

LightHeadEd: Relightable & Editable Head Avatars from a Smartphone

Pranav Manu*
IIIT Hyderabad

pranav.m@research.iiit.ac.in

Astitva Srivastava
IIIT Hyderabad

astitva.s@research.iiit.ac.in

Amit Raj
Google Research

amitrajs@google.com

Varun Jampani
Stability AI

varunjampani@gmail.com

Avinash Sharma
IIT Jodhpur

avinashsharma@iitj.ac.in

P.J. Narayanan
IIIT Hyderabad

pjn@iiit.ac.in

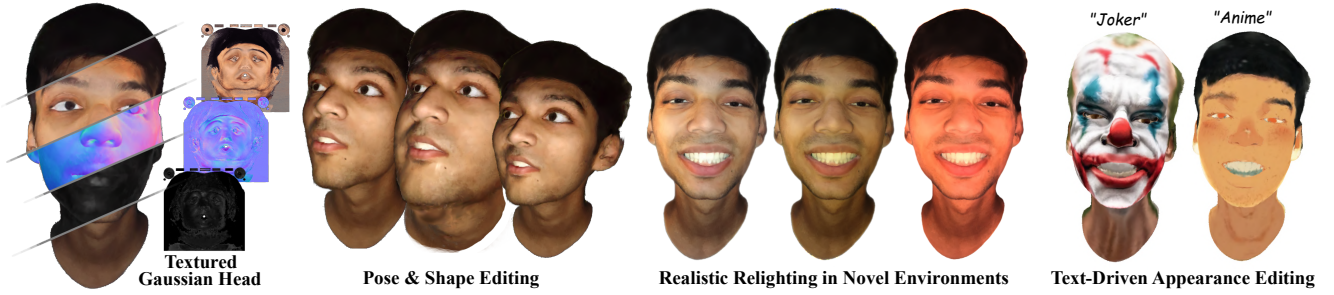


Figure 1. We introduce LightHeadEd to reconstruct realistic head avatars with editing & relighting support.

Abstract

Creating photorealistic, animatable, and relightable 3D head avatars traditionally requires expensive Lightstage with multiple calibrated cameras, making it inaccessible for widespread adoption. To bridge this gap, we present a novel, cost-effective approach for creating high-quality relightable head avatars using only a smartphone equipped with polaroid filters. Our approach involves simultaneously capturing cross-polarized and parallel-polarized video streams in a dark room with a single point-light source, separating the skin’s diffuse and specular components during dynamic facial performances. We introduce a hybrid representation that embeds 2D Gaussians in the UV space of a parametric head model, facilitating efficient real-time rendering while preserving high-fidelity geometric details. Our learning-based neural analysis-by-synthesis pipeline decouples pose and expression-dependent geometrical offsets from appearance, decomposing the surface into albedo, normal, and specular UV texture maps, along with the environment maps. We collect a unique dataset of various subjects performing diverse facial expressions and head movements.

1. Introduction

Relightable, animatable, and editable 3D human head avatars have gained significant popularity in recent years, serving as a versatile and immersive communication medium for the Metaverse, mixed reality platforms, telepresence and beyond. However, creating personalized head avatars with relighting capabilities is challenging due to their dependence on exhaustive multiview data, typically captured using volumetric light-stage systems. These capture systems, being highly sophisticated, expensive, and compute-intensive, are often limited to high-budget production studios. In contrast, enabling affordable 3D head avatars using just a single commodity smartphone could democratize their access to end users, unlocking a wide array of applications in communication and entertainment.

A handful of existing methods [1, 2] attempt to create personalized 3D head avatars from a monocular RGB video captured using a smartphone/webcam. However, they primarily depend on a universal prior model, trained on an exhaustive light-stage dataset of multiple subjects to estimate disentangled facial geometry & reflectance. Such heavy dependence on a dataset/model significantly undermines the affordability and hinders the generalizability to unseen demographics & appearances, while still demanding expensive preprocessing and test-time fine-tuning of the

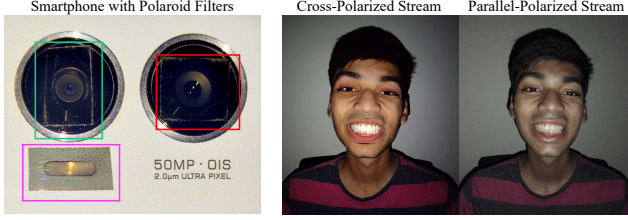


Figure 2. Proposed dynamic capture setup: Smartphone equipped with polaroid filters (left); Cross-Polarized & Parallel-Polarized Monocular Video Streams (right).

model. To eliminate reliance on any prior, several other approaches [3–5] employ inverse physically-based rendering (PBR) to reconstruct a personalized head avatar in a self-supervised analysis-by-synthesis manner. Specifically, [3] uses a parametric head model, FLAME [6], to reconstruct a relightable parametric head mesh from several unconstrained RGB videos of a person. Though FLAME offers animation support via manipulation of low-dimensional pose/expression parameters, it provides only a coarse approximation of the facial geometry and does not model deviations/offsets from the underlying surface (skin deformations, hair, beard, etc.). In order to learn these offsets, recent methods such as [7, 8] propose to combine 3DGS [9] and FLAME head mesh. They propose to sample 3D Gaussian-based splats over a FLAME mesh to learn person-specific details while allowing control over expression/pose. However, these methods are resource intensive, lack relighting support and produce inaccurate geometrical details. In terms of appearance editing, none of the existing monocular methods allow control over the appearance manipulation.

In this paper, we present a novel approach to effortlessly create animatable, editable and relightable head avatars from a commodity smartphone. Instead of relying on a well-calibrated light-stage for the dynamic acquisition of facial performances, we propose a cost-effective and calibration-free capture process to gather realistic head deformations and facial reflectance information. In order to enable relighting, animation and editing support, we propose a novel textured Gaussian head avatar representation, along with an effective self-supervised training methodology to learn a personalized head avatar with high-quality geometry and decomposed albedo, normal and roughness UV maps. Our motivation stems from the fact that the key feature of a light-stage capture setup is multiple RGB cameras equipped with polaroid filters, aiming to decompose skin’s diffuse and specular response to the illumination. Leveraging the same principle, we propose to equip the dual cameras and flashlight of a smartphone with an inexpensive polaroid film, as shown in Figure 2, resulting in a scalable capture process to acquire realistic head deformations of a subject with decomposed surface illumination, in

the form of monocular videos.

Next, we aim to use the captured polarization data to reconstruct the relightable 3D head avatar of the subject with plausible temporal consistency. We first track a FLAME head mesh over the monocular videos, and propose to use 2DGS [10] to model details like skin, hair, beard, etc. Unlike [7, 8] which combine 3DGS with FLAME head mesh, our 2DGS-based representation supports direct surface regularization through normal constraints, achieving high-fidelity facial details while enabling real-time rendering. We embed the 2D Gaussian disks in FLAME’s UV space to learn the Gaussian attributes in the form of UV maps for easy animation and texturing support. Learning appearance in UV space helps eliminate the need for memory-intensive SH-coefficients and more importantly, enables flexible appearance editing by altering the albedo map. We decouple reflectance in the form of albedo, normals and roughness UV maps. Additionally, to handle deformations and appearance changes, we propose to learn expression-dependent residual UV maps. Furthermore, we also learn an environment cubemap to support relighting in arbitrary lighting conditions.

The proposed method for head avatar capture & reconstruction is highly efficient and enables affordable relightable heads with several intriguing features. We demonstrate superior quality results, backed by comprehensive quantitative and ablative analysis, while showcasing several useful applications, such as shape editing & text-guided appearance editing.

In summary, our contributions include:

- We propose an affordable and scalable process for the dynamic acquisition of polarized facial performances from a commodity smartphone equipped with a polaroid film.
- We introduce a novel 2DGS-based head avatar representation with relighting, texture mapping & editing support; and a novel methodology to learn the aforementioned representation from a polarized monocular video sequence.
- We collect a first-of-its-kind dataset of polarized facial performances with diverse facial expressions and head movements.

