# CHASE: 3D-Consistent Human Avatars with Sparse Inputs via Gaussian Splatting and Contrastive Learning

Haoyu Zhao[* 1,2], Hao Wang[* 3], Chen Yang[* 1], Wei Shen[†1]

[1]MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University
[2]School of Computer Science, Wuhan University
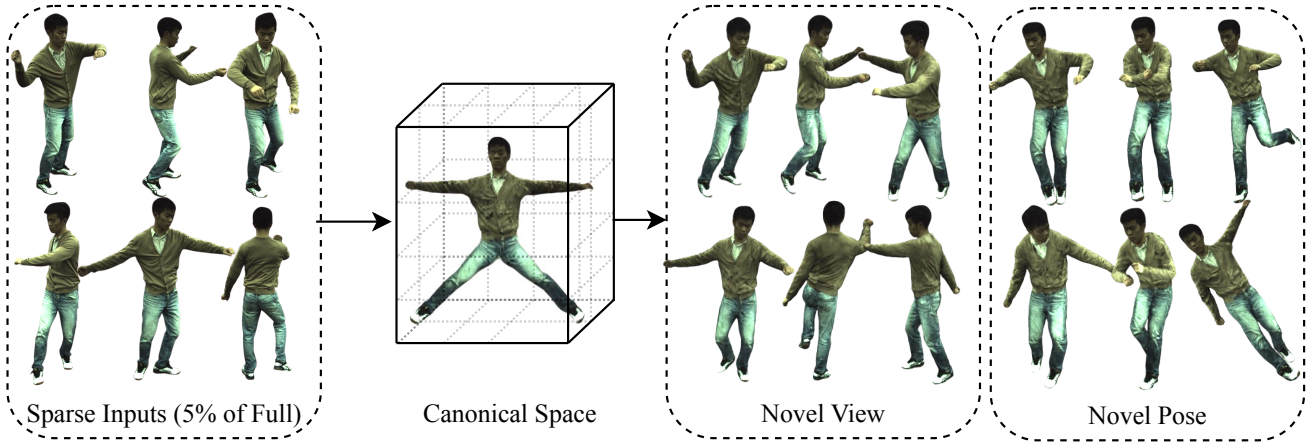[3]Wuhan National Laboratory for Optoelectronics, Huazhong University of Science and Technology

Figure 1. **CHASE.** We propose an efficient method for creating 3D-consistent animatable avatars from just videos. Our method achieve better quality to the most recent SOTA methods [8, 29, 37] in both full and sparse inputs.

## Abstract

*Existing approaches for human avatar generation–both NeRF-based and 3D Gaussian Splatting (3DGS) based– struggle with maintaining 3D consistency and exhibit degraded detail reconstruction, particularly when training with sparse inputs. To address this challenge, we propose CHASE, a novel framework that achieves dense-input-level performance using only sparse inputs through two key innovations: cross-pose intrinsic 3D consistency supervision and 3D geometry contrastive learning. Building upon prior skeleton-driven approaches that combine rigid deformation with non-rigid cloth dynamics, we first establish baseline avatars with fundamental 3D consistency. To enhance 3D consistency under sparse inputs, we introduce a Dynamic Avatar Adjustment (DAA) module, which refines deformed Gaussians by leveraging similar poses from the training set. By minimizing the rendering discrepancy between adjusted Gaussians and reference poses, DAA provides additional supervision for avatar reconstruction. We further maintain global 3D consistency through a novel geometry-aware contrastive learning strategy. While designed for sparse inputs, CHASE surpasses state-of-the-art methods across both full and sparse settings on ZJU-MoCap and H36M datasets, demonstrating that our enhanced 3D consistency leads to superior rendering quality. Project page: https://chaseprojectpage.github.io/.*

## 1. Introduction

Photo-realistic rendering and animation of human bodies is a critical research area with wide-ranging applications in AR/VR, visual effects, virtual try-on, and film production [6]. Early approaches [23] relied on multi-camera setups to capture high-quality data, requiring extensive computational resources and manual effort. While these methods perform well for reconstructing a single scene or object with sufficient input views, they struggle to generalize to

---

* Equal contributions.
†Corresponding Author.
Haoyu Zhao completed this work during an internship at Shanghai Jiao Tong University.

new scenes or objects from limited samples [15].

Recent advancements have explored using neural radiance fields (NeRF) for modeling 3D human avatars [21], typically employing parametric body models to model deformations. Some methods [2, 11, 46] use human template models to facilitate generalizable and robust synthesis. However, NeRF-based methods are less efficient to train and render due to their computationally intensive per-pixel volume rendering process.

Point-based rendering [48] has emerged as an efficient alternative to NeRFs, offering significantly faster rendering. The recently proposed 3D Gaussian Splatting (3DGS) [12] gains popularity for its fast rendering speed. Numerous works have further explored the 3D Gaussian representation for dynamic 3D human avatars [7, 8, 14, 16, 22, 29, 32, 34]. However, these methods often face challenges in maintaining 3D consistency and producing high-quality reconstructions, particularly with sparse inputs.

To address the aforementioned issues, we propose **CHASE**, which is capable of reconstructing 3D **C**onsistent **H**uman **A**vatars with **S**parse inputs via Gaussian Splatting and contrastiv**E** learning, as shown in Fig. 1. We first integrate a skeleton-driven rigid deformation and a non-rigid cloth dynamics deformation to create a human avatar. To enhance 3D consistency under sparse inputs, we utilize the intrinsic 3D consistency of images across different poses within the same person. Specifically, for each training pose/image, we select a similar pose/image from the dataset and then adjust the deformed Gaussians using the proposed Dynamic Avatar Adjustment (DAA), an explicit point-based control graph adjustment strategy, to the selected similar pose. Then, we minimize differences between the rendered image of the adjusted Gaussians and the image corresponding to the selected similar pose, which serves as an additional form of supervision for human avatars. Additionally, we employ 3D geometry contrastive learning, utilizing features from a 3D feature extractor, to further enhance the global 3D consistency of generated human avatars. We conduct extensive experiments on the ZJU-MoCap data [26], H36M [10] and surprisingly find that *CHASE outperforms other SOTAs in both full and sparse inputs setting.* Our work makes the following contributions:

- We propose an explicit point-based control graph adjustment strategy, which introduces a novel 2D image supervision to 3D human body modeling, enhancing the 3D consistency of human avatars.

- We propose a 3D geometry contrastive learning to enforce consistency across different representations of the same pose and enhance global 3D understanding.

- Extensive experiments show that our CHASE achieves SOTA performance quantitatively and qualitatively under full and sparse settings.

