, from which we choose 15 cameras for training and the frontal camera for evaluation. Furthermore, we train our avatars on all sequences, except for the "FREE"-sequence which we keep as a held-out evaluation sequence for the self-reenactment task.

#### 5.1.1 Baselines.

*GaussianAvatars (GA)* [Qian et al. 2023]: GaussianAvatars is a recent method that creates 3DGS-based avatars, by binding them



Fig. 3. **Self-Reenactment**: Qualitative comparison of different methods on the held-out sequence.

to the FLAME 3DMM model [Li et al. 2017]. Therefore, GaussianAvatars can be extremely efficiently animated, since there is no need to evaluate a costly neural component. On the downside, GaussianAvatars are limited to the facial movements lying inside the FLAME expression space.

*GaussianHeadAvatar (GHA)* [Xu et al. 2024]: GHA is another recent 3DGS-based avatar method, which also learns deformation fields from multi-view video and is controlled through their custom multi-view BFM [Paysan et al. 2009] tracking. In general, our approach is similar to GHA. Thus we list the major differences in Section 4.4. Furthermore, we add a version of GHA which is conditioned on our tracked MonoNPHM expression codes, indicated as  $GHA_{NPHM}$ .

*Mixture of Volumetric Primitives (MVP)* [Lombardi et al. 2021]: MVP utilizes a combination of volume rendering and a head-geometry aware CNN that creates a volumetric payload in an amortized fashion. In contrast to our method and the other baselines, MVP utilizes a Variational Auto-Encoder (VAE) [Kingma and Welling 2014] to learn a latent expression encoding based on their 3DMM tracking. Note, however, that we do not provide MVP with view-average textures to obtain a more comparable evaluation setting.

Table 1. **Quantitative Comparison:** We compare against our baselines on self-reenactment using a held-out sequence. For completeness we also report metrics on the held-out camera of the training sequences, denoted as novel-view synthesis (NVS).

Method	NVS			Self-Reenactment		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
MVP	33.42	0.957	0.083	27.19	0.919	0.114
GA	32.95	0.956	0.080	27.77	0.926	0.104
GHA	33.92	0.953	0.045	26.81	0.914	0.077
<b>Ours</b>	<b>36.84</b>	<b>0.971</b>	<b>0.034</b>	<b>30.26</b>	<b>0.934</b>	<b>0.055</b>
$GHA_{NPHM}$	33.09	0.952	0.049	26.60	0.911	0.078
$Ours_{BFM}$	34.97	0.962	0.052	28.82	0.924	0.076

**5.1.2 Metrics.** To evaluate the self-reenactment task we use the Peak Signal-to-Noise Ratio (PSNR), structural similarity index measure (SSIM), and perceptual LPIPS [Zhang et al. 2018] metrics. For the sake of completeness, we also report numbers for a dynamic novel view synthesis (NVS) scenario, where we compare all methods on the held-out camera view of the training sequences. We focus our evaluation on the facial region, since neck and torso are not accurately explained by NPHM and the underlying 3DMMs of our baselines. To this end, we leverage segmentation masks from Facer [Zheng et al. 2022b] to mask out the neck and torso before computing the metrics. Furthermore, we compute metrics at a resolution of 550x802. Since we train GHA on 1024x1024 we downsample and crop the generated images accordingly.

## 5.2 Implementation Details

*Hyper-Parameters.* For both our deformation networks  $\mathcal{F}$  and  $\mathcal{G}$  we use 6-layer MLPs with a hidden dimensionality of 256. In order

to preserve the prior that  $\mathcal{F}$  obtained in our distillation procedure, we set its learning rate to  $4e - 5$ , while  $\mathcal{G}$  is equipped with a much higher learning rate of  $2e - 3$ . Additionally, we freeze the network parameters of  $\mathcal{F}$  for the first 5,000 optimization steps. We decay both learning rates twice by a factor of 2 during the course of 800,000 optimization steps. Furthermore, we employ weight decay on  $\mathcal{F}$  and  $\mathcal{G}$ , using a weight of 0.1 as an additional regularization measure. We perform an ADC step every 5,000 iterations, and multiply the gradient threshold for the densification by a factor of 2, to accommodate for the fact that our loss combines the losses of the RGB rendering and CNN-refined predictions.