We plan to release the code and dataset publicly to drive the advancements in relightable head avatars research.

2. Related Works

Lightstage Capture Systems: Polarization has long been used to decompose scene illumination[11–13], leveraging the fact that the single bounce specular reflection does not alter the polarization of the incoming light. Using polarization filters and controlled lighting, Debevec et al.[14] pioneered the first Lightstage system to estimate human face reflectance, efficiently capturing how the face appears when lit from every possible lighting direction. Though the initial setup was meant for static capture,

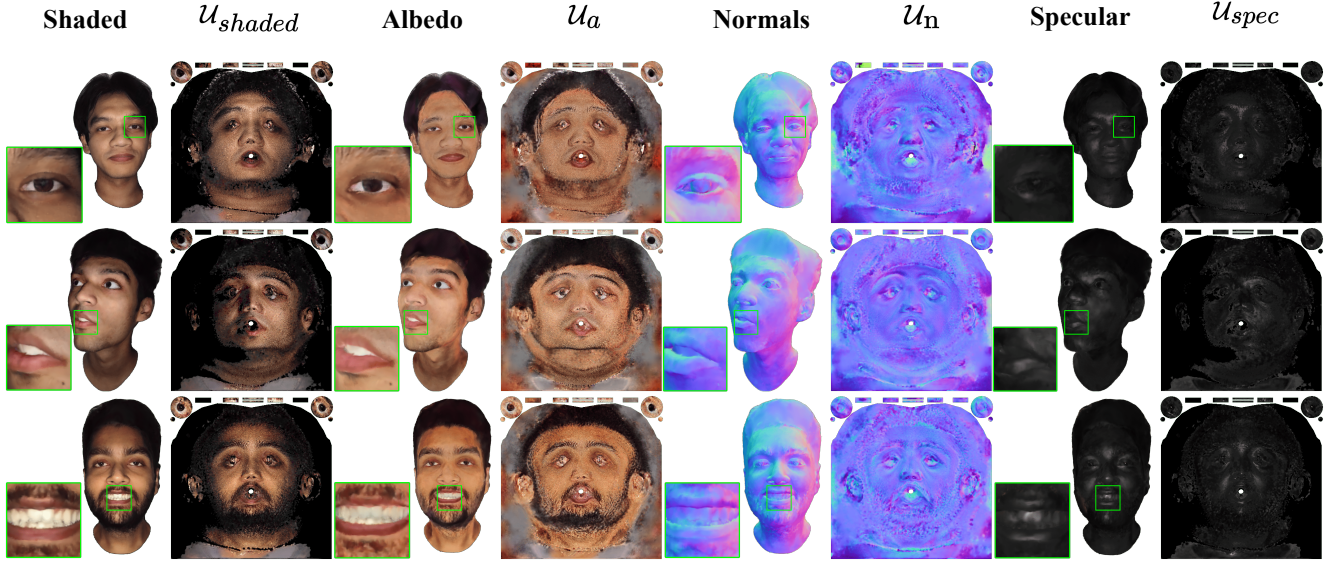


Figure 3. Decomposition of appearance & geometry in UV space.

subsequent advancements[15–17] proposed techniques to compensate for motion, expanding the capture area & improving the surface fidelity. [18] proposed a multi-view setup for dynamic facial texture acquisition without the need for polarized illumination. Following this, [19] reintroduced polarization without active illumination to model subsurface-scattering. Recently proposed [20] proposed further improvements in the system to include global illumination and polarization modeling. Though these volumetric lightstage-based solutions for capturing face reflectance & geometry deliver impressive visual results, they rely on sophisticated hardware, making them highly expensive and bulky, while requiring a significant level of expertise to operate. Hence, they face challenges in scaling to large numbers of identities.

Data-driven Personalized Head Avatars: Facial expression control has advanced from blendshape-based methods in VFX and gaming to data-driven techniques that estimate statistical bases from 3D head scans, enabling nuanced head pose and expression control [21–26]. However, models like FLAME [6] often miss subtle expressions and detailed, person-specific geometry (e.g., hair, beard, skin deformations). For photorealistic avatars, several methods use differential rendering from multiview videos in an analysis-by-synthesis framework [5, 27–32]. While NeRF-based approaches yield personalized avatars, they suffer from slow rendering and coarse geometry, motivating recent methods to employ 3D Gaussian splats (3DGS) [9] for efficient real-time performance—albeit with illumination baked into the appearance, which prevents relighting. A few works [2, 33] achieve relightable avatars using lightstage data, but their dependence on such systems limits

demographic diversity. Diffusion-based techniques [34–36] and neural editing methods [37, 38] enable texture edits yet remain slow and struggle with novel poses. A recent approach [39] uses neural texture maps for flexible editing but also relies on optimization-heavy processes. In this paper, we propose an affordable, scalable method for creating animatable, relightable, and editable head avatars.

Head Avatars from Monocular Video(s): Creating animatable head avatars from monocular videos is yet another interesting & challenging direction. Several existing learning-based methods aim to eliminate the need for multiview information by tracking dynamically evolving head topology over a monocular video sequence, either in the form of a parametric head mesh [3] or a neural head representation [40]. For more detailed and expressive monocular head avatars, existing methods propose to use primitives, such as points[4] or 3D Gaussians[8], on top of FLAME[6] model. Methods like [41] & Gaussian Blendshapes[7] attempt to use Gaussian splats as their primitives, but owing to the nature of ill-defined normals of 3D Gaussians, lack high-quality surface details. These methods also fail to provide any editability support to head meshes. PointAvatar[4] aims to reconstruct a relightable human head avatar in a highly constrained fixed illumination setting. While it achieves impressive rendering quality, the relighting is not Physically-based, and does not take into account the properties of skin reflectance. Moreover, points as primitives are memory intensive in nature, making [4] unsuitable for real-time rendering. Furthermore, none of the existing monocular head avatar methods directly provides texture maps to enable flexible & quick editing appearance editing. Though, FlashAvatar[42] embeds 3D Gaussians within FLAME’s

UV space for initialization, it doesn't learn any UV texture maps and lacks relighting support.

3. Method

To create affordable person-specific relightable head avatars in a monocular setting, we propose an effortless polarized data acquisition process. The captured polarized data is then used to create a novel textured head avatar representation. In order to control head pose & expressions, we propose to use FLAME, combined with 2DGS[10] to model person-specific offset details & appearance. For relighting, we propose a novel self-supervised learning scheme to decompose the polarized information into appearance and shading information in the form of albedo, normal & roughness UV maps in FLAME's parametrization space (as shown in Figure 3), while also learning an environment cubemap. Additionally, to model dynamic changes in the appearance and geometry due to pose/expressions, we learn expression-dependent residual UV maps. We now describe each component detail.

3.1. Polarized Dynamic Data Capture

In order to capture surface illumination information, [43] proposed a static smartphone-based capture setup, utilizing a single back camera with polaroid filter to sequentially capture several cross and parallel-polarized images of a person one-at-a-time, demanding the person to remain stationary for an extended period. However, expanding this setup to a dynamic scenario is non-trivial and requires careful consideration to achieve a scalable and calibration-free capture process. In order to achieve this, we mount a single smartphone on a tripod in a dark room with its flashlight as the only source of illumination (point light) to avoid estimation of complex scene lighting. As shown in Figure 2, we cover the flashlight using a thin linearly polarized filter/film (highlighted in purple), about 4.8 microns in thickness. We also cover the two back cameras of the smartphone with the same polarized film, with a relative angle difference of 0° (red) and 90° (green) w.r.t. the flashlight's film, essentially making one camera parallel-aligned and the other one cross-aligned to the flashlight's polarization. The usage of polaroid filters introduces a tint-shift between the capture stream of the two cameras as can be seen in Figure 2. While [43] precomputes an approximate affine color correction matrix, we delegate the tint-shift correction during the training as discussed in subsection 3.3, to make the capture process seamless.

