# A Survey on 3D Human Avatar Modeling - From Reconstruction to Generation

Ruihe Wang[2*] Yukang Cao[1*†] Kai Han[1] Kwan-Yee K. Wong[1]

[1] The University of Hong Kong [2] Harbin Institute of Technology

## Abstract

*3D modeling has long been an important area in computer vision and computer graphics. Recently, thanks to the breakthroughs in neural representations and generative models, we witnessed a rapid development of 3D modeling. 3D human modeling, lying at the core of many real-world applications, such as gaming and animation, has attracted significant attention. Over the past few years, a large body of work on creating 3D human avatars has been introduced, forming a new and abundant knowledge base for 3D human modeling. The scale of the literature makes it difficult for individuals to keep track of all the works. This survey aims to provide a comprehensive overview of these emerging techniques for 3D human avatar modeling, from both reconstruction and generation perspectives. Firstly, we review representative methods for 3D human reconstruction, including methods based on pixel-aligned implicit function, neural radiance field, and 3D Gaussian Splatting, etc. We then summarize representative methods for 3D human generation, especially those using large language models like CLIP, diffusion models, and various 3D representations, which demonstrate state-of-the-art performance. Finally, we discuss our reflection on existing methods and open challenges for 3D human avatar modeling, shedding light on future research.*

## 1. Introduction

Human avatar modeling has recently shown significant scientific progress, with diverse applications ranging from computer graphics and gaming to virtual reality and medical imaging. While early methods rely on expensive capturing hardware and labor-intensive calibration processes to produce good-looking models [CCS*15], recent advancements have made it much more convenient to reconstruct and generate human avatars from various types of input such as images, videos, or text prompts.

3D human mesh reconstruction methods can roughly be categorized into model-based methods [ASK*05, BSB*07, BTTPM19, APMTM19] and model-free methods [CZ19, CPM20, DLJ*20, PNM*20]. Model-based methods involve fitting an explicit parametric human model (e.g., SMPL [LMR*15]) to an image but they encounter challenges in capturing intricate details like clothing and hair. In contrast, model-free methods overcome these limitations by predicting the occupancy values of a volumetric space. A representative method is PIFu [SHN*19], which exploits a Multi-Layer Perceptron (MLP) to model an implicit function that predicts the occupancy value of a given point by leveraging pixel-aligned features extracted from the input image. However, PIFu does not leverage the human body's structure, thus struggling with challenging poses, self-occlusions, and depth ambiguities. Follow-up works [SSSJ20, HCJS20, HZJ*21, ZYLD21, HXL*20, HXS*21,

CCH*22, XYTB22, CHW23] address these shortcomings by integrating priors such as normal maps, SMPL model, and depth information. However, these methods still face a topological constraint that restricts the model's performance with loose clothing, which is subsequently addressed by ECON [XYC*23] via an explicit method. Besides methods utilizing a single-view image as input, multi-view scenarios [ZSZ*21, SZZ*22b] offer a richer source of information from different viewpoints, leading to improved reconstruction results.

However, the performance of PIFu-based methods is largely limited by the quality of 3D training datasets, which are scarce and difficult to obtain. NeRF [MST*21] enables novel-view synthesis by simply inputting a limited set of images to obtain the RGB color and density value of each 3D point. Building upon NeRF, researchers proposed numerous 3D human reconstruction methods [WCS*22, PZX*21, KKCF21, LCY*23] that represent a 3D human as a neural radiance field without relying on prior knowledge or pre-trained models. Beyond reconstruction, the exploration of synthesizing free-viewpoint animations with user-controlled novel pose sequences is also a promising research direction [LHR*21, PDW*21, LTV*22, NSLH21, WSGT22, WCS*22]. Furthermore, methods [SJL*21, YZG*21, SZZ*22a] have also been proposed to integrate both surface and radiance fields to achieve high-fidelity 3D human novel view synthesis.

Yet, achieving high-quality reconstruction results using neural radiance fields still demands neural networks that are costly to

---

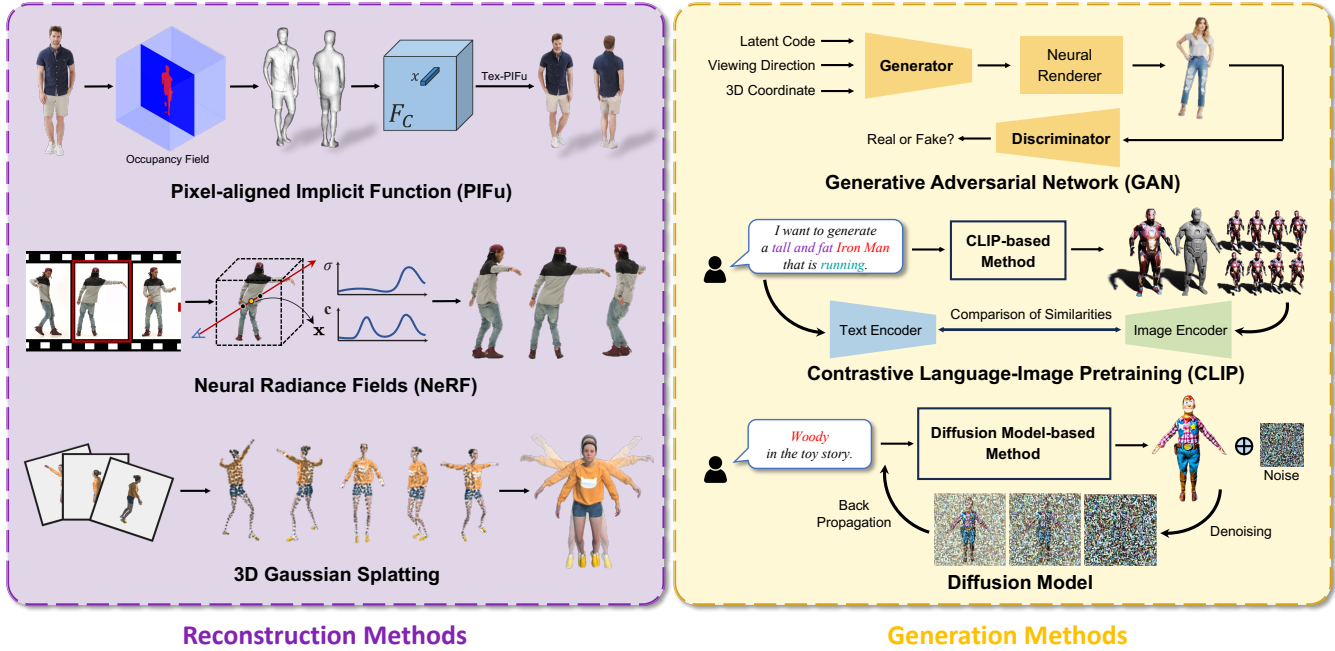† Corresponding author. * Equal contributions.