## 2. Related Work

### 2.1. Contrastive Representation Learning

Contrastive Representation Learning is one of the mainstream self-supervised learning paradigms, which learns potential semantics from constructed invariance or equivariance. In 3D, PointContrast [39] proposes geometric augmentation to generate positive and negative pairs. Cross-Point [1] uses both inter- and intra-modal contrastive learning. PointCLIP [45] achieves image-point alignment by projecting point clouds onto 2D depth images. RECON [28] focuses on single- and cross-modal contrastive learning through discriminative contrast [13] or global feature alignment [30]. Our CHASE introduces a 3D geometry contrastive learning method to enforce consistency across different representations of the same pose.

### 2.2. 3D Editing and Deformation

Traditional deformation methods in computer graphics are typically based on Laplacian coordinates [4], Poisson equations [43], and cage-based methods [41]. However, these methods often rely on implicit and computationally expensive NeRF-based approaches.

Numerous works [3, 47] have proposed techniques for editing 3D Gaussian Splatting (3DGS) [12]. SuGaR [5] introduces a mesh extraction method that produces meshes from 3DGS, which can then be edited. SC-GS [9] proposes deforming Gaussians by transferring the movement of control points. Our CHASE employs a novel explicit point-based control graph deformation strategy, which is more efficient than previous methods.

### 2.3. 3D Human Modeling

Since the high-quality rendering achieved by the seminal work Neural Radiance Fields (NeRF) [21], there has been a surge of research on neural rendering for human avatars [17,18,26,31]. Although NeRF is designed for static objects, HumanNeRF [38] extends NeRF to enable capturing dynamic human motion using just a single monocular video. Neural Body [26] associates a latent code with each SMPL [19] vertex to encode appearance, which is then transformed into observation space based on the human pose. Furthermore, Neural Actor [18] learns a deformable radiance field with SMPL [19] as guidance and utilizes a texture map to improve the final rendering quality. Posevocab [17] designs joint-structured pose embeddings to encode dynamic appearances under different key poses, allowing for more effective learning of joint-related appearances. However, a major limitation of NeRF-based methods is that NeRFs are slow to train and render.

Point-based rendering [48] has proven to be an efficient alternative to NeRFs for fast inference and training. Extending point clouds to 3D Gaussians, 3D Gaussian Splatting

(3DGS) [12] models the rendering process by splatting a set of 3D Gaussians onto the image plane via alpha blending. Given the impressive performance of 3DGS in both quality and speed, numerous works have further explored the 3D Gaussian representation for dynamic 3D human avatar reconstruction [8, 14, 16, 29]. Human Gaussian Splatting [22] showcases 3DGS as an efficient alternative to NeRF. SplattingAvatar [32] and GomAvatar [37] extend lifted optimization to simultaneously optimize the parameters of the Gaussians while walking on the triangle mesh. However, these methods struggle to maintain 3D consistency and produce low-quality reconstructions when applied to human avatar creation with only sparse inputs. Our CHASE introduces a novel 2D image supervision to 3D human body modeling and 3D geometry contrastive learning, enhancing the 3D consistency of human avatars.

## 3. Preliminaries

**SMPL [19].** The SMPL model is a pre-trained parametric human model representing body shape and pose. In SMPL, body shape and pose are controlled by pose and shape. In this work, we apply the Linear Blend Skinning (LBS) algorithm used in SMPL to transform points from a canonical space to a posed space.

**LBS [33].** Linear Blend Skinning (LBS) is a weight-based technique that associates each vertex with one or more joints and uses weight values to describe the influence of each joint on the vertex. Vertex deformation is calculated by linearly interpolating transformations on the associated joints: $\mathcal{X}'_v = \sum_{j=1}^{J} w_j(\mathcal{X}_v) B_j \mathcal{X}_v$, where $J$ represents the number of joints, $N$ represents the number of vertices, $\mathcal{X}'_v \in \mathbb{R}^{N \times 3}$ is the new position of the skinned vertex, $w \in \mathbb{R}^{N \times J}$ is the skinning weight matrix, $B \in \mathbb{R}^{J \times 4 \times 4}$ is the affine transformation matrix of each joint representing rotation and translation, (i.e. bone transforms) and $\mathcal{X}_v \in \mathbb{R}^{N \times 3}$ is the original mesh vertex position.

**3D Gaussian Splatting (3DGS) [12].** 3DGS explicitly represents scenes using point clouds, where each point is modeled as a 3D Gaussian defined by a covariance matrix $\Sigma$ and a center point $\mathcal{X}$, the latter referred to as the mean. The value at point $\mathcal{X}$ is:

$$G(\mathcal{X}) = e^{-\frac{1}{2}\mathcal{X}^T \Sigma^{-1} \mathcal{X}}. \quad (1)$$

For differentiable optimization, the covariance matrix $\Sigma$ is decomposed into a scaling matrix $\mathcal{S}$ and a rotation matrix $\mathcal{R}$, such that $\Sigma = \mathcal{R}\mathcal{S}\mathcal{S}^T\mathcal{R}^T$. $\mathcal{S}$ and $\mathcal{R}$ are stored as the diagonal vector $s \in \mathbb{R}^{N \times 3}$ and a quaternion vector $r \in \mathbb{R}^{N \times 4}$, respectively.

In rendering novel views, differential splatting as introduced by [42], involves using a viewing transform $W$ and the Jacobian matrix $J$ of the affine approximation of the projective transformation to compute the transformed covariance matrix: $\Sigma' = JW\Sigma W^T J^T$. The color and opacity at each pixel are computed from the Gaussian's representations: $G(\mathcal{X}) = e^{-\frac{1}{2}\mathcal{X}^T \Sigma^{-1} \mathcal{X}}$. The blending of $N$ ordered points overlapping a pixel is given by the formula: $\mathcal{C} = \sum_{i \in N} c_i \alpha_i \prod_{j=1}^{i-1}(1 - \alpha_i)$, where $c_i$, $\alpha_i$ represent the density and color of this point computed by a 3D Gaussian $G$ with covariance $\Sigma$ multiplied by an optimizable per-point opacity and SH color coefficients.

## 4. Method

We illustrate the pipeline of our CHASE in Fig. 2. The inputs include images $X = \{x_i\}_{i=1}^N$ obtained from monocular videos, fitted SMPL parameters $P = \{p_i\}_{i=1}^N$, and foreground masks $M = \{m_i\}_{i=1}^N$ of images. CHASE optimizes 3D Gaussians in canonical space, which are then be deformed to match the observation space and be rendered with a given camera view.

