

Drivable 3D Gaussian Avatars

Wojciech Zielonka^{3,1*}, Timur Bagautdinov¹, Shunsuke Saito¹,
Michael Zollhöfer¹, Justus Thies^{2,3}, Javier Romero¹

¹Meta Reality Labs Research ²Technical University of Darmstadt

³Max Planck Institute for Intelligent Systems, Tübingen, Germany

<https://zielon.github.io/d3ga/>



Figure 1. Given a multi-view video, D3GA learns drivable photo-realistic 3D human avatars, represented as a composition of 3D Gaussians embedded in tetrahedral cages. The Gaussians are transformed by those cages, colored with an MLP, and rasterized as splats. We represent the drivable human as a layered set of 3D Gaussians, allowing us to decompose the avatar into its different cloth layers.

Abstract

We present *Drivable 3D Gaussian Avatars (D3GA)*, the first 3D controllable model for human bodies rendered with Gaussian splats. Current photorealistic drivable avatars require either accurate 3D registrations during training, dense input images during testing, or both. The ones based on neural radiance fields also tend to be prohibitively slow for telepresence applications. This work uses the recently presented 3D Gaussian Splatting (3DGS) technique to render realistic humans at real-time framerates, using dense calibrated multi-view videos as input. To deform those primitives, we depart from the commonly used point deformation method of linear blend skinning (LBS) and use a classic volumetric deformation method: cage deformations. Given their smaller size, we drive these deformations with joint angles and keypoints, which are more suitable for communication applications. Our experiments on nine subjects with varied body shapes, clothes, and motions obtain higher-quality results than state-of-the-art methods when using the same training and test data.

*Work done while Wojciech Zielonka was an intern at Reality Labs Research, Pittsburgh, PA, USA

1. Introduction

In the nineteenth century, *the Anonymous Society of Painters, Sculptors, Printmakers, etc.* started the art movement called Impressionism, identified by a technique of “short, broken brushstrokes that barely convey forms”. Our goal, to create photorealistic representations of humans, is one of the things that impressionists ran away from. However, in D3GA¹, we use Gaussian splats as a modern version of those short brushstrokes to conform to the structure and appearance of our real-time, reusable avatars.

Creating drivable (i.e., that can be animated to generate new content) photorealistic humans currently requires dense multi-view data since monocular approaches lack accuracy. Additionally, existing techniques rely on complex pre-processing, including precise 3D registrations [1, 55, 56]. However, obtaining those registrations requires iterative methods that are difficult to integrate into end-to-end pipelines. Other methods that do not require accurate registrations [4] are based on neural radiance fields (NeRFs). They are typically too slow for real-time rendering (with few exceptions [30]) or struggle with garment animations.

¹Referring to Edgar Degas (pronunciation: ed-gr duh-gaa), a French impressionist artist known for his pastel drawings and oil paintings.

In recent work, Kerbl et al. introduced 3D Gaussian Splatting (3DGS) [14] based on the classic rendering approach Surface Splatting [66]. This representation renders higher-quality images at a faster framerate than state-of-the-art methods based on neural radiance fields [30] and does not require any highly accurate 3D initialization.

Unfortunately, 3DGS was designed for static scenes. Time-conditioned Gaussian Splatting [54, 59] are proposed for rendering dynamic scenes. However, similar to [7], these models can only *replay* previously observed content, making them unsuitable for representing novel motion.

Following approaches on drivable NeRFs [53, 54], we model the 3D human appearance and deformations in a canonical space but rely on 3D Gaussians instead of radiance fields. In addition to better performance, Gaussian splats do not require camera ray sampling heuristics.

Drivable NeRFs typically rely on LBS to transform points between canonical and observation spaces. However, D3GA models humans with volumetric primitives in the form of 3D Gaussians and therefore needs to map *volumes* to canonical space. Instead of LBS, our method uses another classic deformation model suitable for transforming volumes: cages [31]. Deforming cages in canonical space entails a deformation gradient, which can be applied directly to the 3D Gaussians in our representation. Our method follows a compositional structure based on separate body, face, and garment cages, allowing us to render those parts independently.

The remaining question is defining the signal that triggers those cage deformations. The current state-of-the-art in drivable avatars [41, 57] requires dense input signals like RGB-D images or even multi-camera setups, which might not be suitable for low-bandwidth connections in telepresence applications. We adopt a more compact input based on the human pose, comprising skeletal joint angles in the form of quaternions and 3D facial keypoints.

We train person-specific models on nine high-quality multi-view sequences with a wide range of body shapes, motion, and clothing (not limited to tight-fitting), which later can be driven with new poses from any subject.

In summary, we present Drivable 3D Gaussian Avatars (D3GA) with the following contributions:

- The first implementation of Drivable 3D Gaussians Splatting (3DGS), applied to digital avatars.
- Tetrahedral cage-based deformations applied to 3DGS.
- State-of-the-art pose-based avatar generation for dense multi-view scenes without ground truth registration.

2. Related Work

Current methods for controllable avatars are primarily based on dynamic Neural Radiance Fields (NeRF) [28,

33, 34], point-based [25, 58, 62] or hybrid representations [1, 4, 22, 64] which are either slow to render or do not correctly disentangle garment from the body; thus, incorrectly generalize to new poses. For a thorough overview, we point the reader to state-of-the-art reports [48, 49, 65] on digital avatars and neural rendering.