**Figure 1:** *An overview of typical 3D human avatar modeling approaches. Example images adapted from [SHN\*19, WCS\*22, HSY\*23, ZJY\*22, HZP\*22, CCH\*23a].*

train and render. 3D Gaussian Splatting (3DGS) [KKLD23] represents and renders complex scenes in a shorter training time without sacrificing speed for quality. By modeling a scene via a set of 3D Gaussians, Kerbl *et al.* [KKLD23] employ an explicit and object-centric method that is different from implicit representations like NeRF and DMTET [SGY\*21]. Following 3DGS, many works [LHQ\*23,ZBS\*23,LTYY23,XCL\*23,HSY\*23,HL23] have been proposed, utilizing its core principles for enhanced 3D human reconstruction, leading to highly animatable and realistic human models.

The advent of Generative Adversarial Networks (GANs) [GPAM\*14] marked an era of human avatar generation. GAN-based generation methods [BKY\*22, ZJY\*22, NSLH22, JJW\*22, HCL\*22] are typically composed of two key components: the StyleGAN [FLJ\*22] architecture and triplane representation proposed in EG3D [CLC\*22], successfully establishing a connection between 3D fields and 2D images. As a result, these methods have achieved remarkable progress in the field of 3D human avatar generation, despite being trained solely on 2D datasets.

While GAN-based methods have yielded impressive results, they still lack the ability to generate unseen characteristics that are not contained in the training dataset. With the recent development of large-language models, 3D generative methods [JMB\*22, MBOL\*22, JYQ\*22, HZP\*22] employ CLIP [RKH\*21] to generate 3D contents directly from text prompts. However, due to CLIP's limitation in fully comprehending textual descriptions, these methods still struggle to generate humans with fine details and complex motions. The advent of diffusion models has significantly ad-

vanced 3D generation by converting Gaussian noise into structured data via a Markov process with denoising steps. Drawing inspiration from DreamFusion's Score Distillation Sampling (SDS) technique [PJBM22], many methods [CCH\*23a, JWZ\*23, HWZ\*23, KAZ\*23,WWY23,LYX\*23,ZLJ\*23,HSZ\*23] are proposed to enhance generation quality from various perspectives, driving advancements of this field. Meanwhile, researchers have also introduced additional possibilities that extend the use of diffusion models to facilitate controllable 3D human editing [SSP\*23, HCH\*23, LWW\*23,PET\*23].

We provide an overview of 3D human modeling via various 3D representations in Fig. 1. In this paper, we present a taxonomy of recent research that maps out the evolution process of 3D human avatar modeling, categorizing it into five key areas: PIFu-based 3D implicit human reconstruction (Sec. 3.1), NeRF-based 3D human novel view synthesis (Sec. 4), 3D Gaussian-based methods (Sec. 5), GAN-based 3D human generation (Sec. 6), and language model-based 3D human generation and editing (Sec. 7). We further examine the models by breaking each category into more detailed subcategories (see Fig. 2 for an overview of this survey). In Sec. 8, we reflect upon existing methods and discuss the open challenges and potential directions for future research, focusing primarily on two categories of approaches: optimization-based methods and feedforward approaches.

