

# 2DGS-Avatar: Animatable High-fidelity Clothed Avatar via 2D Gaussian Splatting

Qipeng Yan  
Academy for Engineering  
and Technology  
Fudan University  
Shanghai, China  
qipyan23@m.fudan.edu.cn

Mingyang Sun  
Academy for Engineering  
and Technology  
Fudan University  
Shanghai, China  
mysun21@m.fudan.edu.cn

Lihua Zhang\*  
Academy for Engineering  
and Technology  
Fudan University  
Shanghai, China  
lihuazhang@fudan.edu.cn

**Abstract**—Real-time rendering of high-fidelity and animatable avatars from monocular videos remains a challenging problem in computer vision and graphics. Over the past few years, the Neural Radiance Field (NeRF) has made significant progress in rendering quality but behaves poorly in run-time performance due to the low efficiency of volumetric rendering. Recently, methods based on 3D Gaussian Splatting (3DGS) have shown great potential in fast training and real-time rendering. However, they still suffer from artifacts caused by inaccurate geometry. To address these problems, we propose 2DGS-Avatar, a novel approach based on 2D Gaussian Splatting (2DGS) for modeling animatable clothed avatars with high-fidelity and fast training performance. Given monocular RGB videos as input, our method generates an avatar that can be driven by poses and rendered in real-time. Compared to 3DGS-based methods, our 2DGS-Avatar retains the advantages of fast training and rendering while also capturing detailed, dynamic, and photo-realistic appearances. We conduct abundant experiments on popular datasets such as AvatarRex and THuman4.0, demonstrating impressive performance in both qualitative and quantitative metrics.

**Index Terms**—animatable avatar, human reconstruction, 2D Gaussian splatting

## I. INTRODUCTION

Creating a high-fidelity and animatable avatar holds significant importance in fields such as AR/VR, entertainment, and film production. Over the past few years, Neural Radiance Fields (NeRF) [1] has been employed by some studies, enabling to reconstruct avatars from videos [2]–[5] and images [6], [7]. Though they achieve photo-realistic rendering, volumetric rendering in NeRF is inefficient and requires expensive training time and computational resources, making them impractical for real-world applications.

Recently, 3D Gaussian Splatting (3DGS) [8] provides a significant solution for fast training and rendering. In contrast to NeRF, 3DGS replaces the hierarchical volume sampling with a depth-based sort along the view direction, named splatting. Some methods [9]–[11] adopt 3DGS to model clothed humans, showing great potential in real-time rendering with high visual quality. However, since 3D Gaussian models use ellipsoids to represent objects, which contradicts the thin nature of surfaces, these approaches may produce fluctuating artifacts and fail to reconstruct accurate geometry. 2D Gaussian Splatting (2DGS)

[12] replaces 3D ellipsoids with 2D ellipses, which are similar to the triangular faces of the mesh, thus it is prone to converge to a more accurate geometry than 3DGS.

Inspired by 2DGS [12], we propose 2DGS-Avatar, a novel method for modeling animatable avatars, achieving fast run-time performance and superior geometry quality. Specifically, the avatar template is first represented by 2D Gaussian primitives in the canonical space, which is initialized by the vertices of the SMPL-X [13]. Then, the Linear Blend Skinning (LBS) [13], [14] is employed to transform these 2D Gaussian primitives from the canonical space to the posed space, with each primitive’s skinning weight assigned by querying a diffused skinning weight field [15]. Finally, we render the RGB images and normal maps with the differentiable rasterizer of 2DGS, which are supervised by the input RGB images and the corresponding normal maps. For better optimization of 2DGS, a self-supervised loss is put forward to ensure that the Gaussian primitives are uniformly distributed on the surface. During the densification phase, we enhance the original strategy in 2DGS with eccentricity filtering that removes the Gaussian primitives with particularly elongated, excessively large, or very low opacity.

To the best knowledge of us, we are the first work that employs 2DGS to model human avatars. In summary, our main contributions are as follows:

- We introduce 2DGS-Avatar, a novel real-time rendering pipeline for modeling animatable high-fidelity clothed avatars based on 2D Gaussian splatting.
- We propose a self-supervised loss that ensures Gaussian primitives are uniformly distributed on the surface, improving rendering results.
- We put forward eccentricity filtering to enhance the adaptive density control by removing particularly elongated Gaussian ellipses, resulting in smoother geometric edges.

## II. RELATED WORK

Recently, the emergence of 3DGS [9] has demonstrated impressive capabilities in 3D reconstruction, real-time rendering, and novel view synthesis. This work is also well-suited for representing avatars, leading to various methods [9]–[11] that extend the 3DGS pipeline to reconstruct human

\*Corresponding author

avatars from monocular RGB images. They represent avatars using Gaussian point clouds, optimizing the parameters of the Gaussians for rendering. These approaches can be categorized into two types: learning Gaussian parameters directly and learning Gaussian parameters through 2D maps.