Dynamic Neural Radiance Fields NeRF [29] is a prevalent appearance model for human avatars. It represents the scene volumetrically by storing density and color information in space using a multi-layer perceptron (MLP). Images can be rendered from this representation by using ray casting with volumetric integration of the sample points [13]. Many methods successfully applied NeRF to dynamic scenes [5, 20, 33, 34, 37, 53, 58, 64] achieving high quality results. However, most of the methods treat avatars as a single layer [19, 28, 35, 43–45, 63], meaning there is no separation between garment and body. This is particularly cumbersome for modeling phenomena like sliding or loose garments. Methods like [3, 4] try to solve this problem using a hybrid representation. They combined explicit geometry from SMPL [23] and implicit dynamic NeRF by integrating the mesh surface into the rendering equation. Despite impressive garment reconstruction, these methods struggle with novel pose prediction. TECA [60] lifts SCARF to a generative framework that enables prompt-based generation of NeRF-based accessories and hairstyles.

Point Based Rendering Before 3DGS, many methods used point-based rendering [25, 44, 62] or sphere splatting [18], which, similarly to 3DGS, have optimizable positions and sizes. NPC by Su et al. [44] defines a point-based body model for avatar representation. Their model requires the evaluation of the nearest neighbor search per ray sample during training, which results in long training times (12h instead of 30 minutes for our model on a similar dataset size), making it impractical for dense multi-view datasets. Ma et al. [25] represent garment as a pose-dependent function that maps a set of points from SMPL [23] to the clothing space. This idea is improved in [38], where a neural deformation field replaces LBS. However, both models restrict themselves to model only geometry, not appearance. Zheng et al. [62] represent the upper part of an avatar as a point cloud, which is progressively grown during the optimization and rasterized using a differentiable point cloud renderer [52]. Despite achieving photorealistic results locally, the avatars suffer from artifacts like holes around low-density regions.

Cage Based Deformations Cages [31] are commonly used for geometry modeling and animation. They serve as a sparse proxy that controls all points in their interior, enabling efficient deformation since only cage nodes have to be controlled to rig the object inside. Yifan et al. [51] introduced the concept of neural cages for detail-preserving shape deformation. The network learns how to rig the

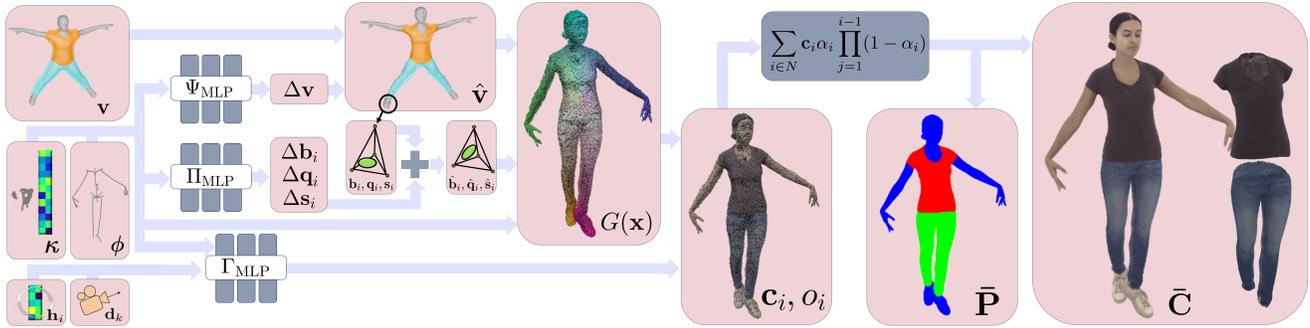


Figure 2. **Overview.** D3GA uses 3D pose ϕ , face embedding κ , viewpoint \mathbf{d}_k and canonical cage \mathbf{v} (as well as auto-decoded color features \mathbf{h}_i) to generate the final render $\bar{\mathbf{C}}$ and auxiliary segmentation render $\bar{\mathbf{P}}$. The inputs in the left are processed through three networks (Ψ_{MLP} , Π_{MLP} , Γ_{MLP}) per avatar part to generate cage displacements $\Delta\mathbf{v}$, Gaussians deformations $\mathbf{b}_i, \mathbf{q}_i, \mathbf{s}_i$ and color/opacity \mathbf{c}_i, o_i respectively. After cage deformations transform canonical Gaussians, they are rasterized into the final images according to Eq. 9.

source object into the target through a proxy regressed by a neural network. Garbin et al. [6] extended dynamic NeRF with tetrahedron cages to facilitate the unposing of ray samples based on tetrahedron intersections. The method is real-time, high-quality, and controllable. However, their results are limited to objects with local deformations like heads, making them not applicable to highly articulatable objects like our full-body avatars. Peng et al. also used a cage to deform a radiance field in CageNeRF [36]. Their low-resolution cages can be applied to full-body avatars but fail to model details like faces or other complex deformations.

Gaussian Splatting D3GA is based on 3D Gaussians Splatting (3DGS) [14], a recent alternative approach to NeRF that achieves high quality and real-time rendering speed. 3DGS is based on 3D Gaussians, a differentiable volumetric representation that can be efficiently rasterized in comparison to expensive ray marching used by NeRF. The recently introduced Dynamic 3DGS [24] enables per-frame dense 6-DOF tracking and novel view synthesis by optimizing the trajectories of the 3D Gaussians from frame t_i to t_{i+1} . Our method extends static or playback 3DGS [14, 24] to drivable applications using a volumetric cage as a deformation proxy, enabling controlling digital avatars.

3. Method

Current methods for dynamic volumetric avatars either map points from deformed to canonical space [9, 10, 21, 39, 39, 64] or they rely on the forward mapping only [2, 19, 35, 44, 50, 62]. Methods based on backward mapping tend to accumulate errors in canonical space since they require an error-prone backward pass and have problems modeling view-dependent effects since mapping the view vector to canonical space uniquely is non-trivial. Therefore, we decided to employ a forward-only mapping. D3GA is built on 3DGS extended by a neural representation and tetrahedral cages to model the color and geometry of each dynamic part of the avatar, respectively. In the following,

we introduce the formulation of 3D Gaussian Splatting and give a detailed description of our method D3GA.