*Runtime.* While we do not focus on efficient training, animation, or rendering of avatars, we acknowledge the importance of fast animation and rendering. In our unoptimized implementation, we can render images at 31 frames per second (FPS) for 550x802 and 18 FPS at 1100x1604 on an NVIDIA RTX3080 graphics card, which includes deformation, rendering, and CNN. When omitting the CNN the speed increases to 43 and 38 FPS, respectively. As a comparison, GHA runs at 22 FPS at 1024x1024 on the same machine. Fig. 8 indicates the number of Gaussians during training, which we limit to a maximum of 250k. We train all our avatars, and baselines, until convergence, which roughly takes 7 hours for GA (on an RTX2080), 30 hours (on an RTX3090) for GHA and our method, and 60 hours (on an RTX2080) for MVP.

*Data Preparation.* We obtain MonoNPHM trackings on the NeRSem-ble dataset using a purely geometric constraint between the MonoNPHM’s predicted surface and a point cloud reconstructed using COLMAP [Schönerberger and Frahm 2016]. We utilize the same tracking algorithm that has been previously used in [Aneja et al. 2024; Kirschstein et al. 2023a]. For training and quantitative evaluation, we use a resolution of 550x802, the same as we use for MVP and GA. For our qualitative results, we fine-tune our avatars on 1100x1604 resolution for another 5 hours of training time. Furthermore, we mask out the torso, since it is neither contained in NPHM’s expression space nor the focus of our work.

## 5.3 Self-Reenactment

Our main evaluation is concerned with the self-reenactment task. For this purpose, all avatars are trained on a set of 21 training sequences alongside their respective tracking results. To evaluate the avatars, they are animated using the tracked expressions from a held-out test sequence. We present qualitative and quantitative results in Fig. 3 and Table 1, respectively, and recommend the reader to consider the supplemental video for temporal results. Our predicted self-reenactments portray the unseen expression more accurately and contain sharper details in relatively static areas like the hair region.

*Importance of the Underlying Tracking.* Interestingly,  $GHA_{NPHM}$  performs slightly worse than GHA, indicating that MonoNPHM expression codes alone do not immediately boost performance. Instead, we hypothesize that without NPHM’s motion prior as initialization, NPHM’s latent expression distribution might provide a more complicated training signal compared to the linear blendshapes of BFM. Additionally, we train our avatars, but utilize the BFM tracking from





Fig. 4. **Cross-Reenactment:** Qualitative comparison of transferring a driving expression from a different identity (left) to an avatar.

GHA, denoted as  $\text{Ours}_{BFM}$ . Indicating that our overall framework still slightly outperforms GHA when using the same tracked expression codes. In Fig. 7 we show that the learned dynamics can still result in accurate self-reenactments, despite flaws in the tracking as long as there is sufficient correlation the associated latent codes.

#### 5.4 Cross-Reenactment

Another crucial task is cross-reenactment, where driving expressions from another person are transferred to the avatar. Since a ground truth for cross-reenactment does not exist, we only report a qualitative comparison, which is presented in Fig. 4. We observe that all methods successfully disentangle identity and expression information, allowing for effective cross-reenactment. Our avatars, however, preserve the most details from the driving expressions. To demonstrate real-world applicability, Fig. 6 depicts cross-reenactment animations of our avatars using monocular RGB videos from a commodity camera under real-world circumstances. To this end, we utilize the monocular MonoNPHM tracker proposed by Giebenhain et al. [2024].

#### 5.5 Ablations Study

In order to verify several important components of NPGA, we perform ablation experiments using three subjects. Quantitative and qualitative results of our ablations can be found in Table 2 and Fig. 5, respectively. First, we run a "vanilla" version of NPGA that serves as a baseline. This version does not utilize per Gaussian features, the CNN, Laplacian smoothness terms, and uses  $e = 1$  in Eq. (13) for the ADC. This model fails to produce sharp renderings for fine-scale details, and areas that are complicated due to frequent occlusions and reflections, like the bottom teeth and eyes.