### 4.1. Non-rigid and Rigid Deformation

Inspired by [29, 38], we deform 3D Gaussians from canonical space $\mathcal{G}_c$ to observation space $\mathcal{G}_o$ by integrating a rigid articulation with a non-rigid transformation. We employ a non-rigid deformation network that takes the canonical positions $\mathcal{X}_c$ of the 3D Gaussians $\mathcal{G}c$ and a pose latent code which encodes SMPL pose $p_i$ using a lightweight hierarchical pose encoder [20]. The network then outputs the offsets for various parameters of the 3D Gaussians $\mathcal{G}_c$: $\Delta(\mathcal{X}, \mathcal{C}, \alpha, s, r)$. The canonical Gaussians are deformed by:

$$\mathcal{X}_d = \mathcal{X}_c + \Delta\mathcal{X}, \mathcal{C}_d = \mathcal{C}_c + \Delta\mathcal{C}, \quad (2)$$

$$\alpha_d = \alpha_c + \Delta\alpha, s_d = s_c \cdot \exp(\Delta s), \quad (3)$$

$$r_d = r_c \cdot [1, \Delta r_1, \Delta r_2, \Delta r_3], \quad (4)$$

where quaternion multiplication $\cdot$ corresponds to multiplying the rotation matrices. With $[1, 0, 0, 0]$ as the identity rotation, $r_d = r_c$ when $\delta r = \mathbf{0}$, thus keeping the original orientation.

We further apply a LBS-based rigid transformation to map the non-rigidly deformed 3D Gaussians $\mathcal{G}_d$ to the observation space $\mathcal{G}_o$. This transformation utilizes LBS weights predicted by a Skinning MLP $f_{\theta_r}$. This process aligns the Gaussians with the target pose in $\mathcal{G}_o$:

$$T = \sum_{j=1}^{J} f_{\theta_r}(\mathcal{X}_d)_j B_j, \mathcal{X}_o = T\mathcal{X}_d, \quad (5)$$

$$\mathcal{R}_o = T_{1:3,1:3}\mathcal{R}_d, \quad (6)$$

where $\mathcal{R}$ is the matrix representations of rotation.

### 4.2. Dynamic Avatar Adjustment

To address extremely sparse inputs, we leverage the intrinsic 3D consistency of human avatars across different
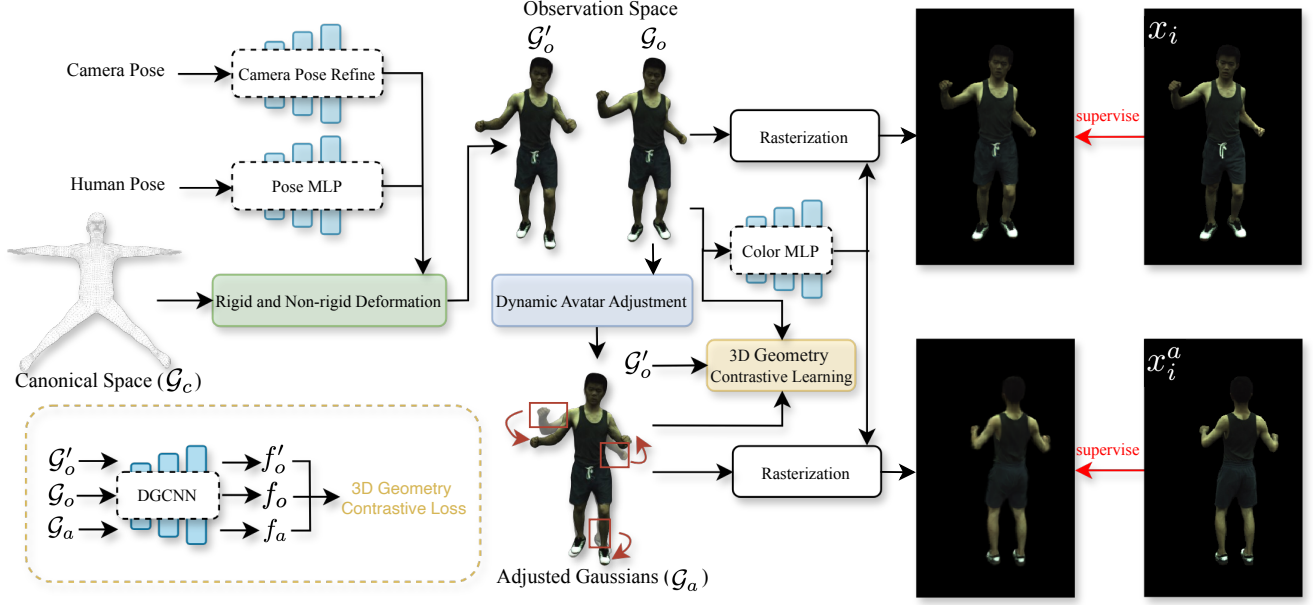
Figure 2. **CHASE Framework.** We first initialize 3D Gaussians in canonical space by randomly sampling 50k points on the SMPL mesh surface. Then, we integrate a rigid human articulation and a non-rigid deformation neural field to deform the 3D Gaussians in canonical space $\mathcal{G}_c$ to the observation space $\mathcal{G}_o$. Next, we select similar poses/images from the dataset for each training pose/image and then adjust the deformed Gaussians $\mathcal{G}_o$ to the similar pose $\mathcal{G}_a$ using Dynamic Avatar Adjustment (DAA). Minimizing the differences between the rendered adjusted Gaussians $\mathcal{G}_a$ and the selected similar images $x_i^a$ serves as an additional supervision. Furthermore, we propose a 3D geometry contrastive learning, which involves comparing features from a 3D feature extractor to improve the avatar's global 3D consistency. Negative pairs consist of the features of the deformed Gaussians $\mathcal{G}_o$ and the adjusted Gaussians $\mathcal{G}_a$. In contrast, positive pairs include the features of $\mathcal{G}_o'$, which is deformed from the canonical space to match the pose adjustments seen in $\mathcal{G}_a$, and $\mathcal{G}_a$.

poses/images, as shown in Fig. 3. Though the same pose may exhibit slight variations in non-rigid deformation, these differences occupy only a small number of pixels from a whole-body perspective and thus have minimal impact on the overall body reconstruction, which we further demonstrate in the experimental Section. 5.3.

Specifically, for each training pose/image, we select a similar pose $p_i^a$ with its paired image $x_i^a$ by computing the orientation and limb angle difference provided by SMPL [19] model from the dataset. Then we use a dense motion field $F_{adj}$ as an additional adjustment to transform deformed Gaussians $\mathcal{G}_o$ into adjusted Gaussians $\mathcal{G}_a$, aligning them with the selected pose/image $(p_i^a/x_i^a)$. In this way, we successfully introduce an additional 2D image supervision, improving the 3D consistency of human avatars.