Using this setup, we simultaneously record two uninterrupted monocular videos (cross & parallel-polarized) of a human subject enacting a predefined set of diverse expressions & poses to ensure adequate deformations & lighting information from different angles is captured. Both the video streams are captured in 1920×1080 resolution at 24

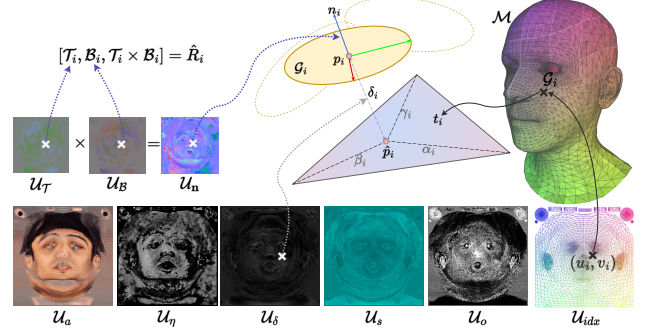


Figure 4. Proposed textured Gaussian head representation with primary UV attribute maps.

FPS. It is important to note that, though we use two back cameras at the same time to ensure both the video streams are in sync, we do not use any stereo-based depth or alignment information due to negligible baseline distance between the two cameras, thereby facilitating avatar creation by treating the two streams as separate monocular RGB videos.

3.2. Textured Gaussian Head Representation

Given the cross-polarized and parallel polarized monocular video sequences, we employ the head-tracking [44] proposed in [45] to obtain the subject-specific FLAME head mesh \mathcal{M} temporally tracked separately over both monocular sequences. We propose to model the remaining subject-specific details on top of FLAME mesh using 2D Gaussian Splats (2DGS)[10], which are essentially ‘flat’ 2D planar elliptical disks embedded in 3D space. 2DGS uses an explicit ray-splat intersection technique, resulting in a perspective-correct splitting and more accurate surface reconstruction. Additionally, this enables direct surface regularization through normal constraints[10], improving the quality of normals for shading. Each 2D Gaussian disk \mathcal{G}_i is characterized using its mean position $p_i \in \mathbb{R}^3$ in 3D space, its 2D scale $S_i \in \mathbb{R}^2$, a quaternion $q_i \in \mathbb{R}^4$ defining its orientation, and its opacity $o_i \in \mathbb{R}^1$. Similar to 3DGS [9], the factorized covariance matrix Σ_i for each 2D Gaussian \mathcal{G}_i is defined as $\Sigma_i = R_i S_i' S_i'^T R_i^T$, where R_i is the rotation matrix, derived from the 6D rotation vector r_i , and $S_i' = [S_i, 0]^T \in \mathbb{R}^3$ (i.e. the third dimension of 3D scale vector is set to zero to represent a flat 2D Gaussian in 3D space). However, we propose to associate all the learnable attributes of 2D Gaussians to FLAME's UV space, by embedding the Gaussians in tangent space of the triangulated mesh, instead of 3D object space.

Embedding 2D Gaussians in UV space: We begin with initializing a 2D Gaussian \mathcal{G}_i for a UV-coordinate (u_i, v_i) in FLAME's canonical UV space \mathcal{U} , as shown in Figure 4. We also estimate a face index map, \mathcal{U}_{idx} , to store the triangle

(face) index m on which a UV-coordinate of \mathcal{U} lies (undefined for empty UV space); this is a one-time computation. To transform Gaussians from canonical pose/expression to the deformed pose/expression directly via posed \mathcal{M} , we first compute the associate face index $m = \mathcal{U}_{idx}(u_i, v_i)$, which gives us a triangle $t_m = (v_{m_0}, v_{m_1}, v_{m_2}) \in 3 \times \mathbb{R}^3$ (3D positions of t_m 's vertices). For each Gaussian G_i , we compute the initial mean position in the posed 3D space as, $\hat{p}_i = (\alpha_i * v_{m_0}) + (\beta_i * v_{m_1}) + (\gamma_i * v_{m_2})$, where $\alpha_i, \beta_i, \gamma_i$ are the barycentric coordinates for the point (u_i, v_i) & $\alpha_i + \beta_i + \gamma_i = 1$. This formulation enables us to define the remaining attributes of the Gaussians (orientation, scale, offsets, appearance, etc.) as learnable UV-maps.

Gaussian Attributes as UV Maps: Given the aforementioned formulation, we now explain how we model the Gaussian attributes as texture maps. We sample 2D Gaussians for every valid UV-textel (excluding empty space) and define a set of *primary* UV-maps for each attribute— \mathcal{U}_a (albedo map), \mathcal{U}_η (roughness map), $\mathcal{U}_\mathcal{T}$ (tangent map), $\mathcal{U}_\mathcal{B}$ (bitangent map), \mathcal{U}_δ (offset map), \mathcal{U}_s (scale map) & \mathcal{U}_o (opacity map). As shown in Figure 4, for each Gaussian G_i with its own UV-coordinate, we can query any attribute $\Omega_i = \mathcal{U}_\Omega(u_i, v_i)$, where Ω_i can be albedo color $a_i \in \mathbb{R}^2$, roughness $\eta_i \in \mathbb{R}$, tangent vector $\mathcal{T}_i \in \mathbb{R}^3$, bitangent vector $\mathcal{B}_i \in \mathbb{R}^3$, 2D scale $s_i \in \mathbb{R}^2$, and opacity $o_i \in \mathbb{R}$. To estimate orientation for the 2D Gaussian G_i , we first query tangent $\mathcal{T}_i = \mathcal{U}_\mathcal{T}(u_i, v_i)$ & bitangent $\mathcal{B}_i = \mathcal{U}_\mathcal{B}(u_i, v_i)$ and perform Gram–Schmidt orthogonalization. The tangent space rotation for Gaussian G_i is defined as $\hat{R}_i = [\mathcal{T}_i, \mathcal{B}_i, \mathcal{T}_i \times \mathcal{B}_i] \in \mathbb{R}^{3 \times 3}$. We then obtain the 3D world space rotation/orientation $R_i = \mathbb{T} * \hat{R}_i$, where \mathbb{T} is a 3×3 transformation matrix for transforming a vector from tangent space to world space. Additionally, we compute the normal vector $\mathcal{N}_i = \mathcal{T}_i \times \mathcal{B}_i$, which is required for shading. To account for geometrical details far from the surface of the parametric head mesh \mathcal{M} , we define offset $\delta_i = \mathcal{U}_\delta(u_i, v_i)$ along the normal direction and shift the initial mean position \hat{p}_i of the 2D Gaussian in 3D space as follows:

$$p_i = \hat{p}_i + \xi_i * \mathbb{T} * \delta_i \quad (1)$$

where ξ_i is the area-adjustment factor, which restricts the Gaussians to go too far from their associated triangles or grow too big in scale, and is computed as:

$$\xi_i = (\alpha_i * a_{m_0}) + (\beta_i * a_{m_1}) + (\gamma_i * a_{m_2}) \quad (2)$$

$$a_{m_j} = \frac{1}{|\mathcal{N}(v_{m_j})|} \sum_{k \in \mathcal{N}(v_{m_j})} \sqrt{\text{Area}(t_k)} \quad (3)$$

here, α, β, γ are barycentric coordinates, and a_{ij} for vertex v_{ij} is the mean-root-sum of area of triangles in v_{m_j} 's neighborhood $\mathcal{N}(v_{m_j})$.