3.1. 3D Gaussian Splatting

3D Gaussian Splatting (3DGS) [14] is designed for real-time novel view synthesis in multi-view static scenes. Their rendering primitives are scaled 3D Gaussians [17, 52] with a 3D covariance matrix Σ and mean μ :

$$G(\mathbf{x}) = e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}. \quad (1)$$

To splat the Gaussians, Zwicker et al. [66] define the projection of 3D Gaussians onto the image plane as:

$$\Sigma' = \mathbf{A} \mathbf{W} \Sigma \mathbf{W}^T \mathbf{A}^T, \quad (2)$$

where Σ' is a covariance matrix in 2D space, \mathbf{W} is the view transformation, and \mathbf{A} is the Jacobian of the affine approximation of the projective transformation. During optimization, enforcing the positive semi-definiteness of the covariance matrix Σ is challenging. To avoid this, Kerbl et al. [14] use an equivalent formulation of a 3D Gaussian as a 3D ellipsoid parameterized with a scale \mathbf{S} and rotation \mathbf{R} :

$$\Sigma = \mathbf{R} \mathbf{S} \mathbf{S}^T \mathbf{R}^T. \quad (3)$$

3DGS uses spherical harmonics [40] to model the view-dependent color of each Gaussian. In practice, appearance is modeled with an optimizable 48 elements vector representing four bands of spherical harmonics.

3.2. Cage Based Deformation Transfer

To deform 3D Gaussians, we leverage tetrahedron cage-based deformations as a coarse proxy for the body, face, and individual garments. To create a cage per garment, we segment all images of a single time instance using an EfficientNet [47] backbone with PointRend [16] refinement, trained on a corpus of similar multi-view captures. The per-image 2D segmentation masks are projected onto a body mesh \mathbf{M}

to obtain per-triangle labels (body, upper, lower). To get the mesh $\hat{\mathbf{M}}$, we fit a low-resolution LBS model to a single 3D scan of the subject and then fit such model to the segmented frame by minimizing the distance to the 3D keypoints, extracted with an EfficientNet trained on similar captures. We transform the body mesh into canonical space with LBS and divide it into body part templates \mathbf{M}_k . The garment meshes are additionally inflated 3cm along the vertex normals. After that, we use TetGen [42] to turn the unposed meshes \mathbf{M}_k into tetrahedral meshes \mathbf{T}_k . Consequently, cages for garments are hollow, containing only their outer layer, while the body cage is solid. The face cage is composed of the body tetrahedra which contains triangles defined as the face on the LBS template. The cage nodes are deformed according to LBS weights transferred from the closest vertex in \mathbf{M}_k . While classic cage methods typically deform the volume according to complex weight definitions [8, 11, 12], using linear weights works well in practice when cage cells are small, making it easier to integrate into an end-to-end training system. Specifically, if we define \mathbf{v}_{ij} as the vertices of tetrahedron i in canonical space, any point \mathbf{x} inside this tetrahedron can be defined in terms of its barycentric coordinates b_j :

$$\mathbf{x} = \sum_{j=1}^4 b_j \mathbf{v}_{ij}. \quad (4)$$

When tetrahedra are transformed to posed space according to $\hat{\mathbf{v}}_{ij} = \text{LBS}(\mathbf{v}_{ij}, \phi, \mathbf{w}_{ij})$, where ϕ is the pose and \mathbf{w}_{ij} the blendweights, the same linear relation holds $\hat{\mathbf{x}} = \sum_{j=1}^4 b_j \hat{\mathbf{v}}_{ij}$. To extend this transformation from points to volumes, we use the deformation transfer [46] as:

$$\mathbf{J}_i \mathbf{E}_i = \hat{\mathbf{E}}_i, \quad (5)$$

$$\mathbf{J}_i = \hat{\mathbf{E}}_i \mathbf{E}_i^{-1}, \quad (6)$$

where $\hat{\mathbf{E}}_i \in \mathbb{R}^{3 \times 3}$ and $\mathbf{E}_i \in \mathbb{R}^{3 \times 3}$ contain three edges from tetrahedron i defined in deformed and canonical spaces, respectively. In the following subsection, we will explore how to use the deformation gradients closed form solution in Eq. 6 for transforming 3D Gaussians.

3.3. Drivable Gaussian Avatars

We initialize a fixed number of Gaussians, whose 3D means μ are sampled on the surface of $\hat{\mathbf{M}}$. The rotation of each Gaussian is initialized so that the first two axes are aligned with the triangle surface and the third one with the normal: this is a good approximation for a smooth surface. The scale is initialized uniformly across a heuristic range depending on inter-point distances as in [14]. Finally, we can assign each sampled position \mathbf{x} to the intersecting tetrahedron and compute its barycentric coordinates $\mathbf{b} \in \mathbb{R}^4$. To deform the tetrahedron volume, we incorporate the deformation gradient \mathbf{J} defined in Eq. 6 into the Gaussian covari-

ance matrix from Eq. 3. The final covariance matrix passed to the rasterizer is denoted as:

$$\hat{\Sigma} = \mathbf{J}_i \Sigma \mathbf{J}_i^T, \quad (7)$$

where \mathbf{J}_i is the deformation gradient of the tetrahedron containing the 3D mean of the Gaussian with covariance Σ . This way, we transfer the deformation into the Gaussians, improving modeling phenomena like garment stretching.