To achieve precise control of the 3D Gaussians, we sample 6,890 points from the SMPL model [19] as our sparse control points in canonical space. Then, we obtain the dense motion field using LBS by locally inheriting the LBS weights from neighboring control points. Specifically, for each 3D Gaussian, we use the k-nearest neighbor (KNN) search to find its nearest neighboring control points

in canonical space. The entire adjustment process is as:

$$w = w_{smpl}[\text{KNN}(xyz_{cano}, xyz_{smpl})], \qquad (7)$$

$$T_o = \sum_{j=1}^{J} w_j B_{oj}, \quad T_o^{'} = \sum_{j=1}^{J} w_j B_{oj}^{'}. \qquad (8)$$

Here, $w_{smpl}$ denotes the LBS weights of the sparse control points, and $T_o(B_o)$ and $T_o'(B_o')$ represent the rigid transformations (bone transformations) from canonical space $\mathcal{G}_c$ to deformed Gaussians $\mathcal{G}_o$, and to the Gaussians with the selected similar pose $\mathcal{G}_o'$, respectively. We then obtain $F_{adj}$, which transforms the deformed Gaussians $\mathcal{G}_o$ into adjusted Gaussians $\mathcal{G}_a$, aligning them with the selected pose $p_i^a$ as:

$$F_{adj} = F_{o'} F_o^{-1}. \qquad (9)$$

We adjust the deformed Gaussians $\mathcal{G}_o$ to adjusted Gaussians $\mathcal{G}_a$ by by adjusting its position and rotation as follow:

$$\mathcal{X}_a = F_{adj} \mathcal{X}_o, \qquad (10)$$

$$\mathcal{R}_a = F_{adj1:3,1:3} \mathcal{R}_o. \qquad (11)$$

### 4.3. 3D Geometry Contrastive learning

Inspired by the success of contrastive learning in 2D image processing and static point cloud analysis, we advocate
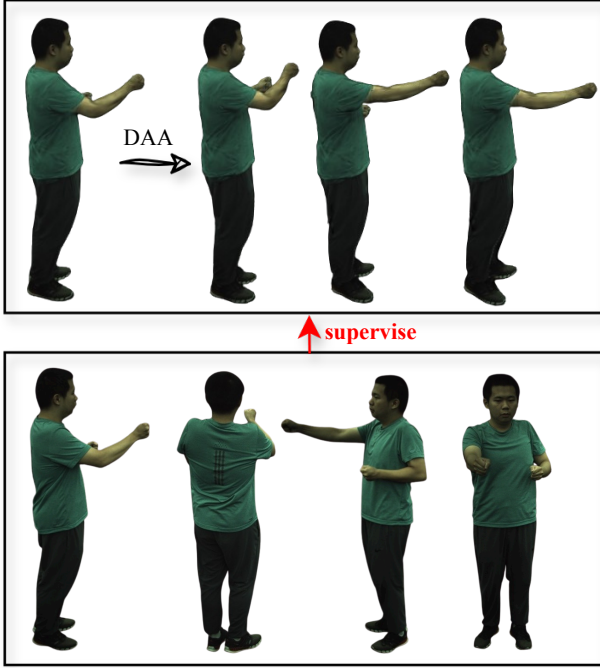
Figure 3. For each training pose/image, we select similar poses/images from the dataset and then adjust the deformed Gaussians using DAA. By minimizing the difference between the rendered image of the adjusted avatar and the selected similar pose image, we introduce additional supervision, thereby refining the creation of photo-realistic and animatable avatars.

for adopting 3D geometry contrastive learning to ensure 3D consistency of avatars. We treat the 3D Gaussians as a 3D point cloud and use DGCNN [36] as the feature extractor. DGCNN is typically trained on general point cloud datasets, which helps it learn geometric structures. The point cloud feature extractor processes the positions of the 3D Gaussians in the observation space $\mathcal{G}_o$, the adjusted Gaussians $\mathcal{G}_a$, and $\mathcal{G}_o'$, which is deformed from the canonical space to match the selected pose $p_i^a$, and outputs their features, creating intermediate graph features to capture global geometric information better. The feature vectors are projected into an invariant space. We denote the projected features of $\mathcal{G}_o$, $\mathcal{G}_a$, and $\mathcal{G}_o'$ as $f_o$, $f_a$, and $f_o'$, respectively.

In the invariant space, we aim to maximize the similarity between $f_a$ and $f_o'$, denoted as $D_{positive}$, and minimize the similarity between $f_a$ and $f_o$, denoted as $D_{negative}$. This feature-level contrastive learning encourages the model to capture pose-aware knowledge, enabling it to differentiate between subtle pose variations while maintaining geometric consistency. Therefore, we compute the 3D geometry

contrastive loss $\mathcal{L}_{contrastive}$ as:

$$D_{positive} = \|f_a - f_o'\|_2, \tag{12}$$
$$D_{negative} = \|f_a - f_o\|_2, \tag{13}$$
$$\mathcal{L}_{contrastive} = \max(0, D_{positive} - D_{negative}). \tag{14}$$

### 4.4. Optimization

We begin by randomly sampling 50k points from the surface of the SMPL mesh to initialize 3D Gaussians in the canonical space.

**Color MLP.** Following [29], we use the inverse rigid transformation to canonicalize the viewing direction: $\hat{d} = T_{1:3,1:3}^{-1}d$, where $T$ and $d$ is the forward transformation matrix defined in LBS and viewing direction, respectively. Theoretically, canonicalizing viewing direction also promotes consistency of the specular component of canonical 3D Gaussians under rigid transformations.

**Pose correction.** Following [29], SMPL [19] parameter fittings from images can be inaccurate. We additionally optimize the per-sequence shape parameter and per-frame translation, global rotation, and local joint rotations.

**Loss function.** Our full loss function consists of several components: an RGB loss $\mathcal{L}_{rgb}$, a mask loss $\mathcal{L}_{mask}$, and a perceptual similarity (LPIPS) loss $\mathcal{L}_{LPIPS}$. We compute these losses on both images rendered from the deformed Gaussians $\mathcal{G}_o$ and the adjusted Gaussians $\mathcal{G}_a$ with their corresponding ground truth images. Additionally, we include a skinning weight regularization loss $\mathcal{L}_{skin}$, as well as isometric regularization losses for both position and covariance, $\mathcal{L}_{isopos}$ and $\mathcal{L}_{isocov}$, following [29]. We also incorporate a 3D geometry contrastive loss $\mathcal{L}_{contrastive}$:

$$\mathcal{L} = \mathcal{L}_{rgb} + \lambda_1\mathcal{L}_{mask} + \lambda_2\mathcal{L}_{LPIPS} + \lambda_3\mathcal{L}_{skin} + $$
$$\lambda_4\mathcal{L}_{isopos} + \lambda_5\mathcal{L}_{isocov} + \lambda_6\mathcal{L}_{contrastive}, \tag{15}$$

where $\lambda$'s are loss weights. For further details of the loss definition and respective weights, please refer to the Supp.Mat.