Residual UV Maps: In addition to primary UV-maps, we also maintain a set of UV-maps to store expression-dependent residuals. We define k residual UV-maps, where each (u_i, v_i) coordinate stores residuals $\Delta\Omega_i$ on top of primary attributes Ω_i to model different geometrical and appearance changes observed throughout the sequence. Specifically, we store $\Delta p_i, \Delta a_i, \Delta \mathcal{T}_i, \Delta \mathcal{B}_i, \Delta s_i, \Delta \eta_i$. In order to provide expression guidance, we take expression parameters of FLAME head mesh, $\psi \in \mathbb{R}^{100}$, and project them onto a k -dimensional space using a projection matrix, $\mathbf{\Pi} \in \mathbb{R}^{k \times 100}$, to obtain linear blend weights $\mathcal{W} = \mathbf{\Pi} \cdot \psi \in \mathbb{R}^k$. The final blended residuals for a Gaussian G_i are obtained as follows:

$$\Delta\Omega_i = \mathcal{U}_\Delta(u_i, v_i) = \sum_{w_j \in \mathcal{W}} w_j * \mathcal{U}_{\Delta_j}(u_i, v_i) \quad (4)$$

The residuals are added to primary attributes Ω_i to obtain final attributes Ω'_i using the following equation:

$$\Omega'_i = \begin{cases} \Omega_i + \Delta\Omega_i & ; \text{for } p_i, \mathcal{T}_i, \mathcal{B}_i \\ \Omega_i * \exp(\Delta\Omega_i) & ; \text{for } a_i, s_i, \eta_i \\ \Omega_i & ; \text{for } o_i. \end{cases} \quad (5)$$

3.3. Learning Facial Geometry & Reflectance

Rendering Equation & BRDF: For appearance modelling, the original 2DGS uses Spherical Harmonics[46] to learn the appearance of each 2D Gaussian. Though SH-coefficients allow for an accurate view-dependent appearance for static scenes, higher-order SH-coefficients are needed to model high-frequency details (e.g. surface normals and roughness), exponentially increasing the memory requirements. More importantly, appearance editing is not directly possible when using such representation. This motivates us to replace SH-coefficients with a single RGB color for appearance c_i . For view-dependent color, we implement a new Physically-based SVBRDF[47] shader for 2DGS, to estimate the Gaussian's appearance given the viewing direction ' ω ' and light direction ' l ' from a point source (flashlight). Furthermore, for learning shading to enable relighting, we disentangle the appearance c_i into albedo, specular and diffused components, namely a_i, f_{s_i} & f_{d_i} . For specular component f_{s_i} , we use Cook-Torrance[48] microfacet specular BRDF defined as follows:

$$f_{s_i}(l, \omega, \eta_i) = k_s * \frac{D(h)F(\omega, h)G(l, \omega, h)}{4(n \cdot l)(n \cdot \omega)} \quad (6)$$

where, h is the half vector, bisecting the angle between l & ω ; k_s is a constant specular gain. Unlike [43], which uses a spatially varying specular gain k_s (which is a learnable constant in our case), we instead use spatially varying roughness η_i associated with each Gaussian \mathcal{G}_i to handle

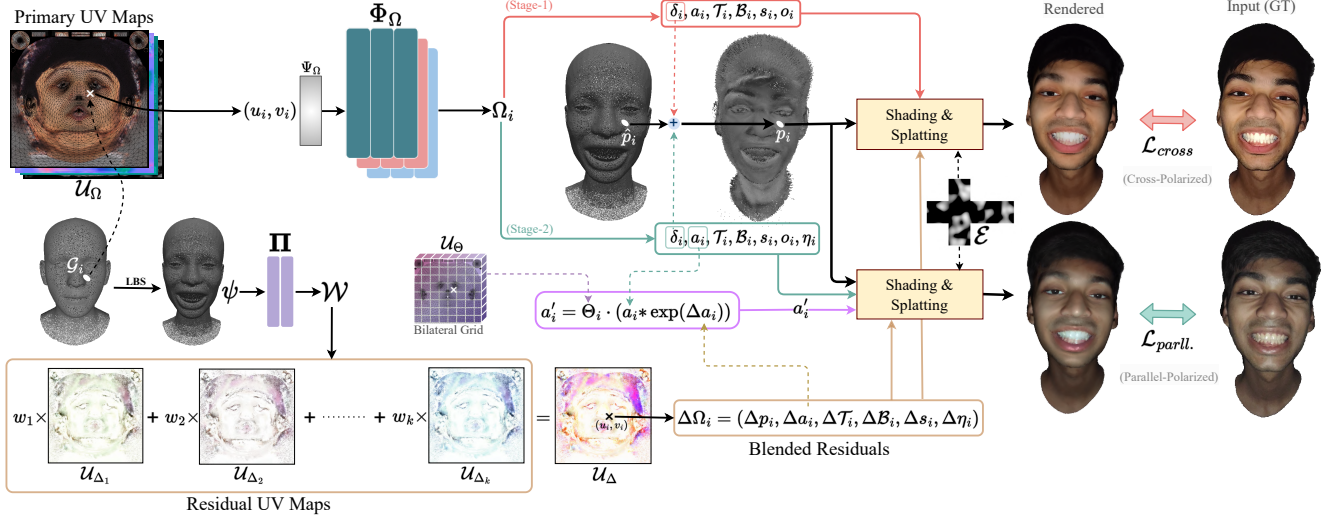


Figure 5. Proposed two-stage training strategy to learn textured Gaussian head avatars with decomposed appearance and geometry.

difference in specularity of skin, hair, teeth etc. Following [43], we use Schlick’s approximation[49] for the Fresnel term F . For the NDF term D , we use an alternative approximation proposed by Trowbridge-Reitz[50]:

$$D(h, \eta_i) = \frac{\eta_i^2}{\pi((n \cdot h)^2(\eta_i^2 - 1) + 1)^2} \quad (7)$$

Our geometric term G uses Smith’s variant of Shlick-GGX approximation[49]. Please note that similar to [43], we assume that $l = \omega$ since the flashlight is very close to the camera lens, making the implementation straightforward.

$$G(l, \omega, h) = G_{Schlick}(l) * G_{Schlick}(\omega) = [G_{Schlick}(\omega)]^2 \quad (8)$$

$$G(l, \omega, h) = G_{Schlick}^2(\omega) = \left[\frac{n \cdot \omega}{(n \cdot \omega)(1 - \lambda) + \lambda} \right]^2 \quad (9)$$

where $\lambda = \eta_i/2$ for remapping Schlick-GGX to match with Smith’s formulation[51].

For the diffuse component, we use the BRDF model proposed by Ashikhmin & Shirley [52]:

$$f_{d_i}(a_i, \omega) = \frac{28a_i}{23\pi} (1 - F_0) (1 - (1 - \frac{n^\top \omega}{2})^5)^2, \quad (10)$$

where a_i is the albedo color and $F_0 = 0.04$ is the reflectance of the skin at normal incidence. Besides shading using point light, we also incorporate ambient shading component f_{env} to account for the light bouncing around the environment. We follow the differentiable version of the *split sum* shading model proposed in [53] to learn environment lighting from image observations through optimization. The final shaded

color is computed as:

$$c_i = f_{d_i} + f_{s_i} + f_{env}. \quad (11)$$

Figure 5 illustrates our novel training methodology for learning a textured Gaussian head avatar. Given cross-polarized and parallel-polarized video sequences, we learn the aforementioned Gaussian attributes in two separate stages. We employ a set of *stacked* MLPs, Φ_Ω , which consists of several hash-encoded MLPs[54] with hash-encoding Ψ_Ω , to predict UV attributes $\Omega_i = \mathcal{U}_\Omega(u_i, v_i)$ for a 2D Gaussian G_i . For residual UV-maps, we initialize learnable UV-tensors \mathcal{U}_{Δ_j} with $z = |\Delta\Omega_i|$ channels, where $j = 1$ to k . To handle the global tint-shift between the cross-polarized and parallel-polarized frames, we use a low-dimensional Bilateral Grid[55] \mathcal{U}_Θ with learnable affine transformation matrix Θ_i for each UV-textel. Additionally, since we allow (u_i, v_i) to optimize over time, we also pre-compute the triangle-index map \mathcal{U}_{idx} for querying the triangle index from UV-coordinates. Finally, we define a learnable cube-map \mathcal{E} to store the environment lighting information. All the values specified within \mathcal{U}_Ω , \mathcal{U}_{Δ_j} , \mathcal{U}_Θ , Π and \mathcal{E} are initialized at random and optimized via differential-rendering losses in a self-supervised manner, using cross-polarized sequence in the first stage and parallel-polarized sequence in the second.