In summary, this paper offers a comprehensive survey for 3D human avatar modeling. Particularly, we make the following contributions: Firstly, we present a thorough and up-to-date review of representative methods for 3D human reconstruction. Secondly, we summarize representative methods for the emerging techniques for
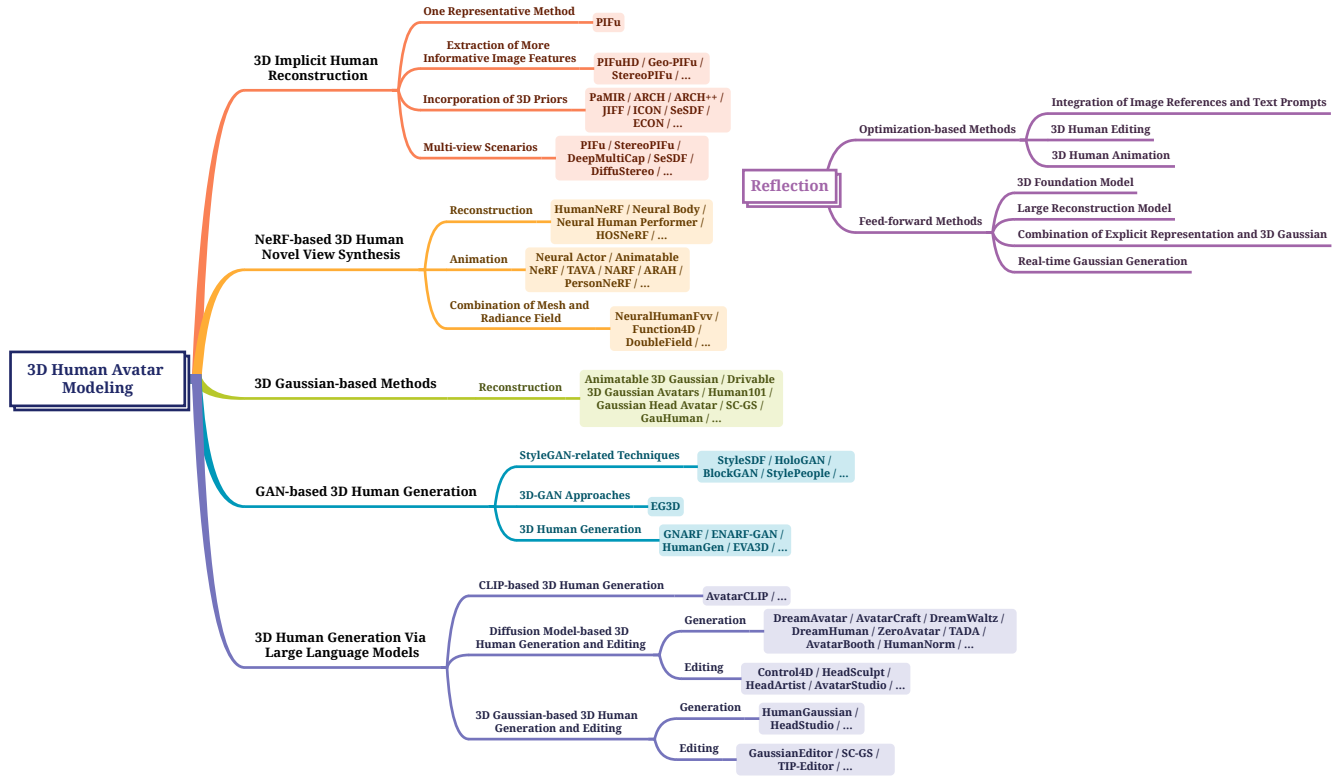
**3D Human Avatar Modeling**

- **3D Implicit Human Reconstruction**
  - One Representative Method → PIFu
  - Extraction of More Informative Image Features → PIFuHD / Geo-PIFu / StereoPIFu / ...
  - Incorporation of 3D Priors → PaMIR / ARCH / ARCH++ / JIFF / ICON / SeSDF / ECON / ...
  - Multi-view Scenarios → PIFu / StereoPIFu / DeepMultiCap / SeSDF / DiffuStereo / ...

- **NeRF-based 3D Human Novel View Synthesis**
  - Reconstruction → HumanNeRF / Neural Body / Neural Human Performer / HOSNeRF / ...
  - Animation → Neural Actor / Animatable NeRF / TAVA / NARF / ARAH / PersonNeRF / ...
  - Combination of Mesh and Radiance Field → NeuralHumanFvv / Function4D / DoubleField / ...

- **3D Gaussian-based Methods**
  - Reconstruction → Animatable 3D Gaussian / Drivable 3D Gaussian Avatars / Human101 / Gaussian Head Avatar / SC-GS / GauHuman / ...

- **GAN-based 3D Human Generation**
  - StyleGAN-related Techniques → StyleSDF / HoloGAN / BlockGAN / StylePeople / ...
  - 3D-GAN Approaches → EG3D
  - 3D Human Generation → GNARF / ENARF-GAN / HumanGen / EVA3D / ...

- **3D Human Generation Via Large Language Models**
  - CLIP-based 3D Human Generation → AvatarCLIP / ...
  - Diffusion Model-based 3D Human Generation and Editing
    - Generation → DreamAvatar / AvatarCraft / DreamWaltz / DreamHuman / ZeroAvatar / TADA / AvatarBooth / HumanNorm / ...
    - Editing → Control4D / HeadSculpt / HeadArtist / AvatarStudio / ...
  - 3D Gaussian-based 3D Human Generation and Editing
    - Generation → HumanGaussian / HeadStudio / ...
    - Editing → GaussianEditor / SC-GS / TIP-Editor / ...

- **Reflection**
  - Optimization-based Methods
    - Integration of Image References and Text Prompts
    - 3D Human Editing
    - 3D Human Animation
  - Feed-forward Methods
    - 3D Foundation Model
    - Large Reconstruction Model
    - Combination of Explicit Representation and 3D Gaussian
    - Real-time Gaussian Generation

**Figure 2:** *Taxonomy of 3D human avatar modeling methods in this survey.*

3D human generation. Finally, we provide our reflections on existing methods for 3D human avatar modeling and discuss insights and potential future research directions for future development of this field.

## 2. Scope of This Survey

This survey delves into recent advancements in 3D human modeling that utilize neural networks. Specifically, we start with illustrating the process of implicit function-based 3D human reconstruction from monocular images. Subsequently, we analyze the impact of the neural radiance field (NeRF) and 3D Gaussian Splatting on 3D human modeling. We then explore the realm of 3D generative AI, with a specific focus on the generative adversarial network (GAN), contrastive language-image pretraining (CLIP), and diffusion models. At the end of this survey, we provide our insights into future directions in this field.