## 5. Experiment

### 5.1. Dataset

**ZJU-MoCap [26].** This dataset features multi-view videos captured by 21 cameras, with human poses recorded using a marker-less motion capture system. For our experiments, we selected six sequences (377, 386, 387, 392, 393, 394). Following the protocol established by HumanNeRF [38] and 3DGS-Avatar [29], we use a single camera for training and the remaining cameras for evaluation. The foreground
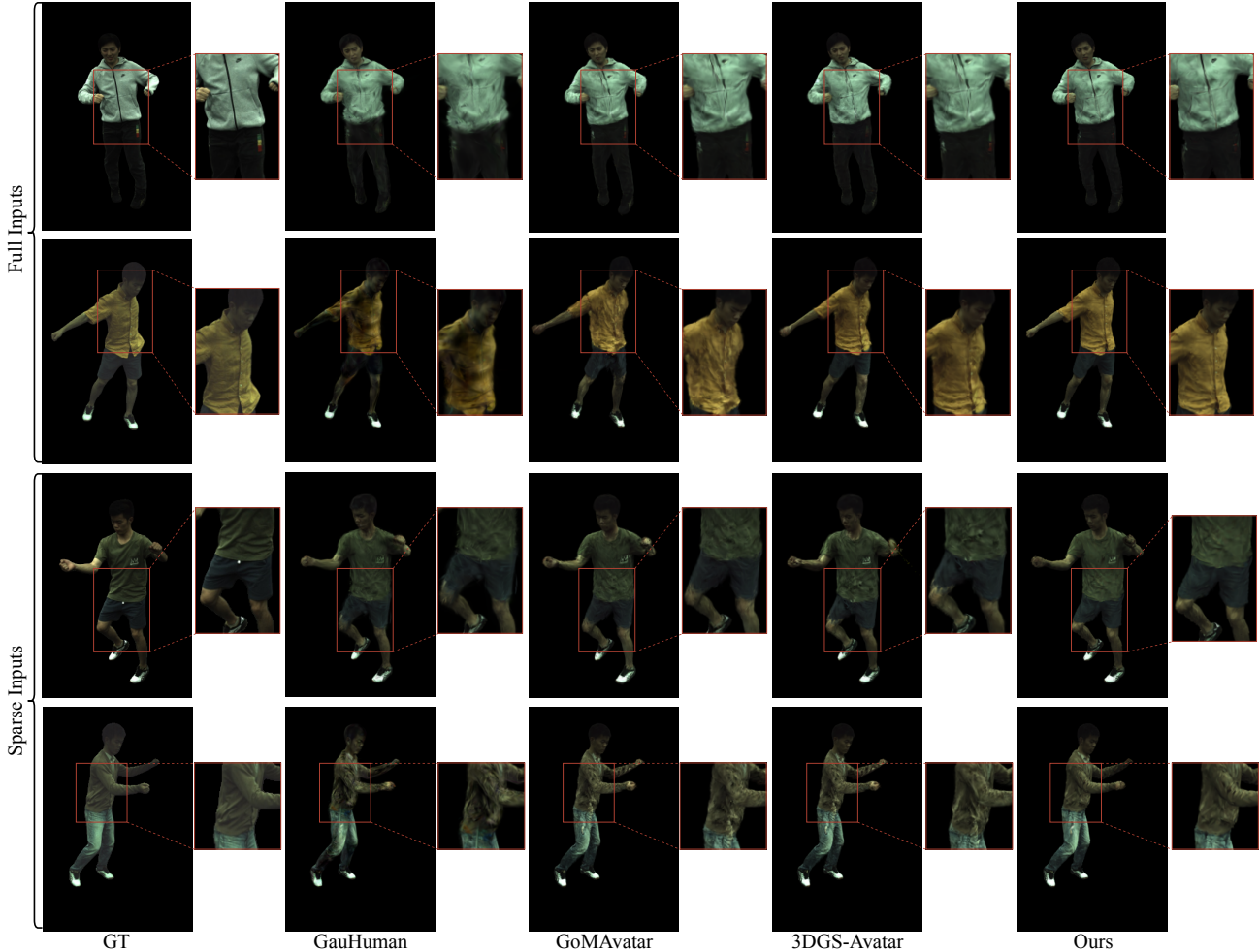
Figure 4. **Qualitative Comparison on ZJU-MoCap [26].** We present results for full and sparse inputs (5% of the full inputs) on the ZJU-MoCap dataset. Results show that our CHASE can produce realistic details with both full and sparse inputs, while other approaches struggle to generate smooth details.

masks, camera, and SMPL parameters provided by the data set are used for evaluation purposes. *We simulate sparse inputs by sampling every 20th frame from the video, between start_frame and end_frame, which corresponds to using 5% of the total frames.* We also specify in the Supp.Mat which images are selected for training.

**H36M [10].** H36M is another widely used dataset for human avatar research, comprising multi-view videos from four cameras and human poses captured via a marker-based motion capture system. We conducted experiments on sequences from subjects S1, S5, S6, S7, S8, S9, and S11, selecting representative actions and dividing the videos into training and test frames. Adhering to the protocol set by ARAH [35], we use three cameras, [54138969, 55011271, 58860488], for training and the remaining camera, [60457274], for testing, and follow their preprocessing steps. We use the SMPL parameters and foreground hu-

mans following [25]. We apply the same approach as mentioned above to simulate sparse inputs in this dataset.

## 5.2. Comparison with State-of-the-art Methods

We compare our CHASE with various SOTA methods for human avatars, including NerF-based methods such as NeuralBody [26], Ani-NeRF [25], HumanNeRF [38], and MonoHuman [44], and 3DGS-based methods such as 3DGS-Avatar [29], GauHuman [8], and GoMAvatar [37] under monocular setup on ZJU-MoCap [26]. The quantitative results are shown in Tab. 1. Overall, our proposed CHASE achieves the best performance in terms of PSNR, SSIM, and LPIPS *with both full and sparse inputs*. Notably, *our CHASE shows only a small performance drop when using only 5% of the data*, especially in LPIPS which is more informative than the other two metrics in our setting [29]. In fact, it declines less compared to other methods and even
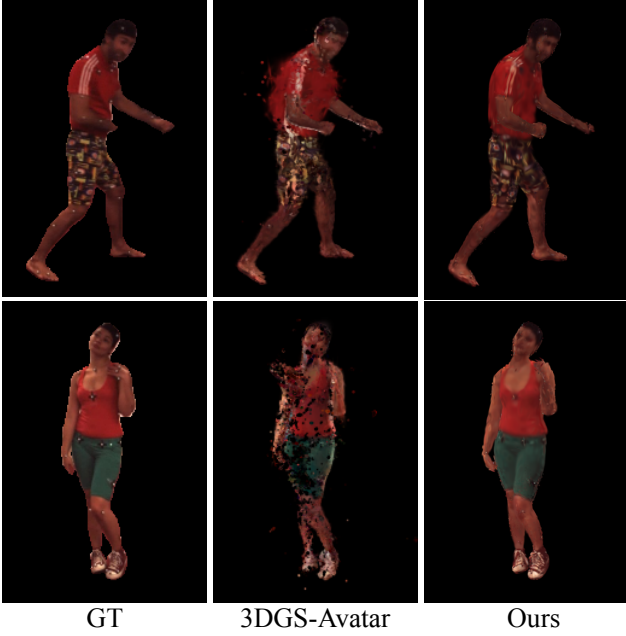
GT      3DGS-Avatar      Ours

Figure 5. **Qualitative Comparison on H36M [10] with sparse inputs.** We demonstrate that our method effectively produces realistic details for novel pose in both rendered images and geometry, whereas other approach struggles to achieve smooth details.