During the first stage, we focus on learning \mathcal{U}_Ω (except the roughness map \mathcal{U}_η), along with \mathcal{U}_{Δ_j} . We feed (u_i, v_i) as input to hash-encoded MLP Φ_Ω to predict $\Omega_i = \Phi_\Omega(\Psi_\Omega(u_i, v_i))$, i.e. albedo a_i , tangent \mathcal{T}_i , bitangent \mathcal{B}_i , scale s_i , opacity o_i , and estimate shifted mean position p_i (Equation 1). For a given pose/expression θ , we compute linear blending weights $\psi = \Pi \cdot \theta$ and use Equation 4 to obtain $\Delta\Omega_i$. We then add the residuals to the primary attributes using Equation 5. At last, we compute final shaded

color c_i via Equation 11 (ignoring the specular component) and use it along with $\mathcal{T}_i, \mathcal{B}_i, s_i, o_i$ to perform 2D Gaussian splatting to rendered the image & minimize loss \mathcal{L}_{cross} .

During the second stage, we freeze the optimization of primary & residual UV-maps (\mathcal{U}_a & $\mathcal{U}_{\Delta a}$), and focus on learning roughness only. We also learn bilateral grid \mathcal{U}_Θ to account for the tint shift between the albedo learned from cross-polarized images (in the previous stage) and parallel-polarized images. We obtain tint-corrected albedo as:

$$a'_i = \Theta_i \cdot (a_i * \exp(\Delta a_i)). \quad (12)$$

We include both diffuse and specular components during second stage while computing shaded color c_i (Equation 11). We query the remaining attributes learned in the first stage and obtain a rendered image via splitting and minimize the loss \mathcal{L}_{parll} . We define both \mathcal{L}_{cross} and \mathcal{L}_{parll} as:

$$\mathcal{L}_{cross/parall} = \mathcal{L}_1 + \mathcal{L}_{SSIM} + \mathcal{L}_{LPIPS} + \mathcal{L}_{scale} \quad (13)$$

where, $\mathcal{L}_{scale} = \sum_i \max(0, s_i - 0.3)$ to prevent Gaussians to grow too large.

In both the stages, we initialize the environment cube map \mathcal{E} at random and optimize its values via Equation 11. During inference, for efficient rendering, we discard the MLP Φ_Ω and use attributes stored in the primary UV-maps \mathcal{U}_Ω , while relying on Π & residual UV-maps \mathcal{U}_Δ to handle poses/expressions-dependent deformations.

4. DuoPolo Dataset

We also introduce ‘‘DuoPolo’’ dataset, which consists of cross-polarized & parallel-polarized video sequences of 10 subjects, captured using the proposed capture setup. Each data sample consists of a 90s video sequence of a human subject enacting diverse head poses and facial expressions, followed by a short utterance of a phrase (to capture subtle lip motions during talking). Each sequence is captured in 1920×1080 resolution at 24 FPS. Along with the video, we also provide per-frame background/foreground segmentation mask, 3D facial landmarks, and parametric head mesh, FLAME[6], tracked over both the videos using [44]. The proposed dataset is first-of-its-kind which aims to make polaroid facial performances accessible to everyone, bridging the gap between expensive light-stage data and affordable head avatars.

5. Training & Implementation Details

We learn 1024×1024 size texture maps via Hash encoded Coordinate MLPs. The Hash-encoded MLP for texture map synthesis has 1 hidden layer with 16 neurons. The Hash Table has 7 levels, each with a size 2^{20} . The base resolution of the hash table was 512, with a growth factor of 1.26. The pose-dependent MLP had 2 hidden layers, with

32 neurons each. The Hash grid for the pose-dependent MLP had 16 levels, each of size 2^{18} with a base resolution of 16 and a growth factor of 1.4. Both the MLP are trained using AdamW[56] optimizer, with a learning rate of $1e^{-3}$.

We use a bilateral grid to color-match the albedo map between the cross-polarized and parallel-polarized streams. The size of the learnable bilateral grid is $16 \times 16 \times 8 \times 12$ which is optimized via AdamW optimizer with a learning rate equal to $1e^{-3}$.

Each sequence is trained for $16k$ iterations in the first stage and $100k$ iterations in the second stage. All our experiments were performed on an Nvidia RTX A6000 GPU.

6. Experiments & Results

Evaluation Datasets & Metrics: For experimental evaluation we use sequences from our proposed DuoPolo dataset. For each sequence, we perform a train-test split where we use the first 80% of the frames for training and the remaining 20% for testing. For quantitative evaluation, we use standard metrics used by existing state-of-the-art methods. Following the evaluation technique of [4], we compute per-sequence (or per-subject) L_1 error, PSNR[57] & SSIM[58] between the rendered frames and input (ground truth) frames from the test split.

6.1. Qualitative Results

Learned Reflectance Maps for Relighting: Figure 3 demonstrates the ability of the proposed method to capture Physically plausible facial reflectance by separating appearance into albedo, roughness, and normals in the form of UV maps, while also reconstructing high-frequency geometrical details. This disentanglement allows relighting the reconstructed head avatars using diverse environment maps as shown in Figure 9.

Editing: Our proposed textured human head representation allows shape editing by changing the shape parameters of underlying FLAME mesh (Figure 6), as well as appearance editing, by modifying only the albedo map, similar to classical textured mesh editing. We show an example of appearance editing in Figure 1 where we manipulate the albedo via text using the approach proposed in CLIP-Head [34].

6.2. Comparison

Since most of the learning-based monocular head avatar methods assume a uniformly lit environment with fixed illumination and do not handle learning on polarization data, for fair evaluation, we compare our training methodology with existing SOTA methods on cross-polarized (diffused) sequences from our captured data. In addition to this, we also use samples from INSTA[54] dataset and 3DGB dataset proposed in [7]. We follow the complete evaluation strategy of GaussianBlenshapes (GBS) [7] and report PSNR, LPIPS & SSIM in Table 1, where we outperform



Figure 6. Shape editing over reconstructed head avatars.

Table 1. Comparison of different methods on INSTA, 3DGB and Our Dataset. Best values are highlighted in green, while second-best values are in yellow.

Method	INSTA[54]			3DGB [7]			DuoPolo (Ours)		
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
Point-Avatar	29.12	0.932	0.094	31.76	0.926	0.145	27.77	0.919	0.128
FlashAvatar	30.14	0.942	0.038	25.92	0.920	0.094	29.79	0.923	0.089
Gaussian Deja-vu	25.56	0.925	0.058	23.45	0.910	0.079	24.73	0.848	0.183
SPARK	23.75	0.873	0.103	24.70	0.863	0.101	23.75	0.873	0.103
GaussianBlendshapes	30.01	0.951	0.084	32.05	0.942	0.144	30.02	0.943	0.094
Ours	30.33	0.952	0.036	32.92	0.943	0.046	31.98	0.955	0.053



Figure 7. Comparison for surface normals of the reconstructed facial geometry.

all the existing methods for head avatar reconstruction from monocular video. We show qualitative comparison in Figure 8 where we show superior quality results compared to others. Our method is able to capture fine details (as highlighted in boxes) while also enabling relighting and editing, which other methods do not directly support.