This survey comprehensively discusses the essential techniques in 3D human modeling, covering 3D reconstruction, generation, and editing, to provide a detailed understanding of the past, present, and future. We have collected papers from major computer vision and computer graphics conferences and journals, as well as preprints available on arXiv. The selection process prioritized relevance to the scope of this survey, aiming to offer an inclusive overview of the rapid advances in this field. However, it is important to note that while this report serves as a compilation of state-of-

the-art methods in a specific domain, it is hard to have a complete coverage due to the vast number of publications and the rapid development of the field. Readers are encouraged to refer to the cited works for more in-depth discussions and additional methodologies.

**Related Survey.** Human mesh recovery (HMR) is fundamental to current 3D human modeling. However, due to space limit, we consider them to be beyond the scope of this report. Interested readers can find more comprehensive insights on HMR in [WTZ*21, ZWC*23, TZLW23]. 3D representations form the foundational basis for both 3D human modeling and the construction of 3D general objects. We encourage readers to learn from [GWH*20, HMZA21] for the details of point clouds, [TFT*20, TTM*22] for NeRF, and [CW24, FXZ*24, WYZ*24] for the most recent 3D Gaussian Splatting. Furthermore, the advancement of large language models has spurred research in another facet of the 3D virtual realm, namely, 3D general objects. Recent progress in 3D generative approaches can be explored through [YZS*23, CTG*24, LZW*23, LZK*24]. We also recommend consulting [CHIS23, XFC*23, PYG*23] for insights into 2D generative AI.

## 3. 3D implicit human reconstruction

In this section, we first discuss methods for 3D human reconstruction based on implicit function, which can accurately capture clothing topologies from single or multiple images. We start with pixel-aligned implicit function (PIFu), which is one of the most repre-
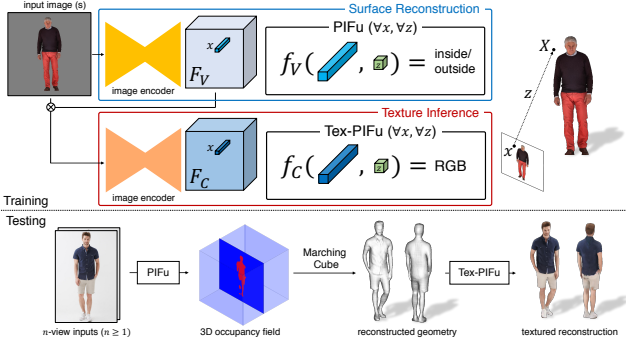
**Figure 3:** *An overview of PIFu network for predicting surface and texture values. Figure obtained from [SHN\*19].*

sentative works in this regime, in Sec. 3.1. We then explore (1) enhancements at the feature level in Sec. 3.2, (2) the usage of 3D priors in Sec. 3.3, and (3) multi-view scenarios in Sec. 3.4.

### 3.1. Pixel-aligned Implicit Function

Given a single image or multiple images as input, PIFu [SHN\*19] applies an hour-glassed image encoder [NYD16] and bi-linear interpolation to extract pixel-aligned image features for each 3D query point $\mathbf{x} \in \mathbb{R}^3$ and then predicts a continuous 3D occupancy field via an implicit function (see Fig. 3). Specifically, the implicit function, implemented as an MLP, defines the 3D human surface as the level set of the function, e.g., $F(\mathbf{x}) = 0.5$. It can be written as:

$$F_V(\mathcal{B}(\psi_{\text{geo}}(I), \pi(\mathbf{x})), z(\mathbf{x})) \mapsto s : s \in \mathbb{R}, \quad (1)$$

where $\mathcal{B}(\psi_{\text{geo}}(I), \pi(\mathbf{x}))$ is the pixel-aligned feature, with $\psi_{\text{geo}}(I) \in \mathbb{R}^{128 \times 128 \times 256}$ denoting the feature map extracted from the image $I \in \mathbb{R}^{512 \times 512 \times 3}$. $\pi(\mathbf{x})$ is the 2D projection of $\mathbf{x}$ on the image, and $z(\mathbf{x})$ is the depth value.

After obtaining the geometry in form of a mesh by applying Marching Cube [LC98] to the occupancy field, PIFu can be extended to texture inference by predicting the RGB value $\mathbf{c}$ for each 3D vertices:

$$F_C(\mathcal{B}(\psi_{\text{tex}}(I, f_{\text{geo}}), \pi(\mathbf{x})), z(\mathbf{x})) \mapsto \mathbf{c} : \mathbf{c} \in \mathbb{R}^3, \quad (2)$$

where $f_{\text{geo}}$ is the geometry feature obtained from the geometry inference stage.

Besides single-view 3D human reconstruction, PIFu can also take calibrated multi-view images as inputs (suppose we have $n$-view images as inputs). For each 3D query point $\mathbf{x}$, we first obtain the image features $\Phi_{\text{view}_i} = \mathcal{B}(\psi(I_i), \pi_i(\mathbf{x}))$ from the $i$-th view. The average feature $\text{avg}\{\cdot\}$ across all the views is utilized to enhance prediction accuracy for both surface and texture:

$$F(\text{avg}\{\Phi_{\text{view}_1}, \dots, \Phi_{\text{view}_n}\}) \mapsto s. \quad (3)$$

While PIFu shows potential in 3D human reconstruction via implicit functions, it encounters two major challenges: (1) The reconstructed 3D mesh still shows sub-optimal geometry details for real-world applications. (2) Due to the lack of 3D information under

single-view scenarios, PIFu falls short in handling complex poses and self-occlusions, and suffers from depth ambiguity issues. Subsequent methods aim to overcome these limitations by extracting more informative image features and incorporating 3D priors.