Table 1. **Quantitative Results on ZJU-MoCap [26].** CHASE achieves state-of-the-art performance across every method. The best and the second best results are denoted by pink and yellow. Frames per second (FPS) is measured on an RTX 3090. We train our model on the dataset that includes only 5% of the origin data for fair quantitative comparison. The metrics are reported in the last four rows of the table. LPIPS$^{\dagger}$ = LPIPS $\times$ 1000.

| Method: | LPIPS$^{\dagger}\downarrow$ | PSNR$\uparrow$ | SSIM$\uparrow$ | FPS |
|---|---|---|---|---|
| NeuralBody [26] | 52.29 | 29.07 | 0.962 | 1.5 |
| Ani-NeRF [25] | 51.98 | 29.17 | 0.961 | 1.1 |
| HumanNeRF [38] | 31.73 | 30.24 | 0.968 | 0.3 |
| MonoHuman [44] | 37.51 | 29.38 | 0.964 | 0.1 |
| DVA [31] | 37.74 | 29.45 | 0.956 | 17 |
| InstantAvatar [11] | 64.41 | 29.73 | 0.938 | 4.2 |
| 3DGS-Avatar [29] | 30.28 | 30.62 | 0.965 | 50 |
| GauHuman [8] | 32.73 | 30.79 | 0.960 | 180 |
| GoMAvatar [37] | 32.53 | 30.37 | 0.969 | 43 |
| Ours | 27.48 | 30.81 | 0.970 | 50 |
| 3DGS-Avatar* [29] | 40.01 | 29.98 | 0.957 | 50 |
| GauHuman* [8] | 35.68 | 30.35 | 0.957 | 180 |
| GoMAvatar* [37] | 42.88 | 30.01 | 0.958 | 43 |
| Ours* | 29.94 | 30.48 | 0.969 | 50 |

surpasses their performance with 100% of the data. It is evidence that our method successfully maintains 3D consistency even with sparse inputs.

Qualitative comparisons on novel view synthesis are

Table 2. **Quantitative Results on H36M [10].** Our CHASE outperforms current SOTA methods in both full and sparse settings

| Method: | Training Poses | | Novel Poses | |
|---|---|---|---|---|
| | PSNR$\uparrow$ | SSIM$\uparrow$ | PSNR$\uparrow$ | SSIM$\uparrow$ |
| NARF [24] | 23.00 | 0.898 | 22.27 | 0.881 |
| NeuralBody [26] | 22.89 | 0.896 | 23.09 | 0.891 |
| Ani-NeRF [25] | 23.00 | 0.890 | 22.55 | 0.880 |
| ARAH [35] | 24.79 | 0.918 | 23.42 | 0.896 |
| 3DGS-Avatar [29] | 32.89 | 0.982 | 32.50 | 0.983 |
| Ours | 33.29 | 0.984 | 32.93 | 0.982 |
| 3DGS-Avatar* [29] | 32.48 | 0.976 | 32.17 | 0.981 |
| Ours* | 32.91 | 0.983 | 32.64 | 0.982 |

Table 3. **Ablation Study on ZJU-MoCap [26].** We both show the result from full input (top group) and sparse input (bottom group).

| Method: | PSNR$\uparrow$ | SSIM$\uparrow$ | LPIPS$^{\dagger}\downarrow$ | FPS |
|---|---|---|---|---|
| w/o non-rigid | 30.32 | 0.968 | 30.41 | 50 |
| w/o contrastive | 30.76 | 0.970 | 27.58 | 50 |
| pointnet | 30.75 | 0.970 | 27.74 | 50 |
| w/o DAA | 30.78 | 0.970 | 27.83 | 50 |
| Top-3 | 30.80 | 0.970 | 27.73 | 50 |
| Top-5 | 30.79 | 0.970 | 27.74 | 50 |
| Top-10 | 30.76 | 0.970 | 27.82 | 50 |
| Full model | 30.81 | 0.970 | 27.48 | 50 |
| w/o con | 30.33 | 0.968 | 29.96 | 50 |
| w/o DAA | 30.42 | 0.968 | 30.17 | 50 |
| Top-3 | 30.44 | 0.969 | 29.98 | 50 |
| Top-5 | 30.43 | 0.968 | 30.19 | 50 |
| Full model | 30.48 | 0.969 | 29.94 | 50 |

shown in Fig. 4. We observe that our method preserves more details compared to other SOTA methods. They often struggle to maintain 3D consistency and deliver suboptimal detail reconstruction in human avatar modeling, particularly when only sparse inputs are available. Please see our project website and supplementary material for more visualization.

For H36M [10], we report the quantitative results against NeRF-based methods such as NARF [24], NeuralBody [26], Ani-NeRF [25] and ARAH [35], and 3DGS-based methods such as 3DGS-Avatar [29] in Tab. 2. Our CHASE significantly outperforms these methods with sparse inputs, showing that our CHASE generalizes well to novel poses with sparse inputs and reconstructs human avatars with better appearance and geometry detail. For qualitative comparisons on novel pose synthesis, as shown in Fig. 5, our method generalizes well to novel pose with just sparse inputs (only 5% of origin data) and reconstruct human avatars with better appearance and geometry detail.
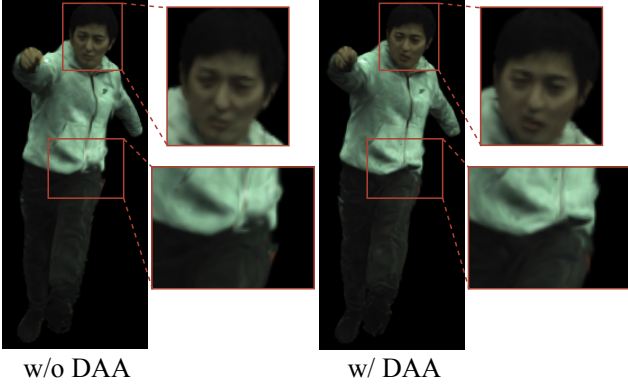
w/o DAA          w/ DAA

Figure 6. **Ablation Study** on DAA, which enhances multi-view 3D consistency, hence improving the overall rendering quality.
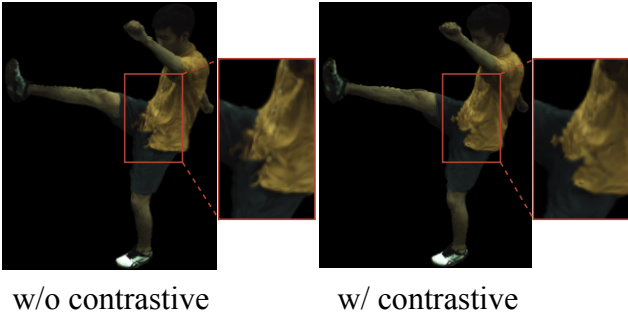


w/o contrastive      w/ contrastive

Figure 7. **Ablation Study** on 3D geometry contrastive learning, which removes the artifacts on highly articulated poses..