We also compare the surface normals of the head avatar with some of the methods which aim to reconstruct facial geometry besides appearance in Figure 7. As shown, PointAvatar produces oversmooth geometry due to inability of point splats to represent detailed geometry. While [4] GBS[7] produces a lot of spiky artifacts, mainly due to the inherent limitation of 3D Gaussians as they are not suitable for surface reconstruction. SPARK[3] models fa-

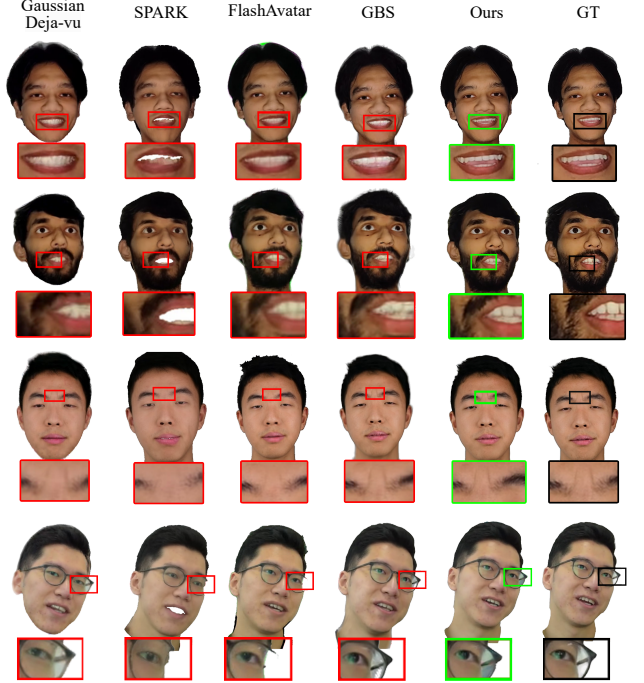


Figure 8. Qualitative comparison with SOTA methods. Our method is able to better capture finer details (such as teeth and eyes) whilst supporting relighting, pose and expression control.

cial geometry as a mesh by learning offsets over the underlying FLAME head mesh, however fails to capture high-frequency details as wrinkles, beard, loose skin etc. On the other hand, our proposed 2DGS-based head representation results in high-quality surface normals capturing subject-specific finer details.

6.3. Ablation Study

We perform an ablative analysis of key components proposed in our method. We demonstrate the effect of using residual UV maps $\Delta\Omega$ in Table 2. We also demonstrate the effect of various choices of low-dimensional projection space for linear basis blend weights (value of k) in Table 3, where we can observe that choosing higher value for k improves the quality of results, but also significantly in-

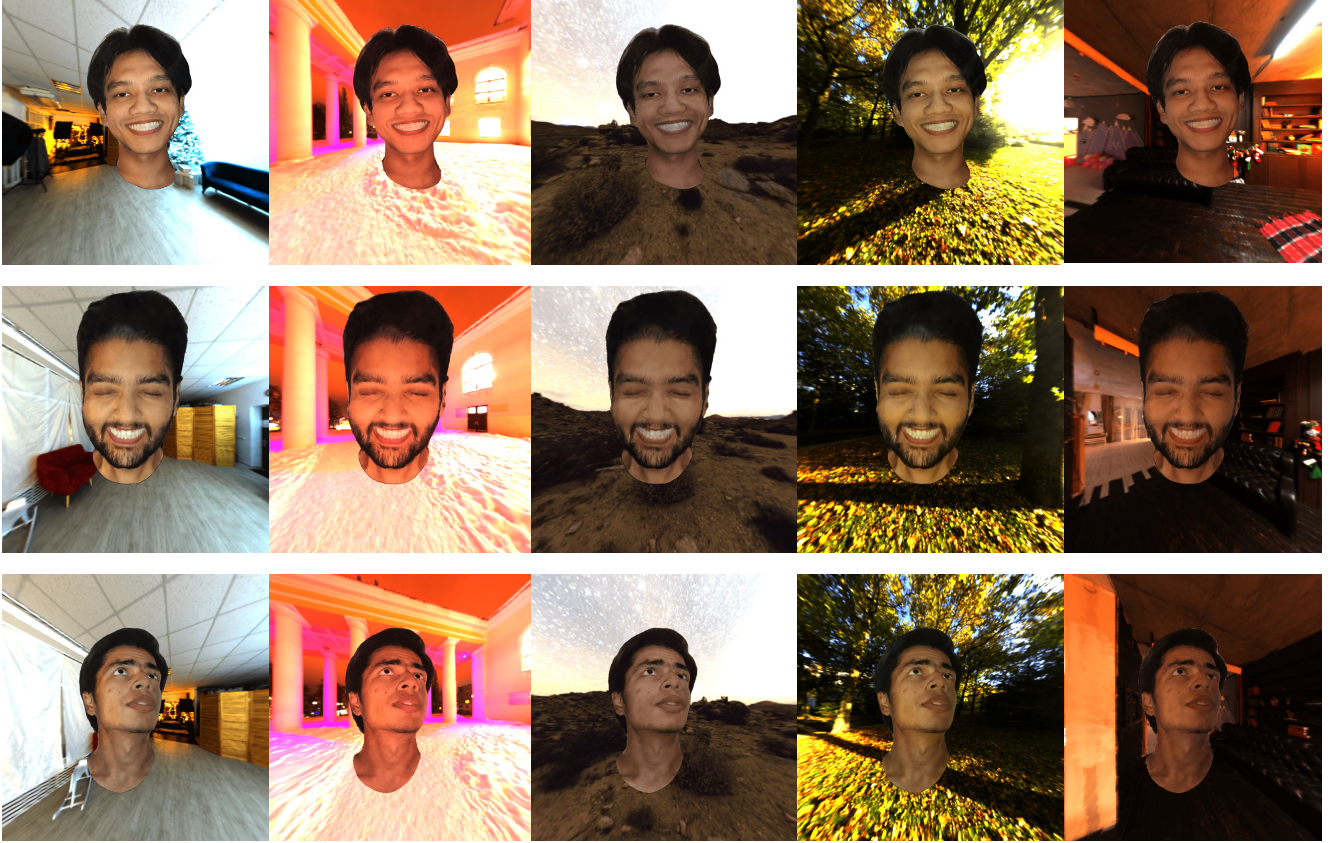


Figure 9. Relighting the reconstructed head avatars in diverse environments.

creases storage. Increasing the value of k will also increase the number of learning parameters, which will increase the training time as well. We also show the effect of resolution of UV maps in Table 4, where we demonstrate the trade-off between rendering quality and training time. Higher resolution maps require more Gaussians, leading to slow convergence. However, the improvements in rendering quality are not drastic; therefore, we use 512×512 as the default resolution.

Table 2. Effect of Residual Maps

Texture Maps	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
w $\Delta\Omega_i$	33.12	0.953	0.048
w/o $\Delta\Omega_i$	30.99	0.945	0.053

Table 3. Ablation over values of k for residual basis

Num. Basis Maps	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Storage(MB) \downarrow
5	32.46	0.953	0.039	52.4
25	33.91	0.960	0.035	262.2
50	33.62	0.960	0.036	524.3
100	33.59	0.960	0.037	1048

Table 4. Effect of TexMap resolution

Texture Resolution	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Time(min) \downarrow
128x128	33.12	0.953	0.048	5.78
256x256	33.65	0.957	0.039	6.2
512x512	33.84	0.959	0.036	11.2
1024x1024	33.92	0.961	0.033	28.8

7. Discussion

7.1. Limitations

Our method has the potential to democratize access to relightable head avatars by making polarized facial performance acquisition more accessible and cost-effective. Using the proposed capture approach, we create a monocular polarized dataset that is unique in its kind. However, its quality cannot be directly compared to data captured using advanced volumetric lightstage systems. The resolution of our dataset is constrained by the limitations of smartphone cameras, which occasionally leads to inaccuracies when modeling fine-grained specular reflections. Additionally, the range of lighting variations is restricted, making it challenging to generalize relighting for complex environment maps.

The proposed representation is further constrained by the parametric head model, which struggles to accurately represent intricate inner-mouth regions such as the tongue and teeth, occasionally resulting in blurry details in those regions. Moreover, the 2DGS representation is unable to capture fine details like individual beard hairs, hair strands, or eyelashes. To address this, we plan to investigate neural strand-based representations for more precise modeling of these facial attributes. Additionally, we aim to explore more accurate techniques for differentiable rendering, such as differentiable ray tracing, as opposed to Gaussian splatting, to better model physics-based reflectance properties.