Each part of the avatar (the garment, body, or face) is controlled by a separate GaussianNet $\mathbb{G}_{\text{Net}} = \{\Gamma_{\text{MLP}}, \Pi_{\text{MLP}}, \Psi_{\text{MLP}}\}$ which is defined as a set of small specialized multi-layer perceptrons (MLP) parametrized as:

$$\begin{aligned} \Psi_{\text{MLP}} &: \{\phi, \text{enc}_{\text{pos}}(\mathbf{v})\} \rightarrow \Delta \mathbf{v}, \\ \Pi_{\text{MLP}} &: \{\phi, \mathbf{b}_i, \mathbf{q}_i, \mathbf{s}_i\} \rightarrow \{\Delta \mathbf{b}_i, \Delta \mathbf{s}_i, \Delta \mathbf{q}_i\}, \\ \Gamma_{\text{MLP}} &: \{\phi, \text{enc}_{\text{view}}(\mathbf{d}_k), \mathbf{h}_i, \mathbf{f}_j\} \rightarrow \{\mathbf{c}_i, o_i\}. \end{aligned} \quad (8)$$

All the networks take joint angles ϕ (or face encodings κ for the face networks) as inputs, in addition to network-specific conditioning. The cage node correction network Ψ_{MLP} takes positional encodings [29] for all canonical vertices to transform them into offsets of the cage node positions similar to SMPL [23] pose-correctives. To adapt our representation further to the pose, the Gaussian correction network Π_{MLP} takes additionally the canonical Gaussian parameters (barycentric coordinates $\mathbf{b}_i \in \mathbb{R}^4$, rotation $\mathbf{q}_i \in \mathbb{R}^4$ and scale $\mathbf{s}_i \in \mathbb{R}^3$) to predict corrections of those same parameters. These two networks are necessary to capture high-frequency details outside the parametric transformation.

In terms of appearance, the shading network Γ_{MLP} transforms information about the encoded view direction and initial color into final color and opacity \mathbf{c}_i, o_i . We depart from 3DGS color representation based on Spherical Harmonics to enable pose-dependent color, which is necessary to model self-shadows and wrinkles in garments. The view angle information is represented as its projection to the first four spherical harmonics bands $\text{enc}_{\text{pos}}(\cdot)$. At the same time, the initial color is an auto-decoded [32] feature vector $\mathbf{h}_i \in \mathbb{R}^{48}$. Moreover, the face region uses as input face embeddings κ instead of pose ϕ . A small auxiliary MLP regresses κ based on 150 3D keypoints \mathbf{k} normalized by their training mean and standard deviations. Finally, we also add an embedding vector with the timeframe of the current sample [27]. This allows D3GA to explain away properties that cannot be modeled (e.g., cloth dynamics) from our input, effectively avoiding excessive blur due to averaging residuals. During testing, the average training embedding is used.

3.4. Training Objectives

As in 3DGS [14], we define the color $\bar{\mathbf{C}}$ of pixel (u, v) :

$$\bar{\mathbf{C}}_{u,v} = \sum_{i \in \mathcal{N}} \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (9)$$



Figure 3. D3GA ablation: shape errors without cage deformations, view-dependent color artifacts with SH, shape smoothness without $\mathcal{L}_{Garment}$, or sliding artifacts with single layer. No Neo-Hookean loss results in reasonable color, but degenerate cage geometry (Fig.7).

where \mathbf{c}_i is the color predicted by Γ_{MLP} , which replaces the spherical harmonics in 3DGS. α_i is computed as the product of the Gaussian density in Eq. 1 with covariance matrix Σ' from Eq. 2 and the learned per-point opacity o_i predicted by Γ_{MLP} . The sum is computed over set \mathcal{N} , the Gaussians with spatial support on (u, v) . The primary loss in D3GA is a weighted sum of three different color losses applied to the estimated image $\bar{\mathbf{C}}$ and the ground truth RGB image \mathbf{C} :

$$\mathcal{L}_{Color} = (1 - \omega)\mathcal{L}_1 + \omega\mathcal{L}_{D-SSIM} + \zeta\mathcal{L}_{VGG}, \quad (10)$$

where $\omega = 0.2$, $\zeta = 0.005$ (after 400k iterations steps and zero otherwise), \mathcal{L}_{D-SSIM} is a structural dissimilarity loss, and \mathcal{L}_{VGG} is the perceptual VGG loss.

To encourage correct garment separation, we introduce a garment loss. Since each Gaussian i is statically assigned to a part, we define \mathbf{p}_i as a constant-per-part color and consequently render $\bar{\mathbf{P}}$ by replacing \mathbf{c}_i by \mathbf{p}_i in Eq. 9. Then, we compute the \mathcal{L}_1 norm between predicted parts $\bar{\mathbf{P}}$ and ground truth segmentations \mathbf{P} , $\mathcal{L}_{Garment} = \mathcal{L}_1(\bar{\mathbf{P}}, \mathbf{P})$. Moreover, we are using the Neo-Hookean loss based on Macklin et al. [26] to enforce the regularization of the predicted tetrahedra for the regions with low supervision signal:

$$\mathcal{L}_{Neo} = \frac{1}{N} \sum_{i=0}^N \frac{\lambda}{2} (\det(\mathbf{J}_i) - 1)^2 + \frac{\mu}{2} (\text{tr}(\mathbf{J}_i^T \mathbf{J}_i) - 3), \quad (11)$$

where \mathbf{J}_i denotes the deformation gradient between a canonical and a deformed tetrahedron (Eq. 6), N is the total

number of tetrahedrons, and λ and μ are the Lamé parameters [26]. The overall loss is defined as:

$$\mathcal{L} = \nu\mathcal{L}_{Color} + \nu\mathcal{L}_{Garment} + \tau\mathcal{L}_{Neo}, \quad (12)$$

where $\nu = 10$ and $\tau = 0.005$ balance the different losses.