#### A. Learning Gaussian Parameters Directly

Methods [9], [11] that directly learn Gaussian parameters typically follow a pipeline that is similar to NeRF-based approaches [2], [16], [17], where the avatar is represented in a canonical space and subsequently transformed into the posed space using LBS, after which the Gaussian primitives are rendered into images through the 3DGS rasterizer. The optimization of Gaussian parameters is performed by minimizing the error between the rendered images and the ground truth, a process that is largely similar to the parameter learning and optimization steps in 3DGS. Additionally, these methods often employ a Multi-Layer Perceptron (MLP) to refine the SMPL pose parameters and LBS skinning weights. Although such methods are characterized by relatively short training times, they are fundamentally limited by the tendency of MLP to prioritize low-frequency information [18], which consequently hampers their ability to accurately capture high-frequency details such as clothing textures, wrinkles, and other intricate geometric essential for achieving a high level of realism.

#### B. Learning Gaussian Parameters through 2D Maps

Methods that learn Gaussian parameters from 2D maps typically representing the human body in the posed space using 2D maps serve as the pose conditions, such as posed position maps [10], UV maps [19], and texture maps [20]. These methods utilize Convolutional Neural Networks (CNN) to directly learn the Gaussian parameters in the posed space. For instance, Animatable Gaussians [10] first learns a parametric template, which is then transformed from the canonical space to the posed space using LBS. For each posed space, the template is mapped into front and back posed position maps, and a StyleUNet [21] is employed to directly learn and optimize the Gaussian parameters. Subsequently, the 3DGS rasterizer is used to render the images. Due to the use of CNN, these methods are able to capture richer image features, leading to improvements compared to those directly optimizing Gaussian parameters. However, because these approaches primarily focus on optimizing CNN, they tend to converge more slowly. Though they can achieve real-time rendering, the training process is more resource-intensive in terms of training time and GPU memory.

### III. PRELIMINARY

#### A. SMPL-X

SMPL-X [13] is a parameterized human model that takes shape parameters  $\beta \in \mathbb{R}^{10}$  and pose parameters  $\theta \in \mathbb{R}^{K \times 3}$  and returns a triangulated mesh  $\mathcal{M}(\beta, \theta)$  by:

$$\mathcal{M}(\beta, \theta) = \text{LBS}(\mathcal{W}, J(\beta), \theta, T(\beta, \theta)), \quad (1)$$

where  $\mathcal{M}(\cdot) \in \mathbb{R}^{N \times 3}$ ,  $\text{LBS}(\cdot)$  denotes the Linear Blend Skinning (LBS) function,  $\mathcal{W}$  is the skinning weights,  $J(\cdot)$  represents the joint locations, and  $T(\cdot)$  is the template mesh with a star-like rest pose. We set  $N = 10475$  and  $K = 127$ , including the body, face, and hands. In our method, LBS is adopted as a transformation that maps the Gaussian primitives from the canonical space to the posed space. Specifically, given a point  $\mathbf{p}_i$  in the canonical space, the LBS function applies a set of linear transformations to map it to the posed space, resulting in the point  $\mathbf{p}'_i$ :

$$\mathbf{p}'_i = \sum_{k=1}^K w_{k,i} G'_k(\theta, J(\beta)) \mathbf{p}_i, \quad (2)$$

where  $w_{k,i}$  is the skinning weight of the  $k$ -th joint for the  $i$ -th point, and  $G'_k$  is the affine transformation matrix of the  $k$ -th joint from the canonical space to the posed space.

#### B. 2D Gaussian Splatting

2DGS [12] is a novel approach for modeling and reconstructing geometrically accurate radiance fields from multi-view images. The Gaussian primitives in 3DGS [9] are represented as 3D ellipsoids, while 2DGS flattens the 3D ellipsoid into a 2D ellipse, named surfels. Each Gaussian primitive is defined by its center point  $\mathbf{p}_c \in \mathbb{R}^3$ , opacity  $\alpha \in \mathbb{R}^1$ , view-dependent color  $\mathbf{c} \in \mathbb{R}^3$  calculated by spherical harmonics coefficients, scaling vector  $\mathbf{s} = (s_u, s_v) \in \mathbb{R}^2$  that controls the 2D Gaussian variance, and a rotation matrix  $\mathbf{R} \in \mathbb{R}^{3 \times 3}$ . The rotation matrix  $\mathbf{R} = [\mathbf{r}_u, \mathbf{r}_v, \mathbf{r}_w]$  is composed of two orthogonal vectors  $\mathbf{r}_u$  and  $\mathbf{r}_v$  on the Gaussian primitive, and the normal vector  $\mathbf{r}_w = \mathbf{r}_u \times \mathbf{r}_v$  of the Gaussian primitive, which is obtained by the cross product of these two vectors. Therefore, a 2D Gaussian is defined in the local tangent plane (also known as the  $uv$  space) in world space, which is represented as:

$$P(u, v) = \mathbf{p}_c + s_u \mathbf{r}_u u + s_v \mathbf{r}_v v = \mathbf{H}(u, v, 1, 1)^T, \quad (3)$$

$$\mathbf{H} = \begin{bmatrix} s_u \mathbf{r}_u & s_v \mathbf{r}_v & \mathbf{0} & \mathbf{p}_c \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} \mathbf{R}\mathbf{S} & \mathbf{p}_c \\ \mathbf{0} & 1 \end{bmatrix}, \quad (4)$$

where  $\mathbf{H} \in \mathbb{R}^{4 \times 4}$  represents the geometry of the Gaussian primitive. For a point  $\mathbf{u} = (u, v)$  in  $uv$  space, its Gaussian value  $\mathcal{G}(\mathbf{u})$  can be simplified and computed as a Gaussian with a mean of 0 and a variance of 1:

$$\mathcal{G}(\mathbf{u}) = \exp\left(-\frac{u^2 + v^2}{2}\right). \quad (5)$$

Then, its coordinates  $P(u, v)$  in the world space can be obtained using (3).

2DGS maps a pixel  $\mathbf{x} = (x, y)$  in screen space to its corresponding point  $\mathbf{u} = (u, v)$  in the  $uv$  space through the Ray-splat Intersection  $F$ , after which each pixel  $\mathbf{c}(\mathbf{x})$  is computed using alpha blending:

$$\mathbf{c}(\mathbf{x}) = \sum_{i=1} \mathbf{c}_i \alpha_i \mathcal{G}_i(F(\mathbf{x})) \prod_{j=1}^{i-1} (1 - \alpha_j \mathcal{G}_j(F(\mathbf{x}))), \quad (6)$$

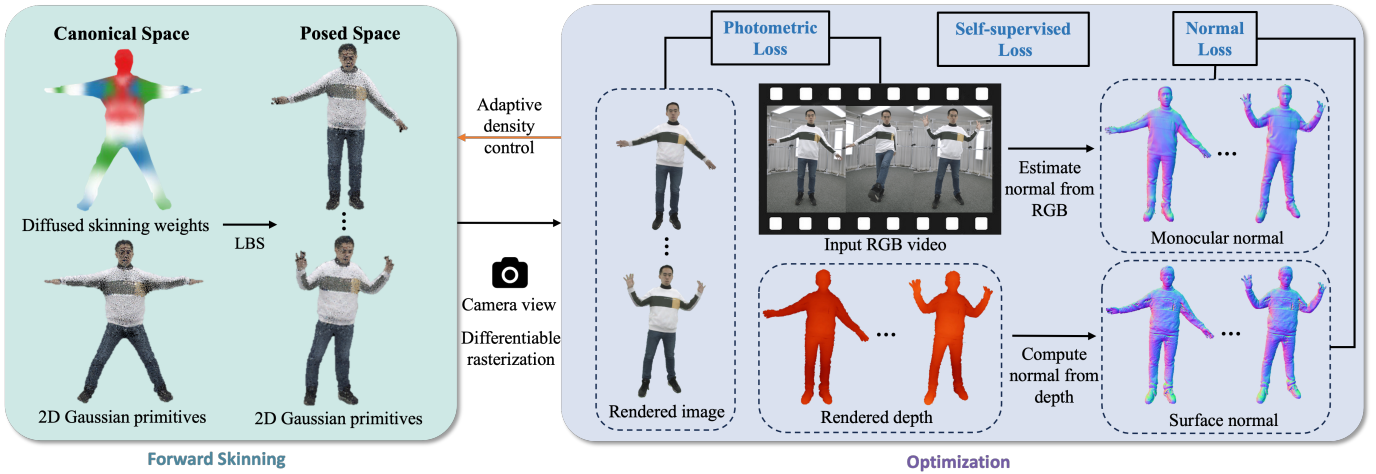


Fig. 1. Illustration of the pipeline. The orange arrows indicate backpropagation. Our method consists of two parts: (1) Transforming Gaussian primitives from the canonical space to the posed space through forward skinning, followed by rasterization to render images and depth maps in the posed space. (2) Optimizing the Gaussian primitives in the canonical space using photometric loss, normal loss, and self-supervised loss.

where  $c_i$  represents the color of the Gaussian primitive calculated by spherical harmonic. In summary, the learnable parameters of  $\mathcal{G}_i$  are  $\Theta_i = \{\mathbf{p}_i, \mathbf{s}_i, \mathbf{R}_i, \alpha_i, \mathbf{c}_i\}$ .