### 3.2. Extraction of More Informative Image Features

**PIFuHD** To enhance the quality of clothing topologies, PIFuHD [SSSJ20] introduces a coarse-to-fine network that integrates both low- and high-resolution image feature maps as well as the predicted front- and back-side normal images. Different from PIFu which takes low-resolution images $I_L \in \mathbb{R}^{512 \times 512 \times 3}$ as input, PIFuHD opts for information from high-resolution images, i.e., $I_H \in \mathbb{R}^{1024 \times 1024 \times 3}$. At the coarse stage, PIFuHD modifies PIFu by concatenating the predicted front and back normal maps ($\mathbf{n}_{L,\text{front}}$ and $\mathbf{n}_{L,\text{back}}$) with the downsampled low-resolution image $I_L$ as inputs:

$$F_L(\mathcal{B}(\psi(I_L \oplus \mathbf{n}_{L,\text{front}} \oplus \mathbf{n}_{L,\text{back}}), \pi(\mathbf{x})), z(\mathbf{x}) \mapsto s_L, \quad (4)$$

where $\oplus$ denotes the concatenation operation. PIFuHD then operates over the original high-resolution input at the fine stage:

$$F_H(\mathcal{B}(\psi(I_H \oplus \mathbf{n}_{H,\text{front}} \oplus \mathbf{n}_{H,\text{back}}), \pi(\mathbf{x})), \Omega(\mathbf{x}) \mapsto s_H, \quad (5)$$

where $\Omega(X)$ denotes the intermediate output from the penultimate layer of the MLP in the coarse branch.

**Geo-PIFu** Geo-PIFu [HCJS20] proposes to extract both 3D space-aligned and 2D pixel-aligned features from the input image, effectively encoding the structural information for occupancy estimation:

$$F(\mathcal{T}(\xi(I), \mathbf{x}), \mathcal{B}(\psi(I), \pi(\mathbf{x})), z(\mathbf{x}))) \mapsto s, \quad (6)$$

where $\xi$ is the 3D U-Net [JBAT17], with $\mathcal{T}(\cdot)$ denoting the multi-scale tri-linear interpolation. However, Geo-PIFu still struggles with capturing high-resolution clothing details since the information employed is extracted only from the image.

**StereoPIFu** StereoPIFu [HZJ\*21] targets at reconstructing 3D implicit human models from binocular images. Specifically, given the multi-view stereo inputs, StereoPIFu proposes a stereo vision-based network to extract voxel-aligned features. This geometric prior, together with a novel relative z-offset and depth maps, are utilized to produce enhanced reconstruction results:

$$F(\mathcal{B}(\psi_l(I), \pi_l(\mathbf{x})), \Phi(\mathbf{x}), \Psi(\mathbf{x}), \eta(Z_E(\mathbf{x}))) \mapsto s, \quad (7)$$

where $\mathcal{B}(\psi_l(I), \pi_l(\mathbf{x}))$ denotes pixel-aligned features from the left image, $\Phi(\mathbf{x})$ and $\Psi(\mathbf{x})$ are voxel-aligned features, and $\eta(\cdot)$ is a transformation function designed to normalize the relative z-offset $Z_E$ to the interval $(-1.0, 1.0)$::

$$Z_E = z(\mathbf{x}) - E(\pi_l(\mathbf{x})), \quad (8)$$

where $E(\cdot)$ is the predicted depth map of the left image.

### 3.3. Incorporation of 3D Priors

Attempts have also been made to incorporate 3D features as priors to enhance robustness, with SMPL [LMR\*15] and SMPL-X [PCG\*19] being frequently employed. Specifically, given the

pose parameter θ and shape parameter β, SMPL can map the canonical model with $n_S$ vertices to observation space[†]:

$$M(\beta, \theta) = \text{lbs}(T(\beta, \theta), J(\beta), \theta, \mathcal{W}), \tag{9a}$$

$$T(\beta, \theta) = \mathbf{T} + B_s(\beta) + B_p(\theta), \tag{9b}$$

where $M$ is the function representing the SMPL model in the observation space, and $T$ gives the transformed vertices. $\mathcal{W}$ is the blend weight, $B_s$ and $B_p$ are the shape blend shape function and pose blend shape function, respectively. $\text{lbs}(\cdot)$ denotes the linear blend skinning function, corresponding to articulated deformation. It poses $T(\cdot)$ based on the pose parameters θ and joint locations $J(\beta)$, using the blend weights $\mathcal{W}$, individually for each body vertex:

$$\mathbf{v}_o = \mathcal{G} \cdot \mathbf{v}_c, \quad \mathcal{G} = \sum_{k=1}^{K} w_k \mathcal{G}_k(\theta, j_k), \tag{10}$$

where $\mathbf{v}_c$ and $\mathbf{v}_o$ respectively are SMPL vertices under the canonical pose and observation space, $w_k$ is the skinning weight, $\mathcal{G}_k(\theta, j_k)$ is the affine deformation that transforms the $k$-th joint $j_k$ from the canonical space to observation space, and $K$ is the number of neighboring joints.