7.2. Ethical Concerns

Our work involves the collection of a dataset, comprising facial information from several individuals. We acknowledge the potential privacy concerns associated with such data and have taken proactive measures to address these issues responsibly. Participation in the dataset collection was entirely voluntary, and we obtained explicit informed consent from all subjects prior to their inclusion. Each participant was fully briefed about the scope, purpose, and potential uses of the dataset, ensuring transparency and understanding. To further safeguard privacy and mitigate risks, we have implemented stringent data access protocols. Access to the dataset is only granted after a thorough review and approval process, requiring researchers to submit a formal application and agree to adhere to ethical guidelines. We emphasize that the dataset will be used solely for academic and research purposes, and its distribution is controlled to prevent misuse.

While our method has been proposed to advance research in relightable head avatars, we recognize that it could potentially be misused in unethical ways, such as creating deepfake content. Deepfakes, while having legitimate applications in entertainment and education, are frequently associated with harmful consequences, including misinformation, identity theft, and digital harassment. We believe that addressing these ethical considerations is essential to balance the benefits of technological progress with its social implications. By fostering transparency, accountability, and oversight, by restricting the usage of our dataset and approach for research purposes only, we aim to minimize risks while enabling the legitimate and responsible use of our research.

7.3. Conclusion

We proposed LightHeadEd, a scalable and affordable method for capturing and reconstructing animatable, relightable, and editable head avatars using a commodity smartphone with polaroid filters. Our approach introduces a seamless dynamic acquisition process for polaroid facial performances, a textured Gaussian head representation with

2D Gaussian attributes embedded in the UV space of a parametric head mesh, and a self-supervised training scheme to reconstruct head avatars from monocular video. LightHeadEd enables applications like relighting and text-based appearance editing, delivering superior geometry, appearance, and physically-accurate shading compared to existing methods. While finer details like hair strands and skin pores remain current limitations, we plan to address these in the future along with support for higher-resolution images (up to 4K) and the collection of a larger, diverse dataset to facilitate a more generalizable universal head model.

References

- [1] 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. Authentic volumetric avatars from a phone scan. *ACM Trans. Graph.*, 41(4), July 2022. ISSN 0730-0301. doi: 10.1145/3528223.3530143. URL <https://doi.org/10.1145/3528223.3530143>. 1
- [2] Junxuan Li, Chen Cao, Gabriel Schwartz, Rawal Khirodkar, Christian Richardt, Tomas Simon, Yaser Sheikh, and Shun-suke Saito. Uravatar: Universal relightable gaussian codec avatars. In *ACM SIGGRAPH 2024 Conference Papers*, 2024. 1, 3
- [3] Kelian Baert, Shrisha Bharadwaj, Fabien Castan, Benoit Maujean, Marc Christie, Victoria Abrevaya, and Adnane Boukhayma. Spark: Self-supervised personalized real-time monocular face capture. 2024. doi: 10.1145/3680528.3687704. 2, 3, 9
- [4] 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 (CVPR)*, 2023. 3, 7, 9
- [5] Guy Gafni, Justus Thies, Michael Zollhöfer, 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 (CVPR)*, pages 8649–8658, June 2021. 2, 3
- [6] Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6):194:1–194:17, 2017. URL <https://doi.org/10.1145/3130800.3130813>. 2, 3, 7
- [7] Shengjie Ma, Yanlin Weng, Tianjia Shao, and Kun Zhou. 3d gaussian blendshapes for head avatar animation. In *ACM SIGGRAPH Conference Proceedings, Denver, CO, United States, July 28 - August 1, 2024*, 2024. 2, 3, 7, 9
- [8] Yufan Chen, Lizhen Wang, Qijing Li, Hongjiang Xiao, Shengping Zhang, Hongxun Yao, and Yebin Liu. Monogaussianavatar: Monocular gaussian point-based head avatar. *arXiv*, 2023. 2, 3
- [9] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time

- radiance field rendering. *ACM Transactions on Graphics*, 42 (4), July 2023. URL <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>. 2, 3, 4
- [10] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2d gaussian splatting for geometrically accurate radiance fields. In *SIGGRAPH 2024 Conference Papers*. Association for Computing Machinery, 2024. doi: 10.1145/3641519.3657428. 2, 4
- [11] Volker Müller. Polarization-based separation of diffuse and specular surface-reflection. In *DAGM-Symposium*, 1995. URL <https://api.semanticscholar.org/CorpusID:51804148>. 2
- [12] S. Rahmann and N. Canterakis. Reconstruction of specular surfaces using polarization imaging. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, volume 1, pages I–I, 2001. doi: 10.1109/CVPR.2001.990468.
- [13] L.B. Wolff and T.E. Boult. Constraining object features using a polarization reflectance model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13(7):635–657, 1991. doi: 10.1109/34.85655. 2
- [14] Paul Debevec, Tim Hawkins, Chris Tchou, Haarm-Pieter Duiker, Westley Sarokin, and Mark Sagar. Acquiring the reflectance field of a human face. In *Proceedings of the 27th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '00, page 145–156, USA, 2000. ACM Press/Addison-Wesley Publishing Co. ISBN 1581132085. doi: 10.1145/344779.344855. URL <https://doi.org/10.1145/344779.344855>. 2
- [15] Abhijeet Ghosh, Graham Fyffe, Borom Tunwattanapong, Jay Busch, Xueming Yu, and Paul Debevec. Multiview face capture using polarized spherical gradient illumination. *ACM Trans. Graph.*, 30(6):1–10, December 2011. ISSN 0730-0301. doi: 10.1145/2070781.2024163. URL <https://doi.org/10.1145/2070781.2024163>. 3
- [16] Cyrus A. Wilson, Abhijeet Ghosh, Pieter Peers, Jen-Yuan Chiang, Jay Busch, and Paul Debevec. Temporal upsampling of performance geometry using photometric alignment. *ACM Trans. Graph.*, 29(2), April 2010. ISSN 0730-0301. doi: 10.1145/1731047.1731055. URL <https://doi.org/10.1145/1731047.1731055>.
- [17] Robert J. Woodham. Photometric method for determining surface orientation from multiple images. 1980. URL <https://api.semanticscholar.org/CorpusID:61691075>. 3
- [18] Paulo Gotardo, Jérémy Riviere, Derek Bradley, Abhijeet Ghosh, and Thabo Beeler. Practical dynamic facial appearance modeling and acquisition. 37(6), December 2018. ISSN 0730-0301. doi: 10.1145/3272127.3275073. URL <https://doi.org/10.1145/3272127.3275073>. 3
- [19] Jérémy Riviere, Paulo Gotardo, Derek Bradley, Abhijeet Ghosh, and Thabo Beeler. Single-shot high-quality facial geometry and skin appearance capture. *ACM Trans. Graph.*, 39(4), August 2020. ISSN 0730-0301. doi: 10.1145/3386569.3392464. URL <https://doi.org/10.1145/3386569.3392464>. 3
- [20] Yingyan Xu, Jérémy Riviere, Gaspard Zoss, Prashanth Chandran, Derek Bradley, and Paulo Gotardo. Improved Lighting Models for Facial Appearance Capture. In Nuria Pelechano and David Vanderhaeghe, editors, *Eurographics 2022 - Short Papers*. The Eurographics Association, 2022. ISBN 978-3-03868-169-4. doi: 10.2312/egs.20221019. 3
- [21] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '99, page 187–194, USA, 1999. ACM Press/Addison-Wesley Publishing Co. ISBN 0201485605. doi: 10.1145/311535.311556. URL <https://doi.org/10.1145/311535.311556>. 3
- [22] Frédéric H. Pighin, Jamie Hecker, Dani Lischinski, Richard Szeliski, and D. Salesin. Synthesizing realistic facial expressions from photographs. *Proceedings of the 25th annual conference on Computer graphics and interactive techniques*, 1998. URL <https://api.semanticscholar.org/CorpusID:74926>.
- [23] Chen Cao, Derek Bradley, Kun Zhou, and Thabo Beeler. Real-time high-fidelity facial performance capture. *ACM Transactions on Graphics (TOG)*, 34:1 – 9, 2015. URL <https://api.semanticscholar.org/CorpusID:207226489>.
- [24] Daniel Vlasic, Matthew Brand, Hanspeter Pfister, and Jovan Popović. Multilinear models for face synthesis. 01 2004. doi: 10.1145/1186223.1186293.
- [25] Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J. Black. Generating 3d faces using convolutional mesh autoencoders. *CoRR*, abs/1807.10267, 2018. URL <http://arxiv.org/abs/1807.10267>.
- [26] Luan Tran and Xiaoming Liu. Nonlinear 3d face morphable model. *CoRR*, abs/1804.03786, 2018. URL <http://arxiv.org/abs/1804.03786>. 3
- [27] Tobias Kirschstein, Shenhan Qian, Simon Giebenhain, Tim Walter, and Matthias Nießner. Nersemble: Multi-view radiance field reconstruction of human heads. *ACM Trans. Graph.*, 42(4), jul 2023. ISSN 0730-0301. doi: 10.1145/3592455. URL <https://doi.org/10.1145/3592455>. 3
- [28] Xiaochen Zhao, Lizhen Wang, Jingxiang Sun, Hongwen Zhang, Jinli Suo, and Yebin Liu. Havatar: High-fidelity head avatar via facial model conditioned neural radiance field. *ACM Trans. Graph.*, oct 2023. ISSN 0730-0301. doi: 10.1145/3626316. URL <https://doi.org/10.1145/3626316>. Just Accepted.
- [29] Tobias Kirschstein, Simon Giebenhain, and Matthias Nießner. Diffusionavatars: Deferred diffusion for high-fidelity 3d head avatars. *arXiv preprint arXiv:2311.18635*, 2023.
- [30] Shenhan Qian, Tobias Kirschstein, Liam Schoneveld, Davide Davoli, Simon Giebenhain, and Matthias Nießner. Gaussianavatars: Photorealistic head avatars with rigged 3d gaussians. *arXiv preprint arXiv:2312.02069*, 2023.
- [31] Simon Giebenhain, Tobias Kirschstein, Martin Rünz, Lourdes Agapito, and Matthias Nießner. Npga: Neural parametric gaussian avatars. In *SIGGRAPH Asia 2024 Conference Papers (SA Conference Papers '24)*, December 3-6,