#### IV. METHOD

Given multi-view RGB videos and the related SMPL-X registrations that include the pose and shape parameters for the character in each frame, our goal is to create an animatable high-fidelity clothed avatar. The pipeline is shown in Fig. 1.

First, we precompute a skinning weight field [15] in the canonical space, diffusing the skinning weights from the SMPL-X surface to the entire canonical space. This allows us to obtain the skinning weight of each Gaussian primitive by querying the skinning weight field. Second, we transform the Gaussian primitives from the canonical space to the posed space using (2), followed by rasterization to render images and depth maps of the Gaussian primitives in the posed space. We optimize the model by minimizing the photometric difference between the rendered images and the corresponding frames of the input RGB videos, as well as the difference between the normals computed from the depth maps and those estimated from the input RGB images. Additionally, we propose a self-supervised loss to constrain the distribution of Gaussian primitives and the smoothness of the normal maps. Finally, we propose an eccentricity filtering algorithm to control the density adaptively.

##### A. Forward Skinning

We initialize a set of Gaussian primitives  $\{\mathcal{G}_i\}_{i=1}^N$  at the vertices of the SMPL-X model in the canonical space according to the shape parameters  $\beta$  from the dataset. We then query the precomputed diffused skinning weight field using the center points  $\mathbf{p}_c$  of the Gaussian primitives  $\mathcal{G}_i$  to obtain the corresponding LBS weights  $w_i$ . After each optimization step, the skinning weights are re-queried. The center point  $\mathbf{p}_c$

is transformed from the canonical space to the corresponding point  $\mathbf{p}'_c$  in the posed space by (2).

##### B. Splatting

Based on the camera’s intrinsic and extrinsic parameters from the input RGB images, the rendered image can be obtained using alpha blending  $c(\mathbf{x}) = \sum_{i=1}^N c_i \alpha_i \mathcal{G}_i(F(\mathbf{x})) T_i$ , where  $T_i = \prod_{j=1}^{i-1} (1 - \alpha_j \mathcal{G}_j(F(\mathbf{x})))$  represents the accumulated transmittance along the ray from  $\mathcal{G}_1$  to  $\mathcal{G}_{i-1}$ , indicating the visibility of  $\mathcal{G}_i$ . Similar to 2DGS, we consider  $T_i = 0.5$  as the surface. Therefore, we only take the depth maps of the visible surface, defined as  $z = \max\{z_i \mid T_i > 0.5\}$ . Based on the depth map, we can also compute the normal vectors. Specifically, for a point  $\mathbf{p}$  in the depth map, with its neighboring points along the x-axis  $\mathbf{p}_x$  and the y-axis  $\mathbf{p}_y$ , the normal vector  $\mathbf{n}$  can be calculated using the following equation:

$$\mathbf{n} = \text{Normalize}((\mathbf{p} - \mathbf{p}_x) \times (\mathbf{p} - \mathbf{p}_y)). \quad (7)$$

##### C. Optimization

To optimize the parameters  $\Theta$  of the Gaussian primitives  $\mathcal{G}$ , our loss function  $\mathcal{L}$  consists of four components: the photometric loss  $\mathcal{L}_p$ , the normal loss  $\mathcal{L}_n$ , the self-supervised loss  $\mathcal{L}_s$ , and the mask loss  $\mathcal{L}_m$ . The total loss  $\mathcal{L}$  is given by:

$$\mathcal{L} = \mathcal{L}_p + \lambda_n \mathcal{L}_n + \lambda_s \mathcal{L}_s + \lambda_m \mathcal{L}_m. \quad (8)$$

**Photometric Loss.** The photometric loss consists of two terms same to 2DGS: an  $L_1$  term and a D-SSIM term. In addition to these, we include an additional term, the Learned Perceptual Image Patch Similarity (LPIPS) [22], to minimize the difference between the rendered image  $\hat{\mathbf{I}}$  and the input image  $\mathbf{I}$ :

$$\mathcal{L}_p = L_1(\hat{\mathbf{I}}, \mathbf{I}) + \lambda_{dssim} L_{dssim}(\hat{\mathbf{I}}, \mathbf{I}) + \lambda_{lrips} L_{lrips}(\hat{\mathbf{I}}, \mathbf{I}). \quad (9)$$