SMPL-X evolves from SMPL to include more face vertices, expression parameters φ and the expression blend shape function $B_e$ into the model:

$$M(\beta, \theta, \phi) = \text{lbs}(T(\beta, \theta, \phi), J(\beta), \theta, \mathcal{W}), \tag{11a}$$

$$T(\beta, \theta, \phi) = \mathbf{T} + B_s(\beta) + B_e(\phi) + B_p(\theta). \tag{11b}$$

**PaMIR** PaMIR [ZYLD21] is one of the first to extract 3D voxel-aligned features and semantic information from SMPL to provide a 3D-aware prior and address complex poses and self-occlusion issues. Specifically, PaMIR first estimates an initial SMPL model from the input image $I$ via a pre-trained GCMR network [KPD19]. For each 3D query point $\mathbf{x}$, its occupancy value is then predicted based on both pixel-aligned and voxel-aligned features:

$$F(\mathcal{B}(\psi(I), \pi(\mathbf{x})), \mathcal{T}(f_{3D}, \mathbf{x})) \mapsto s, \tag{12}$$

where $\mathcal{T}(f_{3D}, \mathbf{x})$ represents the voxel-aligned feature. The 3D feature volume $f_{3D}$ is derived by subsequently (1) converting the SMPL mesh into an occupancy volume $V_s$ through mesh voxelization; (2) encoding the occupancy volume through a 3D encoder $E_{3D}$, i.e., $f_{3D} = E_{3D}(V_s)$. However, PaMIR still performs poorly in reconstructing dynamic and complex movements when hands and clothes are close to each other.

**ARCH** ARCH [HXL*20], on the other hand, utilizes the Semantic Space (SemS) and the Semantic Deformation Field (SemDF) to transform query points from the observation space to the canonical space before calculating the occupancy value. It also propose to extract the spatial feature based on SemS via Radial Basis Function (RBF). An occupancy sub-network $F_s$, a normal sub-network $F_{\mathbf{n}}$, and a color sub-network $F_{\mathbf{c}}$ are then utilized for implicit surface

reconstruction in the canonical space:

$$F_s(\mathcal{B}(f_{2D}, \pi(\mathbf{x})), \mathcal{T}(f_{3D}, \mathbf{x})) \mapsto s, \tag{13a}$$

$$F_{\mathbf{n}}(\mathcal{B}(f_{2D}, \pi(\mathbf{x})), \mathcal{T}(f_{3D}, \mathbf{x}), f_s) \mapsto \mathbf{n}, \tag{13b}$$

$$F_{\mathbf{c}}(\mathcal{B}(f_{2D}, \pi(\mathbf{x})), \mathcal{T}(f_{3D}, \mathbf{x}), f_s, f_{\mathbf{n}}) \mapsto \mathbf{c}, \tag{13c}$$

where $f_{2D}$ is the 2D feature map, $f_{3D}$ is the 3D feature volume, $f_s$ and $f_{\mathbf{n}}$ are feature maps extracted from occupancy and normal sub-networks. Despite adding more information from various poses, ARCH can only reconstruct 3D humans under canonical space. Warping the mesh back to the observation space often results in artifacts such as intersecting surfaces and distorted body parts, leading to the degeneration of the reconstruction fidelity.

**ARCH++** Adopting the deformation field from ARCH, ARCH++ [HXS*21] proposes to learn the joint-space occupancy field in both observation and canonical spaces. It first transforms the posed mesh to the canonical space, and then uniformly samples query points on the mesh surface. For each query point $\mathbf{x}_c$, ARCH++ applies tri-linear interpolation to obtain spatial-aligned features $f_{3D}$, which are encoded by the PointNet++ [QSMG17, QYSG17]. The occupancy value is then jointly predicted in both the observation and canonical spaces to obtain additional constraints on the cross-space consistency:

$$F_o(\mathcal{B}(f_{2D}, \pi(\mathbf{x})), \mathcal{T}(f_{3D}, \mathbf{x})) \mapsto s_o, \tag{14a}$$

$$F_c(\mathcal{B}(f_{2D}, \pi(\mathbf{x})), \mathcal{T}(f_{3D}, \mathbf{x})) \mapsto s_c, \tag{14b}$$

where $s_o$, and $s_c$ stand for the occupancy value respectively in the observation and canonical space.

**JIFF** JIFF [CCH*22] leverages the 3DMM [BV99] as a 3D face prior to extract space-aligned 3D features with detailed geometry and texture information for improving the face quality. Given an input image $I$, JIFF first crops the face region and fits a 3DMM mesh $\mathbf{S}$ based on the cropped image. An encoder, which takes the vertices of 3DMM as input, is then employed to generate the 3D feature volume $\varphi(\mathbf{S})$. Thus, JIFF employs both pixel-aligned 2D features and space-aligned 3D features to complete the reconstruction:

$$F(\mathcal{T}(\varphi(\mathbf{S}), \mathbf{x}), \mathcal{B}(\psi(I), \pi(\mathbf{x})), z(\mathbf{x})) \mapsto s. \tag{15}$$

**ICON** Despite incorporating 3D-aware prior into the implicit function, previous methods still fall short in processing complex poses. To this end, ICON [XYTB22] replaces the global encoder with a more data-efficient local scheme. Specifically, given an image $I$ as input, ICON first uses a normal network $\mathcal{G} = (\mathcal{G}_{\text{front}}, \mathcal{G}_{\text{back}})$ to predict clothed body normal maps $\widehat{\mathcal{N}} = \{\widehat{\mathcal{N}}_{\text{front}}, \widehat{\mathcal{N}}_{\text{back}}\}$ based on $I$ and SMPL's front and back normal renderings:

$$\mathcal{G}(\mathcal{R}_{\mathbf{n}}(\text{SMPL}), I) \to \widehat{\mathcal{N}}, \tag{16}$$

where $\mathcal{R}_{\mathbf{n}}(\cdot)$ is a PyTorch3D [RRN*20] differentiable renderer. An implicit representation of the surface is then regressed based on these local features:

$$F(d(\mathbf{x}), \mathbf{n}(\mathbf{x}), \mathcal{N}(\pi(\mathbf{x}))) \mapsto s, \tag{17}$$

in which $d(\cdot)$ is the signed distance from a point $\mathbf{x}$ to closest vertex $\mathbf{x}'$ of SMPL, $\mathbf{n}(\cdot)$ is the barycentric surface normal of $\mathbf{x}'$, and $\mathcal{N}(\cdot)$