*Per-Gaussian Features.* When adding per-Gaussian features to the vanilla model, denoted as "+p.G.F.", the increased representational capacity results in sharper reconstructions. At the same time, we occasionally observe artifacts of "free-floating" primitives, as highlighted in the second column of Fig. 5. These artifacts can be largely removed using our proposed Laplacian smoothness terms. Table 2 indicates that using this regularization significantly shrinks the generalization gap between training (NVS) and testing (self-reenactment) sequences. Compared to this model, "Ours" also includes a CNN,



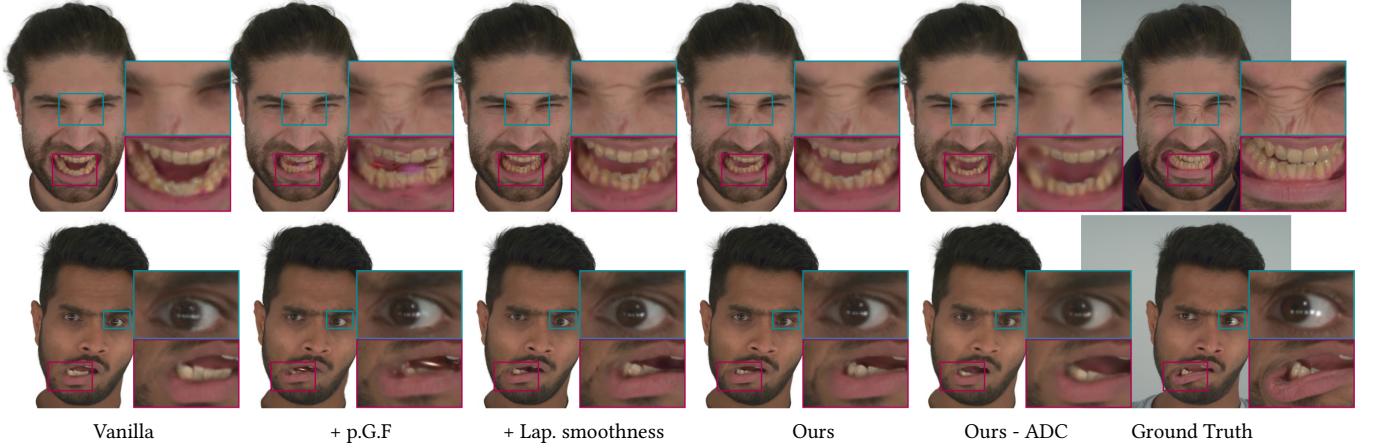


Fig. 5. **Ablation Study:** Without utilizing per Gaussian features ("Vanilla"), the avatars fail to represent fine expression details and complicated regions like the eyes and bottom teeth. Adding per Gaussian features (p.G.F.) results in significantly sharper reconstructions but is prone to artifacts under extreme expressions. Adding our Laplacian regularization ("Lap. smoothness") and a screen-space CNN ("Ours") finally resolves all artifacts. Furthermore, "Ours-ADC" demonstrates that the default densification strategy inhibits detailed reconstructions.

Table 2. **Ablations:** We perform our ablation experiments on a subset of three subjects. We report Novel-View Synthesis (NVS) for completeness.

Method	NVS			Self-Reenactment		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
Vanilla	35.80	0.965	0.048	30.16	0.927	0.067
+PGF	37.04	0.970	0.037	30.54	0.929	0.059
+Lap.smooth	36.85	0.969	0.038	30.56	0.928	0.059
<b>Ours(+CNN)</b>	<b>37.23</b>	<b>0.972</b>	<b>0.033</b>	<b>30.65</b>	<b>0.933</b>	<b>0.053</b>
Ours-ADC	36.12	0.967	0.045	30.49	0.933	0.070
Ours- $\delta_\mu$	36.73	0.968	0.040	30.44	0.931	0.059

which further boosts metrics and visual quality. Additionally, we ablate the importance of letting per-Gaussian features influence the predicting of position offsets, since recent work [Dhamo et al. 2023; Xu et al. 2024] does not allow such interactions. In Table 2 "Ours- $\delta_\mu$ " denotes a version of NPGA that disables such an influence, indicating the benefits of our formulation.