**Normal Loss.** Using only the photometric loss is insufficient for accurately modeling human geometry, especially in

high-frequency regions. Therefore, we use the normal loss as a prior. We apply PIFuHD [23] to infer the normal map  $\mathbf{N}$  from the input RGB image and compute the normal map  $\hat{\mathbf{N}}$  from the rendered depth map. The geometry of the avatar is constrained by minimizing the cosine similarity between two normal maps:

$$\mathcal{L}_n = \lambda_n(1 - \hat{\mathbf{N}} \cdot \mathbf{N}). \quad (10)$$

**Self-supervised Loss.** Gaussian primitives are typically not uniformly distributed, there are more Gaussian primitives in high-frequency regions and fewer in low-frequency regions. Since our LBS weights are diffused from the SMPL model, which is inherently designed for meshes, we propose a self-supervised area loss  $L_{area}$  to constrain the scaling of each Gaussian primitive by minimizing the variance of the product of the two scaling vectors  $(s_u, s_v)$ . This approach ensures that the Gaussian primitives are distributed evenly, similar to triangular faces of mesh. Additionally, we observed that there are some Gaussian primitives with low opacity inside or outside the avatar, which are not meaningful. To address this, we introduce an opacity loss  $L_{opacity}$  inspired by Gaussian surfels [25], encouraging the opacity of the Gaussian primitives to be either close to 1 or 0, thereby ensuring that all Gaussian primitives are distributed on the surface of the avatar:

$$\mathcal{L}_s = \lambda_{area}L_{area} + \lambda_{opacity}L_{opacity}, \quad (11)$$

where  $L_{opacity} = \exp(-(\alpha_i - 0.5)^2/0.05)$ .

**Mask Loss.** Following NeuS [26], we include a mask loss  $\mathcal{L}_m$ , which is obtained by calculating the binary cross-entropy between the alpha map  $\hat{\mathbf{M}} = \sum_{i=1}^N \alpha_i \mathcal{G}_i(F(\mathbf{x}))T_i$  and the mask  $\mathbf{M}$  from the dataset:

$$\mathcal{L}_m = \lambda_m \text{BCE}(\hat{\mathbf{M}}, \mathbf{M}). \quad (12)$$

**Eccentricity Filtering.** The previous area loss only constrains the area of the Gaussian primitives, but this can still result in some very elongated ellipses, which may lead to unsmooth geometry at the edges. Therefore, we propose eccentricity filtering for adaptive density control, which removes Gaussian primitives with an eccentricity higher than a threshold that we set to 9.

## V. EXPERIMENTS

### A. Experimental Setup

**Datasets.** Our experiments are conducted on two popular datasets, AvatarRex [24] and THuman4.0 [27], including 4 sequences with 16 views from the AvatarRex and 3 sequences with 24 views from the THuman4.0. From each dataset, we select 2 sequences for evaluation. Both datasets provide RGB images, masks, and SMPL-X registrations. In addition, we utilize PIFuHD [23] to estimate normal maps from the RGB images for further supervision in our pipeline.

**Metric.** We select Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM) [28], Learned Perceptual Image Patch Similarity (LPIPS) [22] as well as training time on an NVIDIA V100 GPU as quantitative metrics for comparative experiments.

TABLE I  
QUANTITATIVE COMPARISON WITH STATE-OF-THE-ART METHODS.

Method	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Train
Ours	<u>30.93</u>	<u>0.9643</u>	<u>31.37</u>	<b>1 h</b>
GauHuman [11]	29.57	0.9639	35.93	<b>1 h</b>
Animatable Gaussians [10]	<b>31.86</b>	<b>0.9705</b>	<b>29.32</b>	7 h

• The best results are shown in **bold**, while the second best performance is underlined.

**Baselines.** We compare our approach with state-of-the-art 3DGS-based methods [10], [11] from two categories: learning Gaussian parameters directly and learning Gaussian parameters through 2D Maps. For each sequence, we split the data into training and testing sets. All methods are run for 30,000 iterations under their respective default parameter settings. We report the average metric values for the selected sequences across all methods.

### B. Results

**Reconstruction.** Fig. 2 and Table I summarize the comparison results between our 2DGS-Avatar, GauHuman [11], and Animatable Gaussians [19]. Since GauHuman preloads the images into memory before training, a process that takes approximately 30 minutes, rather than using a generator like other methods, we included this time as part of its training time for fairness. As shown in Fig. 2, both our method and Animatable Gaussians produce visibly superior rendering quality compared to GauHuman, with clearer textures in clothing and facial features. This is because GauHuman relies solely on MLP to learn image features, which limits its ability to optimize Gaussian properties effectively. From Table I, it can be observed that both our method and GauHuman require only one hour of training time. However, our approach and Animatable Gaussians outperform GauHuman in all quantitative metrics. While our method falls slightly behind Animatable Gaussians, we achieve similar results in only one-seventh of the training time, with minimal perceptible differences in the rendered images. Additionally, Animatable Gaussians require significantly more GPU memory compared to our approach, further highlighting the efficiency of our method.