We implemented D3GA based on the differentiable 3DGS renderer [14]. The networks Π_{MLP} , Ψ_{MLP} , Γ_{MLP} have three hidden layers with 128 neurons and ReLU activation functions. In our experiments, we train the networks for 700k steps with a multi-step scheduler with a decay rate of 0.33, a batch size of one, and using the Adam optimizer [15] with a learning rate set to $5e - 4$. We ran all experiments on a single Nvidia V100 GPU with 1024×667 images.

4. Dataset

Our dataset consists of nine subjects performing different motions, observed by 200 cameras. We use 12000 frames for training (at 10 FPS) and 1500 for testing (at 30 FPS). The images were captured in a multi-view studio with synchronized cameras at a resolution of 4096×2668 , but they were downsampled to 1024×667 to reduce the computational cost. We use 2D segmentation masks, RGB images, keypoints, and 3D joint angles for training as well as a single registered mesh to create our template $\hat{\mathbf{M}}$.



Figure 4. Qualitative comparisons show that D3GA models garments better than other SOTA approaches, especially loose ones like skirts or sweatpants. FFD stands for free form deformation meshes, which contain a much richer training signal than LBS meshes (see Fig. 9).



Figure 5. Our method allows multilayer garment decompositions.

5. Results

We evaluate and compare our method w.r.t. state-of-the-art multiview-based methods [1, 22, 41]. We compare D3GA to the full-body avatar methods BodyDecoder (BD) [1], MVP-based avatar [22, 41], and DVA [41], which uses dense image conditioning from all cameras. For a fair comparison, we used two types of geometry training input for BD and MVP (see Fig. 9): meshes with a simple LBS model tracked with body keypoints, and detailed registration meshes where vertices are optimized freely to match the 3D reconstruction of each frame (also called free-form deformation, FFD). Note that BodyDecoder also uses an ambient occlusion approximation [1] extracted from FFD meshes.

5.1. Image Quality Evaluation

D3GA is evaluated w.r.t. SSIM, PSNR, and the perceptual metric LPIPS [61]. Table 1 shows that our method is the one that achieves the best PSNR and SSIM among the methods using only LBS (i.e., do not require 3D scans for every frame) and outperforms all FFD methods minus BD FFD in these metrics, despite having poorer training signal and no test images (DVA was tested using all 200 cameras). Moreover, our approach allows us to decompose avatars into drivable garment layers compared to other volumetric methods. Figure 5 shows each separate garment layer, which can be controlled solely by skeleton joint angles, without requiring specific garment registration modules as in [56].

Method	PSNR \uparrow	LPIPS \downarrow	SSIM \uparrow
Ours	30.634	0.054	0.965
MVP LBS [22]	28.795	0.051	0.955
BD LBS [1]	29.919	0.044	0.960
BD FFD [1]	30.999	0.039	0.964
MVP FFD [22]	30.072	0.043	0.960
DVA [41]	30.239	0.042	0.963

Table 1. Our method scores the best in terms of PSNR and SSIM for LBS-based methods. However, it lacks the sharpness of the mesh-based method. Moreover, our method outperforms MVP [22], which uses FFD meshes, scoring second in total for PSNR error and the best in SSIM. ■ First, ■ second, ■ third place.



Figure 6. The additional supervision $\mathcal{L}_{Garment}$ improves the garment’s shape by reducing semitransparent effects at the boundary.

Experiment	PSNR \uparrow	LPIPS \downarrow	SSIM \uparrow
Ours	29.825	0.056	0.960
w/o cage	28.279	0.066	0.955
w/ SH	28.960	0.058	0.957
w/o $\mathcal{L}_{Garment}$	30.140	0.057	0.961
w/o \mathcal{L}_{Neo}	29.911	0.056	0.960
Single layer	29.740	0.057	0.959

Table 2. Lack of cage proxy significantly increases reconstruction error. Moreover, single-layer avatars incorrectly model sliding garments, and using SH for color modeling struggles with wrinkles and self-shadows. ■ First, ■ second, ■ third place.

5.2. Ablation Studies

In this section, the influence of different pipeline components (deformation and appearance representations, number of Gaussians) and losses ($\mathcal{L}_{Garment}$, \mathcal{L}_{Neo}) are analyzed.

Cage based deformation To evaluate the relevance of deforming Gaussians with cages, we performed an experiment where Gaussian positions were transformed with LBS directly. When initializing the Gaussian positions, instead of assigning them to a tetrahedral mesh proxy, we use LBS to transform their 3D mean according to the closest point in $\hat{\mathbf{M}}$. The network Π_{MLP} predict updates to the rotation, scale and Gaussian mean $\Delta \mathbf{q}_i, \Delta s_i, \Delta \mu_i$, while network Ψ_{MLP} is disabled. The results presented in the third column of Fig. 3 show significant artifacts, especially for the highly dynamic parts of the avatar. One of the reasons is the lack of stretching, which is implicitly provided by the deformation gradient from the cage.

Garment loss The garment loss $\mathcal{L}_{Garment}$ serves two primary purposes: it improves garment separation and reduces erroneously translucent regions. Figure 6 depicts the effect of the loss on a t-shirt. It can be seen that quality degrades without it, especially on the edges. Although image metrics PSNR and SSIM are best without $\mathcal{L}_{Garment}$, we can observe qualitatively that regions between garments’ boundary are blurry and have erroneous opacity.