*Adaptive Density Control.* Finally, we show the importance of using an adjusted ADC strategy. Surprisingly, using the default ADC settings with a densification interval of 5000 steps and an opacity reset interval of 50.000 steps almost completely diminishes the improvements of our other contributions. While we do not claim that using  $e = 2$  in Eq. (13) is necessary for great performance, we simply note that finding a setting that lets enough Gaussians appear is important, especially for fine-scaled wrinkles and teeth. The progress of the number of Gaussian during the optimization is illustrated in Fig. 8.

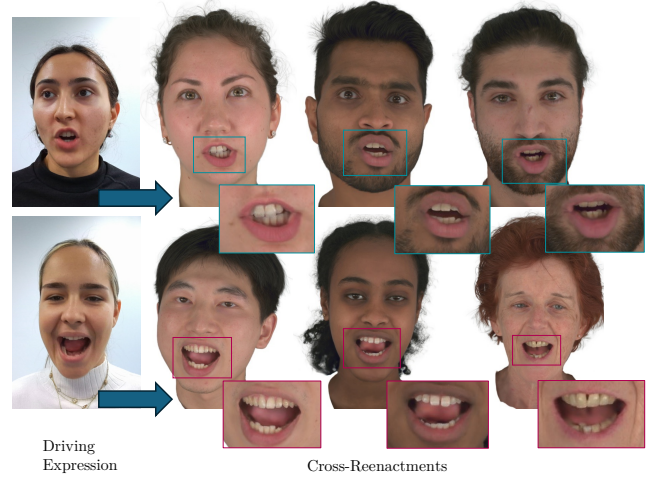


Fig. 6. **Real-World Application:** We utilize the monocular RGB tracking from MonoNPHM to animate our high-fidelity avatars, demonstrating the applicability of our avatars outside of multi-view capture studios.

## 6 LIMITATIONS AND FUTURE WORK

In our experiments, we show that NPGA can create controllable and high-fidelity virtual head avatars from multi-view video data. However, both controllability and reconstruction quality of our avatars are fundamentally restricted to what the underlying 3DMM can explain. Therefore, regions like the neck, torso, tongue, and eyeball rotation, which are not explained by NPHM's expression codes, cannot be animated as reliably or might even lead to artifacts due to overfitting. Possible solutions are extensions of the underlying 3DMM to provide a more complete description of a person's state, e.g. the inclusion of the neck [Zhang et al. 2023] or even torso and complete bodies [Pavlakos et al. 2019].

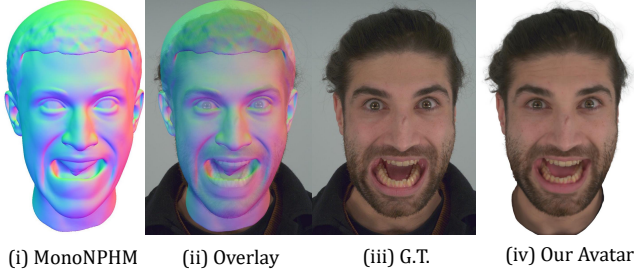


Fig. 7. **Tracking Failure:** While our MonoNPHM-based tracking occasionally fails for extreme facial expressions (see normals (i) and overlay (ii)), NPGA still produces good animations for held-out expressions (see (iv)) due to the smooth nature of expression latent space.