**Animation.** As shown in Fig. 2 and Fig. 3, the reconstructed avatar can be driven by LBS with novel pose sequences sampled from AMASS [29] and THuman4.0\_POSE [27]. The run-time performance can reach 60 FPS on a single NVIDIA RTX-4070s GPU, which is competent in real-world applications.

### C. Ablation Study

We study the effect of different components proposed in our method, including the area loss, normal loss, and eccentricity filtering strategy. The metrics are presented in the Table II. The full model achieves the best results across all quantitative metrics, demonstrating that all the proposed modules are effective, and optimal performance is only achieved when all modules work together. In addition, we also visualize the effect of  $\mathcal{L}_{area}$  in Fig. 4. The results show that our proposed  $\mathcal{L}_{area}$



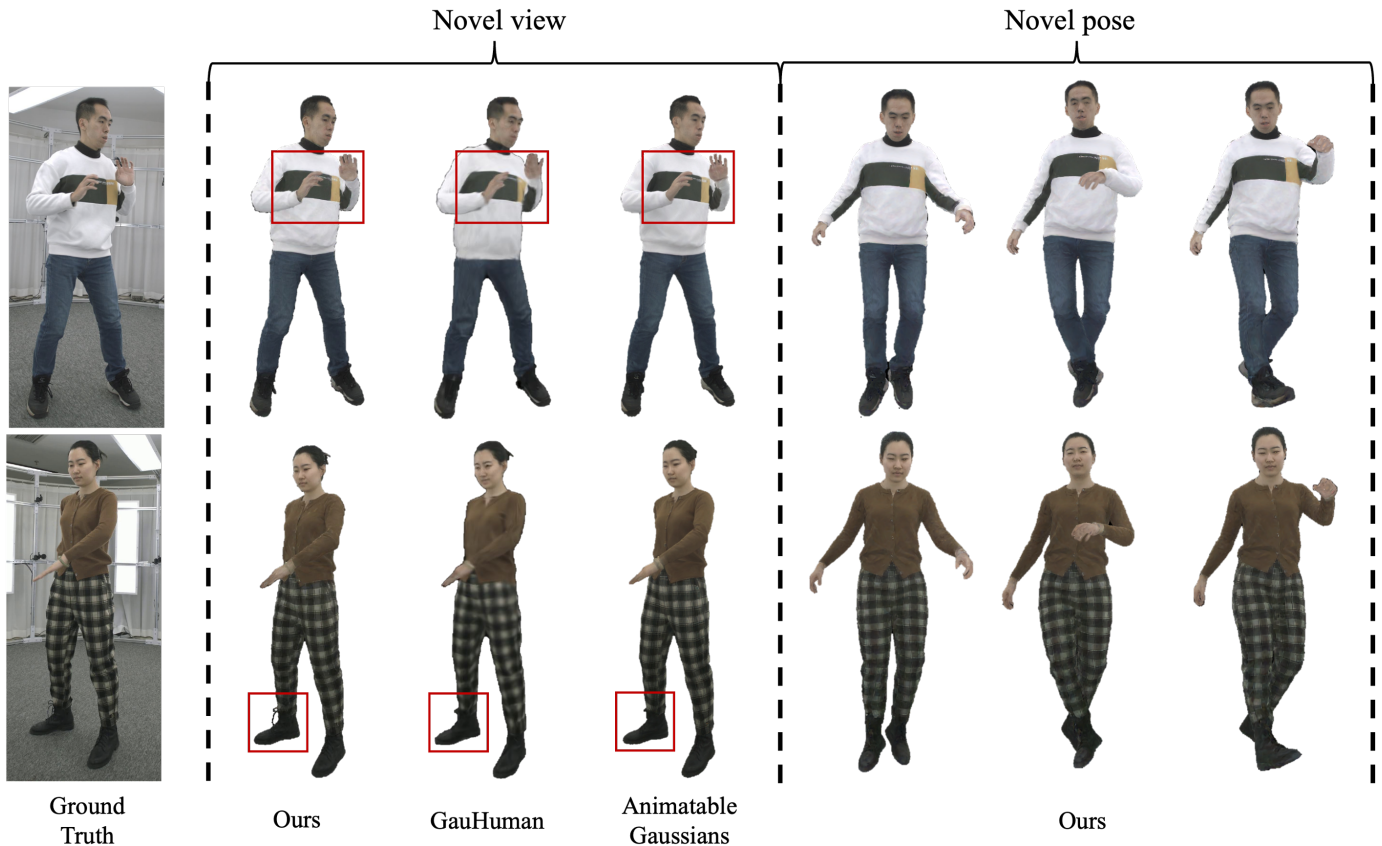


Fig. 2. Qualitative comparison on AvatarRex [24]. We show the results for both novel view and novel pose on sequences of “avatarrex\_zzr” and “avatarrex\_lbn2” in AvatarRex. Our method reaches comparable visual effects to Animatable Gaussians [19] while surpassing GauHuman [11] in terms of surface details, such as hands, clothes and shoes.