Spherical harmonics appearance 3DGS uses spherical harmonics (SH) to model the color of each Gaussian in the scene. The fourth column of Fig. 3 shows results where

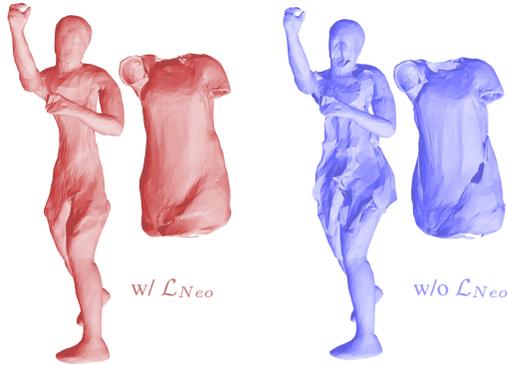


Figure 7. The effect of the tetrahedra regularization \mathcal{L}_{Neo} is mostly visible in the regions which lack supervision or undergo sliding, which covers them for most of the time.

Γ_{MLP} is replaced by the SH layer in 3DGS. As can be seen, it struggles to capture phenomena like self-shadows and wrinkles, which are pose-dependent.

Tetrahedral regularization We introduced \mathcal{L}_{Neo} (Equation 11) to avoid geometry artifacts that could potentially misplace the Gaussians. It prevents tetrahedra from losing too much volume, flipping, or diverging in size from the canonical shape. Optimization of layered garments will naturally struggle for regions that are either permanently or temporarily covered, resulting in geometric artifacts, which can be alleviated by \mathcal{L}_{Neo} regularization (Fig. 7).

Single layer avatar D3GA supports a single-layer training for the garment and body, which struggles to model proper garment sliding. The results are presented in the last column of Fig. 3. It can be observed that the edges between the T-shirt and jeans are oversmoothed.

Number of Gaussians As shown in Table 3, the runtime of D3GA depends on the number of Gaussians. Generally, the best quality ratio to rendering time is between 25k and 100k Gaussians. We chose 100k for our experiments.

6. Discussion

While D3GA shows better quality and competitive rendering speed w.r.t. the state of the art, there are still particular challenges. High-frequency patterns, like stripes, may result in blurry regions. One way of improving image quality would be using a variation autoencoder to regress Gaussian parameters per texel of a guide mesh similar to [22].



Figure 8. Reposing eight avatars with a pose from another subject

Experiment	PSNR \uparrow	LPIPS \downarrow	SSIM \uparrow	FPS \uparrow	Δ FPS
25k Gaussians	29.938	0.058	0.960	28	107%
100k Gaussians	29.825	0.056	0.960	26	100%
200k Gaussians	29.864	0.056	0.960	23	88%
300k Gaussians	29.864	0.056	0.960	20	77%

Table 3. Average frame rate per second at 1024×667 resolution w.r.t to the amount of Gaussian measured on a Nvidia V100 GPU. 100k Gaussians provide the best rendering-time-to-quality ratio.



Figure 9. Comparison of LBS (in red) and FFD meshes (in blue)

Despite using the $\mathcal{L}_{Garment}$ loss, self-collisions for loose garments are still challenging, and the sparse controlling signal does not contain enough information about complex wrinkles or self-shading. An exciting follow-up work direction would be replacing the appearance model in D3GA with a relightable one. In order to extend D3GA repositability (see Fig. 8) to reshapability and cloth transfer, we would like to upgrade our LBS model to a more general one (e.g. SMPL [23]) and integrate a cloth shape space. Finally, we would like to point out that D3GA is very flexible and can be adapted to specific applications, e.g. by using more Gaussians to capture high-frequency details (to the detriment of rendering speed) or removing garment supervision to reduce PSNR (if good cage geometry decomposition is not needed).

7. Conclusion

We have proposed D3GA, a novel approach for reconstructing animatable human avatars enabled by drivable 3D Gaussians rigged with tetrahedral cages. Our method shows high-quality results, better than the state of the art using similar input and comparable to approaches using richer information like FFD meshes or images at test time. Moreover, our solution shows promising results in geometry and appearance modeling for dynamic sequences without requiring ground truth geometry, thus shortening the data processing time.

Acknowledgement The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting WZ. We also want to thank Giljoo Nam for the help with Gaussian visualizations, Anka Chen for very useful conversations about tetrahedrons, and Shou-I Yu and Robbin Xu for their invaluable help with data processing.