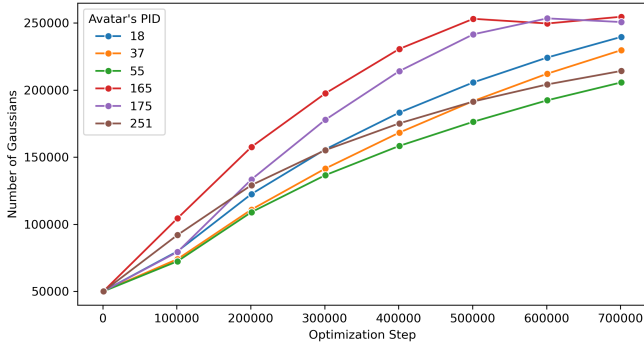


Fig. 8. **Number of Gaussians:** We illustrate the growing number of Gaussians for each of our avatars, where PID denotes the personal identifiers in the Nersemble dataset. We limit the maximum number to 250k.

Furthermore, as a data-driven approach to avatar creation, our method is limited, to some degree, to the available training data per person. We believe that recent large-scale multi-view video dataset of human heads [Kirschstein et al. 2023b; Pan et al. 2024] open up opportunities to learn a generalized head model, such as [Cao et al. 2022], with much higher fidelity than NPHM and other available 3DMMs, through the use of photometric optimization and efficient rendering, like 3DGS.

## 7 ETHICAL CONSIDERATIONS

The creation of photo-realistic avatars has potential for a wide range of malevolent activities, e.g. privacy violation, identity theft and spreading of deceptive content through deepfakes. We condemn any malicious or unauthorized uses of such technology, and emphasize the need for reliable detection methods that maintain the authenticity of media content and avert negative societal impact.

## 8 CONCLUSION

In this work, we have proposed Neural Parametric Gaussian Avatars (NPGA), a method for creating accurately controllable and high-fidelity virtual head avatars. The main focus of our work is the usage of MonoNPHM’s rich expression space and motion prior. To this end, we leverage a cycle-consistency strategy to distill a forward deformation field from MonoNPHM, such that it becomes compatible with 3D Gaussian Splatting. We proposed an effective

regularization strategy and adjust the adaptive density control strategy, both of which are important for ideal avatar quality. We show that utilizing per-Gaussian features to conditioning the complete dynamics module helps, compared to limiting the influence to appearance changes. In our experiments, we significantly outperform the previous state-of-the-art avatars on self-reenactment. Finally, we showed the applicability of our avatars beyond a controlled multi-view set-up by animating them from monocular RGB video trackings.

**Acknowledgements.** This work was funded by Synthesia and supported by the ERC Starting Grant Scan2CAD (804724), the German Research Foundation (DFG) Research Unit “Learning and Simulation in Visual Computing”. We would like to thank our research assistant Mohak Mansharamani, and Angela Dai for the video voice-over.