Fig. 3. More results on sequences of “subject00” and “subject02” in THuman4.0 [27] with novel pose.

TABLE II  
ABLATION STUDY ON AVATARREX [24].

	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
Full model	<b>31.36</b>	<b>0.9781</b>	<b>34.75</b>
w/o $L_{area}$	31.35	0.9775	36.56
w/o $L_{normal}$	31.26	0.9773	35.22
w/o eccentricity filtering	31.19	0.9772	35.83

• The best results are shown in **bold**.

effectively leads to a more uniform distribution of Gaussian primitives.

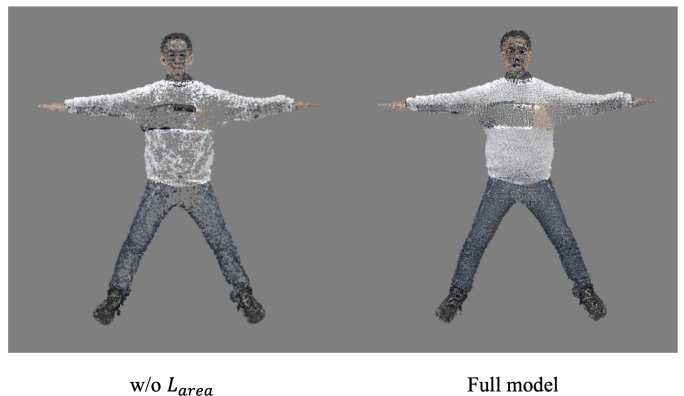


Fig. 4. The visualization of the ablation study on  $L_{area}$ . With  $L_{area}$ , the Gaussian primitives are prone to converge towards a more uniform distribution around the surface.

## VI. DISCUSSION

**Conclusion.** In this paper, we introduced 2DGS-Avatar, which, to the best of our knowledge, is the first method to represent clothed avatars using 2DGS. Our approach efficiently reconstructs high-fidelity clothed avatars from monocular RGB videos and enables real-time rendering. Experimental results

demonstrate that our method strikes an effective trade-off between rendering quality, memory consumption and training time, achieving near state-of-the-art performance with significantly reduced computational resources and time.

**Limitations.** (1) When the input RGB videos are relatively blurry, the reconstruction quality tends to degrade, particularly in shadowed areas where lighting is insufficient. Gaussian-DK [30] alleviates this problem by extracting a lightness map. (2) Although 2DGS-Avatar is capable of reconstructing high-fidelity avatars, the animation of the avatar, particularly in simulating clothing wrinkles, lacks a high level of realism. IF-Garments [31] shows great potential in modeling clothing details by combining neural fields with XPBD [32]. (3) It still remains a challenge to model the less frequently observed regions, such as underarms or shoe soles.