References

- [1] Timur M. Bagautdinov, Chenglei Wu, Tomas Simon, Fabián Prada, Takaaki Shiratori, Shih-En Wei, Weipeng Xu, Yaser Sheikh, and Jason M. Saragih. Driving-signal aware full-body avatars. *ACM Transactions on Graphics (TOG)*, 40:1 – 17, 2021. [1](#), [2](#), [6](#), [7](#)
- [2] Xu Chen, Yufeng Zheng, Michael J. Black, Otmar Hilliges, and Andreas Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11574–11584, 2021. [3](#)
- [3] Yao Feng, Weiyang Liu, Timo Bolkart, Jinlong Yang, Marc Pollefeys, and Michael J. Black. Learning disentangled avatars with hybrid 3d representations. *arXiv*, 2023. [2](#)
- [4] Yao Feng, Jinlong Yang, Marc Pollefeys, Michael J. Black, and Timo Bolkart. Capturing and animation of body and clothing from monocular video. *SIGGRAPH Asia 2022 Conference Papers*, 2022. [1](#), [2](#)
- [5] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8645–8654, 2020. [2](#)
- [6] Stephan J. Garbin, Marek Kowalski, Virginia Estellers, Stanislaw Szymanowicz, Shideh Rezaeifar, Jingjing Shen, Matthew Johnson, and Julien Valentin. Voltmorph: Real-time, controllable and generalisable animation of volumetric representations. *CoRR*, abs/2208.00949, 2022. [3](#)
- [7] Mustafa Isik, Martin Rünz, Markos Georgopoulos, Taras Khakhulin, Jonathan Starck, Lourdes Agapito, and Matthias Nießner. Humanrf: High-fidelity neural radiance fields for humans in motion. *ACM Trans. Graph.*, 42(4):160:1–160:12, 2023. [2](#)
- [8] Alec Jacobson, Ilya Baran, Jovan Popović, and Olga Sorkine-Hornung. Bounded biharmonic weights for real-time deformation. *ACM SIGGRAPH 2011 papers*, 2011. [4](#)
- [9] Timothy Jeruzalski, Boyang Deng, Mohammad Norouzi, J. P. Lewis, Geoffrey E. Hinton, and Andrea Tagliasacchi. Nasa: Neural articulated shape approximation. *ArXiv*, abs/1912.03207, 2019. [3](#)
- [10] Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel, and Anurag Ranjan. Neuman: Neural human radiance field from a single video. In *European Conference on Computer Vision*, 2022. [3](#)
- [11] Pushkar Joshi, Mark Meyer, Tony DeRose, Brian Green, and Tom Sanocki. Harmonic coordinates for character articulation. *ACM Trans. Graph.*, 26(3):71, 2007. [4](#)
- [12] Tao Ju, Scott Schaefer, and Joe D. Warren. Mean value coordinates for closed triangular meshes. *ACM SIGGRAPH 2005 Papers*, 2005. [4](#)
- [13] James T. Kajiya. The rendering equation. *Proceedings of the 13th annual conference on Computer graphics and interactive techniques*, 1986. [2](#)
- [14] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkuehler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (TOG)*, 42:1 – 14, 2023. [2](#), [3](#), [4](#), [5](#)
- [15] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014. [5](#)
- [16] Alexander Kirillov, Yuxin Wu, Kaiming He, and Ross B. Girshick. Pointrend: Image segmentation as rendering. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 9796–9805. Computer Vision Foundation / IEEE, 2020. [3](#)
- [17] Georgios Kopanas, Julien Philip, Thomas Leimkühler, and George Drettakis. Point-based neural rendering with per-view optimization. *Computer Graphics Forum*, 40, 2021. [3](#)
- [18] Christoph Lassner and Michael Zollhöfer. Pulsar: Efficient sphere-based neural rendering. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1440–1449, 2021. [2](#)
- [19] Ruilong Li, Julian Tanke, Minh Vo, Michael Zollhofer, Jurgen Gall, Angjoo Kanazawa, and Christoph Lassner. Tava: Template-free animatable volumetric actors. *ArXiv*, abs/2206.08929, 2022. [2](#), [3](#)
- [20] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6494–6504, 2020. [2](#)
- [21] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor. *ACM Transactions on Graphics (TOG)*, 40:1 – 16, 2021. [3](#)
- [22] Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhofer, Yaser Sheikh, and Jason M. Saragih. Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (TOG)*, 40:1 – 13, 2021. [2](#), [6](#), [7](#), [8](#)
- [23] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Smpl: A skinned multi-person linear model. *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, 2015. [2](#), [4](#), [8](#)
- [24] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. *ArXiv*, abs/2308.09713, 2023. [3](#)
- [25] Qianli Ma, Jinlong Yang, Siyu Tang, and Michael J. Black. The power of points for modeling humans in clothing. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10954–10964, 2021. [2](#)
- [26] Miles Macklin and Matthias Müller. A constraint-based formulation of stable neo-hookean materials. *Proceedings of the 14th ACM SIGGRAPH Conference on Motion, Interaction and Games*, 2021. [5](#)
- [27] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7206–7215, 2020. [4](#)
- [28] Marko Mihajlovic, Aayush Bansal, Michael Zollhofer, Siyu Tang, and Shunsuke Saito. Keypointnerf: Generalizing