## REFERENCES

- Shivangi Aneja, Justus Thies, Angela Dai, and Matthias Nießner. 2024. FaceTalk: Audio-Driven Motion Diffusion for Neural Parametric Head Models. arXiv:2312.08459 [cs.CV]
- ShahRukh Athar, Zexiang Xu, Kalyan Sunkavalli, Eli Shechtman, and Zhixin Shu. 2022. RIGNeRF: Fully Controllable Neural 3D Portraits. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 20364–20373.
- Benjamin Attal, Jia-Bin Huang, Christian Richardt, Michael Zollhoefer, Johannes Kopf, Matthew O’Toole, and Changil Kim. 2023. HyperReel: High-Fidelity 6-DoF Video with Ray-Conditioned Sampling. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Volker Blanz and Thomas Vetter. 1999. A morphable model for the synthesis of 3D faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*. 187–194.
- Samuel Rota Buló, Lorenzo Porzi, and Peter Kotschieder. 2024. Revising Densification in Gaussian Splatting. arXiv:2404.06109 [cs.CV]
- Chen Cao, Tomas Simon, Jin Kyu Kim, Gabe Schwartz, Michael Zollhoefer, Shun-Suke Saito, Stephen Lombardi, Shih-En Wei, Danielle Belko, Shou-I Yu, Yaser Sheikh, and Jason Saragih. 2022. Authentic Volumetric Avatars from a Phone Scan. *ACM Trans. Graph.* 41, 4, Article 163 (jul 2022), 19 pages. <https://doi.org/10.1145/3528223.3530143>
- Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. 2022. Efficient Geometry-aware 3D Generative Adversarial Networks. In *CVPR*.
- Xu Chen, Tianjian Jiang, Jie Song, Max Rietmann, Andreas Geiger, Michael J. Black, and Otmar Hilliges. 2023. Fast-SNARF: A Fast Deformer for Articulated Neural Fields. *Pattern Analysis and Machine Intelligence (PAMI)* (2023).
- Xu Chen, Yufeng Zheng, Michael J Black, Otmar Hilliges, and Andreas Geiger. 2021. SNARF: Differentiable Forward Skinning for Animating Non-Rigid Neural Implicit Shapes. In *International Conference on Computer Vision (ICCV)*.
- Helisa Dhamo, Yinyu Nie, Arthur Moreau, Jifei Song, Richard Shaw, Yiren Zhou, and Eduardo Pérez-Pellitero. 2023. HeadGaS: Real-Time Animatable Head Avatars via 3D Gaussian Splatting. arXiv:2312.02902 [cs.CV] <https://arxiv.org/abs/2312.02902>
- Guy Gafni, Justus Thies, Michael Zollhoefer, and Matthias Nießner. 2021. Dynamic Neural Radiance Fields for Monocular 4D Facial Avatar Reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 8649–8658.
- Simon Giebenhain, Tobias Kirschstein, Markos Georgopoulos, Martin Rünz, Lourdes Agapito, and Matthias Nießner. 2023. Learning Neural Parametric Head Models. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*.
- Simon Giebenhain, Tobias Kirschstein, Markos Georgopoulos, Martin Rünz, Lourdes Agapito, and Matthias Nießner. 2024. MonoNPHM: Dynamic Head Reconstruction from Monocular Videos. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*.
- Philip-William Grassal, Malte Prinzler, Titus Leistner, Carsten Rother, Matthias Nießner, and Justus Thies. 2022. Neural head avatars from monocular RGB videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18653–18664.
- Mustafa İşık, Martin Rünz, Markos Georgopoulos, Taras Khakhulin, Jonathan Starck, Lourdes Agapito, and Matthias Nießner. 2023. HumanRF: High-fidelity neural radiance fields for humans in motion. arXiv preprint arXiv:2305.06356 (2023).
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023a. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics* 42, 4 (July 2023). <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023b. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on*