- Tokyo, Japan, 2024. ISBN 979-8-4007-1131-2/24/12. doi: 10.1145/3680528.3687689.
- [32] 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 (CVPR)*, 2024. 3
 - [33] Shunsuke Saito, Gabriel Schwartz, Tomas Simon, Junxuan Li, and Giljoo Nam. Relightable gaussian codec avatars. In *CVPR*, 2024. 3
 - [34] Pranav Manu, Astitva Srivastava, and Avinash Sharma. Clip-head: Text-guided generation of textured neural parametric 3d head models. In *SIGGRAPH Asia 2023 Technical Communications*, SA '23, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400703140. doi: 10.1145/3610543.3626169. URL <https://doi.org/10.1145/3610543.3626169>. 3, 7
 - [35] Astitva Srivastava, Pranav Manu, Amit Raj, Varun Jampani, and Avinash Sharma. Wordrobe: Text-guided generation of textured 3d garments, 2024. URL <https://arxiv.org/abs/2403.17541>.
 - [36] Dave Zhenyu Chen, Yawar Siddiqui, Hsin-Ying Lee, Sergey Tulyakov, and Matthias Nießner. Text2tex: Text-driven texture synthesis via diffusion models. *arXiv preprint arXiv:2303.11396*, 2023. 3
 - [37] Ayaan Haque, Matthew Tancik, Alexei A. Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions, 2023. URL <https://arxiv.org/abs/2303.12789>. 3
 - [38] Mohit Mendiratta, Xingang Pan, Mohamed Elgharib, Kartik Teotia, Mallikarjun B R, Ayush Tewari, Vladislav Golyanik, Adam Kortylewski, and Christian Theobalt. Avatarstudio: Text-driven editing of 3d dynamic human head avatars, 2023. URL <https://arxiv.org/abs/2306.00547>. 3
 - [39] Cong Wang, Di Kang, He-Yi Sun, Shen-Han Qian, Zi-Xuan Wang, Linchao Bao, and Song-Hai Zhang. Mega: Hybrid mesh-gaussian head avatar for high-fidelity rendering and head editing. *arXiv preprint arXiv:2404.19026*, 2024. 3
 - [40] Simon Giebenhain, Tobias Kirschstein, Markos Georgopoulos, Martin Rünz, Lourdes Agapito, and Matthias Nießner. Monophm: Dynamic head reconstruction from monocular videos. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2024. 3
 - [41] 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 (CVPR)*, 2024. 3
 - [42] Jun Xiang, Xuan Gao, Yudong Guo, and Juyong Zhang. Flashavatar: High-fidelity head avatar with efficient gaussian embedding, 2024. URL <https://arxiv.org/abs/2312.02214>. 3
 - [43] Dejan Azinović, Olivier Maury, Christophe Hery, Matthias Nießner, and Justus Thies. High-res facial appearance capture from polarized smartphone images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023. 4, 5, 6
 - [44] Shenhan Qian. Versatile head alignment with adaptive appearance priors. September 2024. URL <https://github.com/ShenhanQian/VHAP>. 4, 7
 - [45] 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. 4
 - [46] 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 (CVPR)*, pages 5501–5510, June 2022. 5
 - [47] Fred E. Nicodemus. Directional reflectance and emissivity of an opaque surface. *Applied Optics*, 4(7):767, July 1965. doi: 10.1364/AO.4.000767. 5
 - [48] R. L. Cook and K. E. Torrance. A reflectance model for computer graphics. *ACM Trans. Graph.*, 1(1):7–24, January 1982. ISSN 0730-0301. doi: 10.1145/357290.357293. URL <https://doi.org/10.1145/357290.357293>. 5
 - [49] Christophe Schlick. An inexpensive brdf model for physically-based rendering. *Computer Graphics Forum*, 13(3):233–246, 1994. doi: <https://doi.org/10.1111/1467-8659.1330233>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/1467-8659.1330233>. 6
 - [50] Bruce Walter, Stephen R. Marschner, Hongsong Li, and Kenneth E. Torrance. Microfacet models for refraction through rough surfaces. In *Proceedings of the 18th Eurographics Conference on Rendering Techniques*, EGSR'07, page 195–206, Goslar, DEU, 2007. Eurographics Association. ISBN 9783905673524. 6
 - [51] B. Smith. Geometrical shadowing of a random rough surface. *IEEE Transactions on Antennas and Propagation*, 15(5):668–671, 1967. doi: 10.1109/TAP.1967.1138991. 6
 - [52] Michael Ashikhmin and Peter Shirley. An anisotropic phong light reflection model. *Journal of Graphics Tools*, 5, 01 2001. 6
 - [53] Jacob Munkberg, Jon Hasselgren, Tianchang Shen, Jun Gao, Wenzheng Chen, Alex Evans, Thomas Müller, and Sanja Fidler. Extracting Triangular 3D Models, Materials, and Lighting From Images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8280–8290, June 2022. 6
 - [54] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022. doi: 10.1145/3528223.3530127. URL <https://doi.org/10.1145/3528223.3530127>. 6, 7, 9
 - [55] Yuehao Wang, Chaoyi Wang, Bingchen Gong, and Tianfan Xue. Bilateral guided radiance field processing, 2024. URL <https://arxiv.org/abs/2406.00448>. 6
 - [56] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>. 7
 - [57] Fernando A. Fardo, Victor H. Conforto, Francisco C. de Oliveira, and Paulo S. Rodrigues. A formal evaluation of

psnr as quality measurement parameter for image segmentation algorithms, 2016. URL <https://arxiv.org/abs/1605.07116>. 7

- [58] Alain Horé and Djemel Ziou. Image quality metrics: Psnr vs. ssim. In *2010 20th International Conference on Pattern Recognition*, pages 2366–2369, 2010. doi: 10.1109/ICPR.2010.579. 7