---

[†] In this survey, we refer to the space with poses of the SMPL model that differ from canonical pose as the observation space.

is a normal vector extracted from $\widehat{\mathcal{N}}^{c}$ depending on the visibility of $\mathbf{x}'$:

$$\mathcal{N}(\pi(\mathbf{x})) = \begin{cases} \mathcal{N}_{\text{front}}(\pi(\mathbf{x})), & \text{if } \mathbf{x}' \text{ is visible}, \\ \mathcal{N}_{\text{back}}(\pi(\mathbf{x})), & \text{otherwise.} \end{cases} \quad (18)$$

**SeSDF** SeSDF [CHW23] aims to flexibly and robustly extract detailed information of clothed humans from either a single image or uncalibrated multi-view images. Specifically, it proposes a self-evolved signed distance module (SeSDF), which refines the signed distance field derived from SMPL-X using both 2D pixel-aligned and space-aligned 3D image features. This approach enhances the signed distance field with additional clothing information that is consistent with the image features:

$$F_{\text{sdf}}\left(\mathcal{T}(\varphi(S), \mathbf{x}), \mathcal{B}(\psi(I), \pi(\mathbf{x})), \mathcal{D}(d(\mathbf{x})), \mathbf{n}(\mathbf{x})\right) \mapsto (d', \mathbf{n}'), \quad (19)$$

$$\mathcal{D}(d) = (d, \sin(2^0\pi d), \cos(2^0\pi d), \ldots, \sin(2^L\pi d), \cos(2^L\pi d)), \quad (20)$$

where $d' \in \mathbb{R}$ denotes the SDF derived from the SMPL-X model. Afterwards, given a 3D point $X$ and its features, the implicit function of SeSDF can be formulated as:

$$F\left(\mathcal{T}(\varphi(S), \mathbf{x}), \mathcal{B}(\psi(I), \pi(\mathbf{x})), \mathcal{D}(d'(\mathbf{x})), \mathbf{n}'(\mathbf{x}), z(\mathbf{x})\right) \mapsto s. \quad (21)$$

**ECON** Summarizing the above implicit methods reveals two challenges: using only image information can cause depth ambiguity and human pose inaccuracy, while incorporating 3D priors often results in missing clothing details as SMPL or SMPL-X contains only minimal clothes. Consequently, ECON [XYC*23] introduces an explicit method instead of the implicit function. It first follows ICON to predict 2D front and back normal and depth maps from the input image and the predicted SMPL-X model. It then uses a depth-aware silhouette-consistent bilateral normal integration (d-BiNI) optimizer [CSS*22] to recover the 3D front and back surfaces separately. Based on these partial surface estimates, it applies IF-Nets+ [CPM20] to implicitly complete the body. With optional face or hands from the SMPL-X model, screened poisson [KH13] is finally employed to combine all the 3D surfaces to form the complete 3D human body.

Other works including IP-Net [BSTPM20], S-PIFu [CLZL22a], and IntegratedPIFu [CLZL22b] also incorporate priors like SMPL model and depth for better human reconstruction, which readers can refer to on their own for detailed illustration.

### 3.3.1. Qualitative and Quantitative Evaluations

We first select typical implicit function-based methods that use only 2D features (PIFu [SHN*19], PIFuHD [SSSJ20]), as well as PaMIR [ZYLD21], ICON [XYTB22] that utilize 3D features, and the explicit method ECON [XYC*23] for quantitative evaluations. By referring to Tab. 1, we can observe that the incorporation of 3D priors leads to largely enhanced accuracy in capturing the human poses, resulting in much lower Chamfer, Point-to-Surface (P2S) and normal errors for both CAPE [MYR*20] and RenderPeople [Ren] test sets. We further provide qualitative evaluations in Fig. 4. Due to the reliance on the SMPL/SMPL-X model, ICON and PaMIR present degenerated clothing details than PI-FuHD which includes only 2D image features.
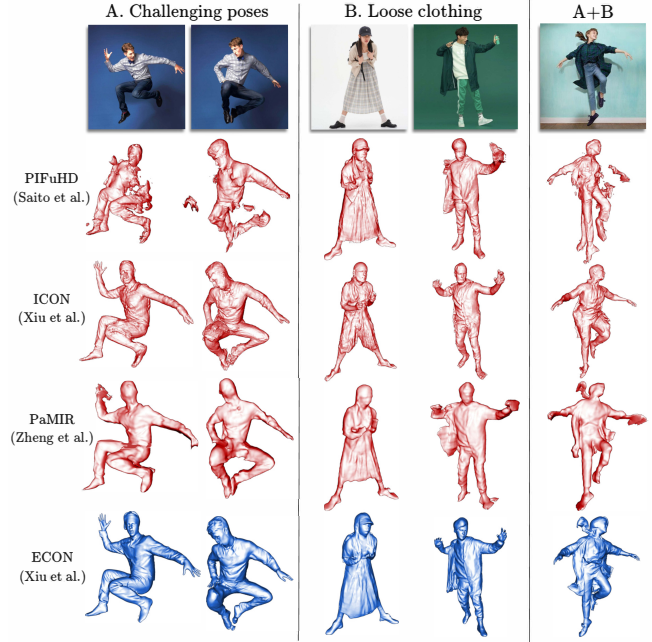