## REFERENCES

- [1] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” in *Computer Vision – ECCV 2020, Lecture Notes in Computer Science*, Jan 2020, p. 405–421.
- [2] S. Peng, Y. Zhang, Y. Xu, Q. Wang, Q. Shuai, H. Bao, and X. Zhou, “Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9054–9063.
- [3] Y. Feng, J. Yang, M. Pollefeys, M. J. Black, and T. Bolkart, “Capturing and animation of body and clothing from monocular video,” in *SIGGRAPH Asia 2022 Conference Papers*, 2022, pp. 1–9.
- [4] W. Jiang, K. M. Yi, G. Samei, O. Tuzel, and A. Ranjan, “Neuman: Neural human radiance field from a single video,” in *European Conference on Computer Vision*. Springer, 2022, pp. 402–418.
- [5] S.-Y. Su, F. Yu, M. Zollhöfer, and H. Rhodin, “A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose,” *Advances in neural information processing systems*, vol. 34, pp. 12 278–12 291, 2021.
- [6] S. Cha, K. Seo, A. Ashtari, and J. Noh, “Generating texture for 3d human avatar from a single image using sampling and refinement networks,” in *Computer Graphics Forum*, vol. 42, no. 2. Wiley Online Library, 2023, pp. 385–396.
- [7] F. Zhao, W. Yang, J. Zhang, P. Lin, Y. Zhang, J. Yu, and L. Xu, “Humannerf: Efficiently generated human radiance field from sparse inputs,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 7743–7753.
- [8] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, “3d gaussian splatting for real-time radiance field rendering,” *ACM Trans. Graph.*, vol. 42, no. 4, pp. 139–1, 2023.
- [9] Z. Qian, S. Wang, M. Mihajlovic, A. Geiger, and S. Tang, “3dgs-avatar: Animatable avatars via deformable 3d gaussian splatting,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 5020–5030.
- [10] Z. Li, Z. Zheng, L. Wang, and Y. Liu, “Animatable gaussians: Learning pose-dependent gaussian maps for high-fidelity human avatar modeling,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 19 711–19 722.
- [11] S. Hu, T. Hu, and Z. Liu, “Gauhuman: Articulated gaussian splatting from monocular human videos,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 20 418–20 431.
- [12] B. Huang, Z. Yu, A. Chen, A. Geiger, and S. Gao, “2d gaussian splatting for geometrically accurate radiance fields,” in *ACM SIGGRAPH 2024 Conference Papers*, 2024, pp. 1–11.
- [13] G. Pavlakos, V. Choutas, N. Ghorbani, T. Bolkart, A. A. A. Osman, D. Tzionas, and M. J. Black, “Expressive body capture: 3D hands, face, and body from a single image,” in *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 10 975–10 985.
- [14] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black, “SMPL: A skinned multi-person linear model,” *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, vol. 34, no. 6, pp. 248:1–248:16, Oct. 2015.
- [15] S. Lin, H. Zhang, Z. Zheng, R. Shao, and Y. Liu, “Learning implicit templates for point-based clothed human modeling,” in *European Conference on Computer Vision*. Springer, 2022, pp. 210–228.
- [16] S. Peng, J. Dong, Q. Wang, S. Zhang, Q. Shuai, X. Zhou, and H. Bao, “Animatable neural radiance fields for modeling dynamic human bodies,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14 314–14 323.
- [17] C.-Y. Weng, B. Curless, P. P. Srinivasan, J. T. Barron, and I. Kemelmacher-Shlizerman, “Humannerf: Free-viewpoint rendering of moving people from monocular video,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 16 210–16 220.
- [18] M. Tancik, P. P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghuvaran, U. Singhal, R. Ramamoorthi, J. T. Barron, and R. Ng, “Fourier features let networks learn high frequency functions in low dimensional domains,” *NeurIPS*, 2020.
- [19] L. Hu, H. Zhang, Y. Zhang, B. Zhou, B. Liu, S. Zhang, and L. Nie, “Gaussianavatar: Towards realistic human avatar modeling from a single video via animatable 3d gaussians,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [20] H. Pang, H. Zhu, A. Kortylewski, C. Theobalt, and M. Habermann, “Ash: Animatable gaussian splats for efficient and photoreal human rendering,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024, pp. 1165–1175.
- [21] L. Wang, X. Zhao, J. Sun, Y. Zhang, H. Zhang, T. Yu, and Y. Liu, “Styleavatar: Real-time photo-realistic portrait avatar from a single video,” in *ACM SIGGRAPH 2023 Conference Proceedings*, 2023, pp. 1–10.
- [22] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, “The unreasonable effectiveness of deep features as a perceptual metric,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.
- [23] S. Saito, T. Simon, J. Saragih, and H. Joo, “Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization,” in *CVPR*, 2020.
- [24] Z. Zheng, X. Zhao, H. Zhang, B. Liu, and Y. Liu, “Avatarrex: Real-time expressive full-body avatars,” *ACM Transactions on Graphics (TOG)*, vol. 42, no. 4, 2023.
- [25] P. Dai, J. Xu, W. Xie, X. Liu, H. Wang, and W. Xu, “High-quality surface reconstruction using gaussian surfels,” in *ACM SIGGRAPH 2024 Conference Papers*, 2024, pp. 1–11.
- [26] P. Wang, L. Liu, Y. Liu, C. Theobalt, T. Komura, and W. Wang, “Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction,” *arXiv preprint arXiv:2106.10689*, 2021.
- [27] Z. Zheng, H. Huang, T. Yu, H. Zhang, Y. Guo, and Y. Liu, “Structured local radiance fields for human avatar modeling,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022.
- [28] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli, “Image quality assessment: from error visibility to structural similarity,” *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [29] N. Mahmood, N. Ghorbani, N. F. Troje, G. Pons-Moll, and M. J. Black, “Amass: Archive of motion capture as surface shapes,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 5442–5451.
- [30] S. Ye, Z.-H. Dong, Y. Hu, Y.-H. Wen, and Y.-J. Liu, “Gaussian in the dark: Real-time view synthesis from inconsistent dark images using gaussian splatting,” *arXiv preprint arXiv:2408.09130*, 2024.
- [31] M. Sun, Q. Yan, Z. Liang, D. Kou, D. Yang, R. Yuan, X. Zhao, M. Li, and L. Zhang, “If-garments: Reconstructing your intersection-free multi-layered garments from monocular videos,” in *ACM Multimedia 2024*.
- [32] M. Macklin, M. Müller, and N. Chentanez, “Xpbd: position-based simulation of compliant constrained dynamics,” in *Proceedings of the 9th International Conference on Motion in Games*, 2016, pp. 49–54.