### 5.3. Ablation Study

In this section, we conduct ablation experiments using the ZJU-MoCap [26] dataset with both full inputs and sparse inputs to evaluate the effectiveness of our proposed modules. We also conduct experiments to evaluate different backbones in 3D geometry contrastive learning. *Notably, our CHASE maintains 3D consistency effectively without increasing any extra inference time.*

**Non-rigid deformation.** Non-rigid Deformation is designed for complex deformations, such as cloth bending and stretching. As shown in Tab. 3, non-rigid deformation is required to achieve optimal performance, demonstrating non-rigid regions are well rendered.

**Dynamic Avatar Adjustment.** As shown in Tab.3, incorporating DAA results in our full model outperforming the baseline in terms of LPIPS which is particularly informative compared to other metrics in our setting [29,40]. Fig. 6 illustrates that DAA serves as an effective 2D image supervision for 3D human body modeling, enhancing 3D consistency and reducing artifacts while improving multi-view consistency. Additionally, DAA reduces artifacts and inaccuracies in the geometry caused by sparse inputs, further

enhancing the robustness and visual realism of the model in practical applications.

**Similar (pose/image) selection.** Our method selects a similar (pose/image) to supervise the generated avatar as described in Section 4.2. Ablation studies in Tab. 3 show that even when selecting less similar (poses/images), our approach significantly improves performance. Top-n refers to choosing the n-th most similar pose/image. These results highlight CHASE 's effectiveness in handling limited data or irregular large-scale human movements.

**3D geometry contrastive learning.** Our core idea is that cross-3D-modal contrastive learning can facilitate communication between 3D models for obtaining powerful representations. To verify this, we further do ablation studies to show qualitative comparisons in Tab. 3 and Fig. 7. We can find the full model (w/ contrastive) preserves finer details and provides a more realistic and detailed reconstruction of clothing, demonstrating that 3D geometry contrastive learning enhances 3D consistency.

**Backbone for 3D contrastive learning.** In Tab. 3, we show the ablation study on different backbones, including Point-Net [27] (pointnet) and DGCNN [36] (full model). This indicates that DGCNN captures local geometric details better by building a dynamic graph structure, while PointNet relies on global feature learning, which may cause some local information to be missing.

## 6. Conclusion

In this paper, we present CHASE, a 3D-consistent human modeling framework utilizing Gaussian Splatting with both full and sparse inputs. We first integrates a skeleton-driven rigid deformation and a non-rigid cloth dynamics deformation to create human avatar. To improve 3D consistency under sparse inputs, we use the intrinsic 3D consistency of images across poses. For each training image, we select similar pose/image from the dataset and adjust the deformed Gaussians to selected pose by Dynamic Avatar Adjustment (DAA). Minimizing the difference between the image rendered by adjusted Gaussians and image paired with selected similar pose serves as an additional supervision, hence enhancing the 3D consistency of human avatars. Furthermore, to enforce global 3D consistency across different representations of the same pose, we propose a 3D geometry contrastive learning. Extensive experiments on two popular datasets demonstrate that CHASE not only achieves superior fidelity in generating human avatars compared to current SOTA methods but also excels in handling both monocular and sparse input scenarios. We hope that our method could foster further research in high-quality clothed human avatar synthesis from monocular views.

**Limitations.** 1). CHASE lacks the capability to extract 3D

meshes. Developing a method to extract meshes from 3D Gaussians is an important direction for future research. 2). CHASE performs less effectively in reconstructing humans with complex clothing, as similar poses may still exhibit substantial differences in non-rigid deformations. This remains a challenge we plan to address in future work.

# References

[1] Mohamed Afham, Isuru Dissanayake, Dinithi Dissanayake, Amaya Dharmasiri, Kanchana Thilakarathna, and Ranga Rodrigo. Crosspoint: Self-supervised cross-modal contrastive learning for 3d point cloud understanding. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 9902–9912, 2022. 2

[2] Mingfei Chen, Jianfeng Zhang, Xiangyu Xu, Lijuan Liu, Yujun Cai, Jiashi Feng, and Shuicheng Yan. Geometry-guided progressive nerf for generalizable and efficient neural human rendering. In *Proc. of European Conf. on Computer Vision*, pages 222–239, 2022. 2

[3] Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng Yang, Yikai Wang, Zhongang Cai, Lei Yang, Huaping Liu, and Guosheng Lin. Gaussianeditor: Swift and controllable 3d editing with gaussian splatting. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 21476–21485, 2024. 2

[4] Lin Gao, Yu-Kun Lai, Jie Yang, Ling-Xiao Zhang, Shihong Xia, and Leif Kobbelt. Sparse data driven mesh deformation. *IEEE transactions on visualization and computer graphics*, 27(3):2085–2100, 2019. 2

[5] Antoine Guédon and Vincent Lepetit. Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 5354–5363, 2024. 2

[6] Jennifer Healey, Wang, and et al. A mixed-reality system to promote child engagement in remote intergenerational storytelling. In *International Symposium on Mixed and Augmented Reality Adjunct*, pages 274–279, 2021. 1

[7] Liangxiao Hu, Hongwen Zhang, Yuxiang Zhang, Boyao Zhou, Boning Liu, Shengping Zhang, and Liqiang Nie. Gaussianavatar: Towards realistic human avatar modeling from a single video via animatable 3d gaussians. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 634–644, 2024. 2

[8] Shoukang Hu et al. GauHuman: Articulated gaussian splatting from monocular human videos. In *cvpr*, pages 20418–20431, 2024. 1, 2, 3, 6, 7

[9] Yi-Hua Huang, Yang-Tian Sun, Ziyi Yang, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. SC-GS: Sparse-controlled gaussian splatting for editable dynamic scenes. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 4220–4230, 2024. 2

[10] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6M: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Trans. on Pattern Anal. and Mach. Intell.*, 36(7):1325–1339, 2013. 2, 6, 7

[11] Tianjian Jiang, Xu Chen, Jie Song, and Otmar Hilliges. Instantavatar: Learning avatars from monocular video in 60 seconds. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 16922–16932, 2023. 2, 7

[12] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023. 2, 3

[13] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Proc. of Advances in Neural Information Processing Systems*, 33:18661–18673, 2020. 2

[14] Muhammed Kocabas, Jen-Hao Rick Chang, James Gabriel, Oncel Tuzel, and Anurag Ranjan. HUGS: Human gaussian splats. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 505–515, 2024. 2, 3