- Graphics 42, 4 (July 2023). <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>
- Hyeonwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Nießner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. 2018. Deep Video Portraits. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 163.
- Diederik P. Kingma and Max Welling. 2014. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14–16, 2014, Conference Track Proceedings*. arXiv:https://arxiv.org/abs/1312.6114v10 [stat.ML]
- Tobias Kirschstein, Simon Giebenhain, and Matthias Nießner. 2023a. Diffusion-Avatars: Deferred Diffusion for High-fidelity 3D Head Avatars. *arXiv preprint arXiv:2311.18635* (2023).
- Tobias Kirschstein, Shenhan Qian, Simon Giebenhain, Tim Walter, and Matthias Nießner. 2023b. NeRsemble: Multi-View Radiance Field Reconstruction of Human Heads. *ACM Trans. Graph.* 42, 4, Article 161 (jul 2023), 14 pages. <https://doi.org/10.1145/3592455>
- Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. 2017. Learning a model of facial shape and expression from 4D scans. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)* 36, 6 (2017), 194:1–194:17. <https://doi.org/10.1145/3130800.3130813>
- Tianye Li, Mira Slavcheva, Michael Zollhöfer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, and Zhao Yang Lv. 2022. Neural 3D Video Synthesis From Multi-View Video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 5521–5531.
- Zhe Li, Zerong Zheng, Lizhen Wang, and Yebin Liu. 2024. Animatable Gaussians: Learning Pose-dependent Gaussian Maps for High-fidelity Human Avatar Modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. 2019. Neural Volumes: Learning Dynamic Renderable Volumes from Images. *ACM Trans. Graph.* 38, 4, Article 65 (July 2019), 14 pages.
- Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. 2021. Mixture of Volumetric Primitives for Efficient Neural Rendering. *ACM Trans. Graph.* 40, 4, Article 59 (jul 2021), 13 pages. <https://doi.org/10.1145/3450626.3459863>
- Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. 2024. Dynamic 3D Gaussians: Tracking by Persistent Dynamic View Synthesis. In *3DV*.
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tanic, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM* 65, 1 (2021), 99–106.
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM Trans. Graph.* 41, 4, Article 102 (July 2022), 15 pages. <https://doi.org/10.1145/3528223.3530127>
- Dongwei Pan, Long Zhuo, Jingtian Piao, Huiwen Luo, Wei Cheng, Yuxin Wang, Siming Fan, Shengqi Liu, Lei Yang, Bo Dai, Ziwei Liu, Chen Change Loy, Chen Qian, Wayne Wu, Dahua Lin, and Kwan-Yee Lin. 2024. RenderMe-360: A Large Digital Asset Library and Benchmarks Towards High-fidelity Head Avatars. *Advances in Neural Information Processing Systems* 36 (2024).
- Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. 2021a. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5865–5874.
- Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. 2021b. HyperNeRF: A Higher-Dimensional Representation for Topologically Varying Neural Radiance Fields. *ACM Trans. Graph.* 40, 6, Article 238 (dec 2021).
- Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. 2019. Expressive Body Capture: 3D Hands, Face, and Body from a Single Image. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 10975–10985.
- Pascal Paysan, Reinhard Knecht, Brian Amberg, Sami Romdhani, and Thomas Vetter. 2009. A 3D face model for pose and illumination invariant face recognition. In *2009 sixth IEEE international conference on advanced video and signal based surveillance*. Ieee, 296–301.
- Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. 2020. Convolutional Occupancy Networks. In *European Conference on Computer Vision (ECCV)*.
- Shenhan Qian, Tobias Kirschstein, Liam Schoneveld, Davide Davoli, Simon Giebenhain, and Matthias Nießner. 2023. GaussianAvatars: Photorealistic Head Avatars with Rigid 3D Gaussians. *arXiv preprint arXiv:2312.02069* (2023).
- Shunsuke Saito, Gabriel Schwartz, Tomas Simon, Junxuan Li, and Giljoo Nam. 2024. Relightable Gaussian Codec Avatars. In *CVPR*.
- Shunsuke Saito, Jinlong Yang, Qianli Ma, and Michael J. Black. 2021. SCANimate: Weakly Supervised Learning of Skinned Clothed Avatar Networks. In *Proceedings IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*.