- image-based volumetric avatars using relative spatial encoding of keypoints. *ArXiv*, abs/2205.04992, 2022. 2
- [29] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf. *Communications of the ACM*, 65:99–106, 2020. 2, 4
- [30] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multi-resolution hash encoding. *ACM Transactions on Graphics (TOG)*, 41:1–15, 2022. 1, 2
- [31] Jesús R Nieto and Antonio Susín. Cage based deformations: a survey. In *Deformation Models: Tracking, Animation and Applications*, pages 75–99. Springer, 2012. 2
- [32] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 165–174, 2019. 4
- [33] Keunhong Park, U. Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B. Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5845–5854, 2020. 2
- [34] Keunhong Park, U. Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B. Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf. *ACM Transactions on Graphics (TOG)*, 40:1–12, 2021. 2
- [35] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9050–9059, 2020. 2, 3
- [36] Yicong Peng, Yichao Yan, Shengqi Liu, Yuhao Cheng, Shanyan Guan, Bowen Pan, Guangtao Zhai, and Xiaokang Yang. Cagenerf: Cage-based neural radiance field for generalized 3d deformation and animation. In *NeurIPS*, 2022. 3
- [37] Malte Prinzler, Otmar Hilliges, and Justus Thies. Diner: Depth-aware image-based neural radiance fields. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12449–12459, 2022. 2
- [38] Sergey Prokudin, Qianli Ma, Maxime Raafat, Julien Valentin, and Siyu Tang. Dynamic point fields. *arXiv preprint arXiv:2304.02626*, 2023. 2
- [39] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10313–10322, 2020. 3
- [40] Ravi Ramamoorthi and Pat Hanrahan. An efficient representation for irradiance environment maps. *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, 2001. 3
- [41] Edoardo Remelli, Timur M. Bagautdinov, Shunsuke Saito, Chenglei Wu, Tomas Simon, Shih-En Wei, Kaiwen Guo, Zhe Cao, Fabián Prada, Jason M. Saragih, and Yaser Sheikh. Drivable volumetric avatars using texel-aligned features. *ACM SIGGRAPH 2022 Conference Proceedings*, 2022. 2, 6, 7
- [42] Hang Si. Tetgen: A quality tetrahedral mesh generator and a 3d delaunay triangulator (version 1.5 — user’s manual). 2013. 4
- [43] Shih-Yang Su, Timur M. Bagautdinov, and Helge Rhodin. Danbo: Disentangled articulated neural body representations via graph neural networks. In *European Conference on Computer Vision*, 2022. 2
- [44] Shih-Yang Su, Timur M. Bagautdinov, and Helge Rhodin. Npc: Neural point characters from video. *ArXiv*, abs/2304.02013, 2023. 2, 3
- [45] Shih-Yang Su, Frank Yu, Michael Zollhoefer, and Helge Rhodin. A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose. In *Neural Information Processing Systems*, 2021. 2
- [46] Robert W. Sumner and Jovan Popović. Deformation transfer for triangle meshes. *ACM SIGGRAPH 2004 Papers*, 2004. 4
- [47] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 6105–6114. PMLR, 2019. 3
- [48] Ayush Tewari, Otto Fried, Justus Thies, Vincent Sitzmann, S. Lombardi, Z. Xu, Tanaba Simon, Matthias Nießner, Edgar Tretschk, L. Liu, Ben Mildenhall, Pranatharthi Srinivasan, R. Pandey, Sergio Orts-Escolano, S. Fanello, M. Guang Guo, Gordon Wetzstein, J y Zhu, Christian Theobalt, Manju Agrawala, Donald B. Goldman, and Michael Zollhöfer. Advances in neural rendering. *Computer Graphics Forum*, 41, 2021. 2
- [49] Ayush Kumar Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason M. Saragih, Matthias Nießner, Rohit Pandey, S. Fanello, Gordon Wetzstein, Jun-Yan Zhu, Christian Theobalt, Maneesh Agrawala, Eli Shechtman, Dan B. Goldman, and Michael Zollhofer. State of the art on neural rendering. *Computer Graphics Forum*, 39, 2020. 2
- [50] Shaofei Wang, Marko Mihajlovic, Qianli Ma, Andreas Geiger, and Siyu Tang. Metaavatar: Learning animatable clothed human models from few depth images. In *Neural Information Processing Systems*, 2021. 3
- [51] Yifan Wang, Noam Aigerman, Vladimir G. Kim, Siddhartha Chaudhuri, and Olga Sorkine-Hornung. Neural cages for detail-preserving 3d deformations. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 72–80. Computer Vision Foundation / IEEE, 2020. 2
- [52] Yifan Wang, Felice Serena, Shihao Wu, Cengiz Öztireli, and Olga Sorkine-Hornung. Differentiable surface splatting for point-based geometry processing. *ACM Transactions on Graphics (TOG)*, 38:1–14, 2019. 2, 3
- [53] Chung-Yi Weng, Brian Curless, Pratul P. Srinivasan, Jonathan T. Barron, and Ira Kemelmacher-Shlizerman. Humannerf: Free-viewpoint rendering of moving people from

- monocular video. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 16189–16199. IEEE, 2022. 2
- [54] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. *CoRR*, abs/2310.08528, 2023. 2
- [55] Donglai Xiang, Timur M. Bagautdinov, Tuur Stuyck, Fabián Prada, Javier Romero, Weipeng Xu, Shunsuke Saito, Jingfan Guo, Breannan Smith, Takaaki Shiratori, Yaser Sheikh, Jessica K. Hodgins, and Chenglei Wu. Dressing avatars. *ACM Transactions on Graphics (TOG)*, 41:1 – 15, 2022. 1
- [56] Donglai Xiang, Fabián Prada, Timur M. Bagautdinov, Weipeng Xu, Yuan Dong, He Wen, Jessica K. Hodgins, and Chenglei Wu. Modeling clothing as a separate layer for an animatable human avatar. *ACM Transactions on Graphics (TOG)*, 40:1 – 15, 2021. 1, 7
- [57] Donglai Xiang, Fabián Prada, Zhe Cao, Kaiwen Guo, Chenglei Wu, Jessica K. Hodgins, and Timur M. Bagautdinov. Drivable avatar clothing: Faithful full-body telepresence with dynamic clothing driven by sparse rgb-d input. 2023. 2
- [58] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Point-based neural radiance fields. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5428–5438, 2022. 2
- [59] Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. *arXiv preprint arXiv:2309.13101*, 2023. 2
- [60] Hao Zhang, Yao Feng, Peter Kulits, Yandong Wen, Justus Thies, and Michael J. Black. Text-guided generation and editing of compositional 3d avatars. *ArXiv*, abs/2309.07125, 2023. 2
- [61] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 586–595, 2018. 7
- [62] Yufeng Zheng, Yifan Wang, Gordon Wetzstein, Michael J. Black, and Otmar Hilliges. Pointavatar: Deformable point-based head avatars from videos. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21057–21067, 2022. 2, 3
- [63] Zerong Zheng, Han Huang, Tao Yu, Hongwen Zhang, Yandong Guo, and Yebin Liu. Structured local radiance fields for human avatar modeling. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15872–15882, 2022. 2
- [64] Wojciech Zielonka, Timo Bolkart, and Justus Thies. Instant volumetric head avatars. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4574–4584, 2022. 2, 3
- [65] Michael Zollhöfer, Justus Thies, Pablo Garrido, Derek Bradley, Thabo Beeler, Patrick Pérez, Marc Stamminger, Matthias Nießner, and Christian Theobalt. State of the art on monocular 3d face reconstruction, tracking, and applications. *Computer Graphics Forum*, 37, 2018. 2
- [66] Matthias Zwicker, Hans Rüdiger Pfister, Jeroen van Baar, and Markus H. Gross. Surface splatting. *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, 2001. 2, 3