[15] Youngjoong Kwon, Baole Fang, Yixing Lu, Haoye Dong, Cheng Zhang, Francisco Vicente Carrasco, Albert Mosella-Montoro, Jianjin Xu, Shingo Takagi, Daeil Kim, et al. Generalizable human gaussians for sparse view synthesis. *arXiv preprint arXiv:2407.12777*, 2024. 2

[16] Jiahui Lei, Yufu Wang, Georgios Pavlakos, Lingjie Liu, and Kostas Daniilidis. GART: Gaussian articulated template models. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 19876–19887, 2024. 2, 3

[17] Zhe Li, Zerong Zheng, Yuxiao Liu, Boyao Zhou, and Yebin Liu. Posevocab: Learning joint-structured pose embeddings for human avatar modeling. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–11, 2023. 2

[18] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. *ACM transactions on graphics (TOG)*, 40(6):1–16, 2021. 2

[19] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *Acm Transactions on Graphics*, 34(248), 2015. 2, 3, 4, 5

[20] Marko Mihajlovic, Yan Zhang, Michael J Black, and Siyu Tang. LEAP: Learning articulated occupancy of people. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 10461–10471, 2021. 3

[21] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 2

[22] Arthur Moreau, Jifei Song, Helisa Dhamo, Richard Shaw, Yiren Zhou, and Eduardo Pérez-Pellitero. Human gaussian splatting: Real-time rendering of animatable avatars. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 788–798, 2024. 2, 3

[23] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 3504–3515, 2020. 1

[24] Atsuhiro Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. Neural articulated radiance field. In *Proc. of IEEE Intl. Conf. on Computer Vision*, pages 5762–5772, 2021. 7

[25] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *Proc. of IEEE Intl. Conf. on Computer Vision*, pages 14314–14323, 2021. 6, 7

[26] 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. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 9054–9063, 2021. 2, 5, 6, 7, 8

[27] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 652–660, 2017. 8

[28] Zekun Qi, Runpei Dong, Guofan Fan, Zheng Ge, Xiangyu Zhang, Kaisheng Ma, and Li Yi. Contrast with reconstruct: Contrastive 3d representation learning guided by generative pretraining. In *Proc. of Intl. Conf. on Machine Learning*, pages 28223–28243, 2023. 2

[29] Zhiyin Qian, Shaofei Wang, Marko Mihajlovic, Andreas Geiger, and Siyu Tang. 3dgs-avatar: Animatable avatars via deformable 3d gaussian splatting. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 5020–5030, 2024. 1, 2, 3, 5, 6, 7, 8

[30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proc. of Intl. Conf. on Machine Learning*, pages 8748–8763, 2021. 2

[31] Edoardo Remelli, Timur Bagautdinov, Shunsuke Saito, Chenglei Wu, Tomas Simon, Shih-En Wei, Kaiwen Guo, Zhe Cao, Fabian Prada, Jason Saragih, et al. Drivable volumetric avatars using texel-aligned features. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–9, 2022. 2, 7

[32] 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 *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 1606–1616, 2024. 2, 3

[33] Robert W Sumner, Johannes Schmid, and Mark Pauly. Embedded deformation for shape manipulation. In *ACM siggraph 2007 papers*, pages 80–es, 2007. 3

[34] Hongsheng Wang, Weiyue Zhang, Sihao Liu, Xinrui Zhou, Shengyu Zhang, Fei Wu, and Feng Lin. Gaussian control with hierarchical semantic graphs in 3d human recovery. *arXiv preprint arXiv:2405.12477*, 2024. 2

[35] Shaofei Wang, Katja Schwarz, Andreas Geiger, and Siyu Tang. Arah: Animatable volume rendering of articulated human sdfs. In *Proc. of European Conf. on Computer Vision*, pages 1–19, 2022. 6, 7

[36] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *ACM Transactions on Graphics (tog)*, 38(5):1–12, 2019. 5, 8

[37] Jing Wen, Xiaoming Zhao, Zhongzheng Ren, Alexander G Schwing, and Shenlong Wang. GomavAtar: Efficient animatable human modeling from monocular video using gaussians-on-mesh. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2059–2069, 2024. 1, 3, 6, 7

[38] 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 *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 16210–16220, 2022. 2, 3, 5, 6, 7

[39] Saining Xie, Jiatao Gu, Demi Guo, Charles R Qi, Leonidas Guibas, and Or Litany. Pointcontrast: Unsupervised pretraining for 3d point cloud understanding. In *Proc. of European Conf. on Computer Vision*, pages 574–591, 2020. 2

[40] Chen Yang, Sikuang Li, Jiemin Fang, Ruofan Liang, Lingxi Xie, Xiaopeng Zhang, Wei Shen, and Qi Tian. Gaussianobject: Just taking four images to get a high-quality 3d object with gaussian splatting. *arXiv preprint arXiv:2402.10259*, 2024. 8

[41] Wang Yifan, Noam Aigerman, Vladimir G Kim, Siddhartha Chaudhuri, and Olga Sorkine-Hornung. Neural cages for detail-preserving 3d deformations. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 75–83, 2020. 2

[42] Wang Yifan, Felice Serena, Shihao Wu, Cengiz Öztireli, and Olga Sorkine-Hornung. Differentiable surface splatting for point-based geometry processing. *ACM Transactions on Graphics*, 38(6):1–14, 2019. 3

[43] Yizhou Yu, Kun Zhou, Dong Xu, Xiaohan Shi, Hujun Bao, Baining Guo, and Heung-Yeung Shum. Mesh editing with poisson-based gradient field manipulation. In *ACM SIGGRAPH*, pages 644–651, 2004. 2

[44] Zhengming Yu, Wei Cheng, xian Liu, Wayne Wu, and Kwan-Yee Lin. MonoHuman: Animatable human neural field from monocular video. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 16943–16953, 2023. 6, 7

[45] Renrui Zhang, Ziyu Guo, Wei Zhang, Kunchang Li, Xupeng Miao, Bin Cui, Yu Qiao, Peng Gao, and Hongsheng Li. Pointclip: Point cloud understanding by clip. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 8552–8562, 2022. 2

[46] Fuqiang Zhao, Wei Yang, Jiakai Zhang, Pei Lin, Yingliang Zhang, Jingyi Yu, and Lan Xu. HumanNeRF: Efficiently generated human radiance field from sparse inputs. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 7743–7753, 2022. 2

[47] Haoyu Zhao, Xingyue Zhao, Lingting Zhu, Weixi Zheng, and Yongchao Xu. HFGS: 4d gaussian splatting with emphasis on spatial and temporal high-frequency components for endoscopic scene reconstruction. *arXiv preprint arXiv:2405.17872*, 2024. 2

[48] Yufeng Zheng, Wang Yifan, Gordon Wetzstein, Michael J Black, and Otmar Hilliges. Pointavatar: Deformable point-based head avatars from videos. In *Proc. of IEEE Conf.*

*on Computer Vision and Pattern Recognition*, pages 21057–21067, 2023. 2