- Johannes Lutz Schönberger and Jan-Michael Frahm. 2016. Structure-from-Motion Revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Ruizhi Shao, Zerong Zheng, Hanzhang Tu, Boning Liu, Hongwen Zhang, and Yebin Liu. 2023. Tensor4D: Efficient Neural 4D Decomposition for High-fidelity Dynamic Reconstruction and Rendering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Liangchen Song, Anpei Chen, Zhong Li, Zhang Chen, Lele Chen, Junsong Yuan, Yi Xu, and Andreas Geiger. 2023. NeRFPlayer: A Streamable Dynamic Scene Representation with Decomposed Neural Radiance Fields. *IEEE Transactions on Visualization and Computer Graphics* 29, 5 (2023), 2732–2742. <https://doi.org/10.1109/TVCG.2023.3247082>
- J. Thies, M. Zollhöfer, M. Stamminger, C. Theobalt, and M. Nießner. 2016. Face2Face: Real-time Face Capture and Reenactment of RGB Videos. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE.
- Lizhen Wang, Zhiyua Chen, Tao Yu, Chenguang Ma, Liang Li, and Yebin Liu. 2022. FaceVerse: a Fine-grained and Detail-controllable 3D Face Morphable Model from a Hybrid Dataset. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR2022)*.
- Cheng-hsin Wu, Ningyuan Zheng, Scott Ardisson, Rohan Bali, Danielle Belko, Eric Brockmeyer, Lucas Evans, Timothy Godisart, Hyowon Ha, Xuhua Huang, Alexander Hypes, Taylor Koska, Steven Krenn, Stephen Lombardi, Xiaomin Luo, Kevyn McPhail, Laura Millerschoen, Michal Perdoch, Mark Pitts, Alexander Richard, Jason Saragih, Junko Saragih, Takaaki Shiratori, Tomas Simon, Matt Stewart, Autumn Trimble, Xinshuo Weng, David Whitewolf, Chenglei Wu, Shouo-I Yu, and Yaser Sheikh. 2022. Multiface: A Dataset for Neural Face Rendering. In *arXiv*. <https://doi.org/10.48550/ARXIV.2207.11243>
- Yuelang Xu, Benwang Chen, Zhe Li, Hongwen Zhang, Lizhen Wang, Zerong Zheng, and Yebin Liu. 2024. Gaussian Head Avatar: Ultra High-fidelity Head Avatar via Dynamic Gaussians. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Yuelang Xu, Hongwen Zhang, Lizhen Wang, Xiaochen Zhao, Huang Han, Qi Guojun, and Yebin Liu. 2023. LatentAvatar: Learning Latent Expression Code for Expressive Neural Head Avatar. In *ACM SIGGRAPH 2023 Conference Proceedings*.
- Gengshan Yang, Minh Vo, Natalia Neverova, Deva Ramanan, Andrea Vedaldi, and Hanbyul Joo. 2022. BANMo: Building Animatable 3D Neural Models From Many Casual Videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2863–2873.
- Haotian Yang, Hao Zhu, Yanru Wang, Mingkai Huang, Qiu Shen, Ruigang Yang, and Xun Cao. 2020. FaceScape: A Large-Scale High Quality 3D Face Dataset and Detailed Riggable 3D Face Prediction. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. 2023. Deformable 3D Gaussians for High-Fidelity Monocular Dynamic Scene Reconstruction. *arXiv preprint arXiv:2309.13101* (2023).
- Tarun Yenamandra, Ayush Tewari, Florian Bernard, Hans-Peter Seidel, Mohamed Elgharib, Daniel Cremers, and Christian Theobalt. 2021. i3DMM: Deep Implicit 3D Morphable Model of Human Heads. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12803–12813.
- Longwen Zhang, Zijun Zhao, Xinzhou Cong, Qixuan Zhang, Shuqi Gu, Yuchong Gao, Rui Zheng, Wei Yang, Lan Xu, and Jingyi Yu. 2023. HACK: Learning a Parametric Head and Neck Model for High-Fidelity Animation. *ACM Trans. Graph.* 42, 4, Article 41 (jul 2023), 20 pages. <https://doi.org/10.1145/3592093>
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *CVPR*.
- Xiaochen Zhao, Lizhen Wang, Jingxiang Sun, Hongwen Zhang, Jinli Suo, and Yebin Liu. 2023. HAvatar: High-fidelity Head Avatar via Facial Model Conditioned Neural Radiance Field. *ACM Trans. Graph.* 43, 1, Article 6 (nov 2023), 16 pages. <https://doi.org/10.1145/3626316>
- Mingwu Zheng, Hongyu Yang, Di Huang, and Liming Chen. 2022a. ImFace: A Nonlinear 3D Morphable Face Model with Implicit Neural Representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Yinglin Zheng, Hao Yang, Ting Zhang, Jianmin Bao, Dongdong Chen, Huangyu Huang, Lu Yuan, Dong Chen, Ming Zeng, and Fang Wen. 2022b. General facial representation learning in a visual-linguistic manner. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18697–18709.
- Yufeng Zheng, Wang Yifan, Gordon Wetzstein, Michael J. Black, and Otmar Hilliges. 2023. PointAvatar: Deformable Point-based Head Avatars from Videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Wojciech Zielonka, Timur Bagautdinov, Shunsuke Saito, Michael Zollhöfer, Justus Thies, and Javier Romero. 2023. Drivable 3D Gaussian Avatars. (2023). arXiv:2311.08581 [cs.CV]
- Wojciech Zielonka, Timo Bolkart, and Justus Thies. 2022. Instant Volumetric Head Avatars. arXiv:2211.12499 [cs.CV]