**Figure 4:** *Qualitative comparison on the single-view setting. Figure obtained from [XYC*23].*

Among them, ECON [XYC*23] introduces an explicit method that is capable of maintaining the robustness of explicit shape models for unseen poses without sacrificing the topological flexibility of implicit functions for loose clothing. However, it is important to note that single-view image only provides information about the front view of the human body, which inherently limits the effectiveness of these reconstruction methods.
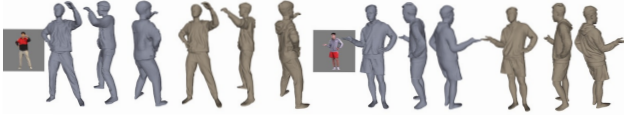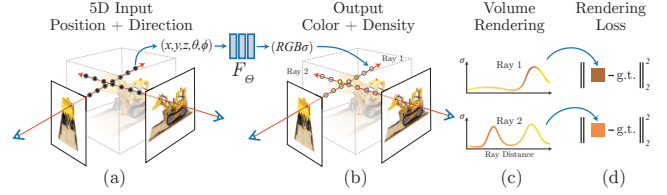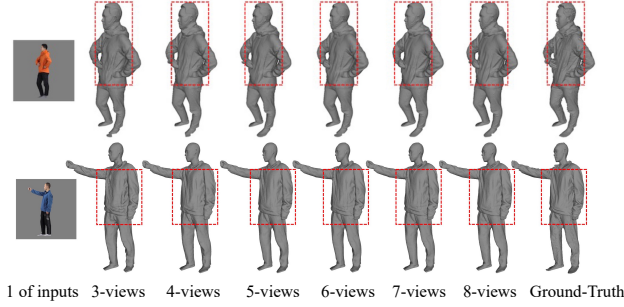
### 3.4. Multi-view Scenarios

Multi-view images can offer a more comprehensive set of 3D information that can compensate each other for getting better reconstruction results in terms of (1) finer clothing details, and (2) enhanced accuracy in capturing motion. However, the key challenge lies in effectively integrating features from different viewpoints, especially when occlusions occur in multiple views or when information for the same 3D query point varies across viewpoints.

PIFu [SHN*19] incorporates multi-view features through average pooling that treats multi-view features equally. However, this approach overlooks the different quality of predictions from different views. For example, the image feature of a 3D point extracted from a non-occluded view likely yields the most accurate prediction, while features extracted from an occluded or lateral view should have minimal impact on the prediction. Obviously, average pooling is not an optimal method for feature fusion.

Researchers therefore explore different feature fusion strategies. StereoPIFu [HZJ*21] considers a pair of stereo images as input for depth-aware reconstruction, DeepMultiCap [ZSZ*21] leverages a self-attention mechanism for multi-view fusion, SeSDF [CHW23] proposes an occlusion-aware feature fusion strategy to fuse features

**Table 1:** *Quantitative comparison on the single-view setting. Results obtained from [XYC\*23].*

| Method | Feature Included | | | | | Quantitative Number | | | | | |
| | | | | | | CAPE [MYR\*20] | | | RenderPeople [Ren] | | |
| | 2D Feature | 3D Feature | SMPL(-X) | Normal | Method Type | Chamfer↓ | P2S↓ | Normals↓ | Chmafer↓ | P2S↓ | Normals↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| PIFu [SHN\*19] | ✓ | | | | implicit | 1.722 | 1.548 | 0.0674 | 1.706 | 1.642 | 0.0709 |
| PIFuHD [SSSJ20] | ✓ | | | | implicit | 3.767 | 3.591 | 0.0994 | 1.946 | 1.983 | 0.0658 |
| PaMIR [ZYLD21] | ✓ | ✓ | ✓ | ✓ | implicit | 0.989 | 0.992 | 0.0422 | 1.296 | 1.430 | 0.0518 |
| ICON [XYTB22] | ✓ | | ✓ | ✓ | implicit | 0.971 | 0.909 | 0.0409 | 1.373 | 1.522 | 0.0566 |
| ECON [XYC\*23] | | | ✓ | ✓ | explicit | 0.996 | 0.967 | 0.0413 | 1.401 | 1.422 | 0.0516 |



**Figure 5:** *Single-view reconstruction (grey) vs multi-view reconstruction (yellow). Figure obtained from [CHW23].*



**Figure 7:** *NeRF overview. Figure obtained from [MST\*21].*



1 of inputs   3-views   4-views   5-views   6-views   7-views   8-views   Ground-Truth

**Figure 6:** *Reconstruction with different numbers of input views. Figure obtained from [CHW23] and model is trained with three views.*

from different views effectively, DiffuStereo [SZZ\*22b] introduces a multi-level diffusion-based stereo network to produce highly accurate depth maps, which are then converted into a high-quality 3D human model through an efficient multi-view fusion strategy.

### 3.4.1. Qualitative and Quantitative Evaluations

We begin by presenting a comparison between single-view and multi-view 3D human reconstruction and show the visualizations in Fig. 5. Based on these qualitative results, we observe that multi-view 3D human reconstruction consistently achieves more accurate 3D human poses and exhibits enhanced clothing details by leveraging additional information in the 3D space.

To further evaluate the impact of the number of input images, we present qualitative assessments in Fig. 6. By examining these visualizations, we can observe an enhancement in the reconstruction quality when the number of input views increased from three to five. Nevertheless, the enhancement might become minimal while the number of views continues to increase. This situation primarily arises due to two factors: (1) the input views are already saturated,

and (2) the performance might be correlated with the number of views used in model training.

## 4. NeRF-based 3D Human Novel View Synthesis

### 4.1. Neural Radiance Fields (NeRF)

Unlike 3D implicit human reconstruction methods that show high demands on the quantity and quality of the 3D hard-to-obtain dataset, NeRF [MST\*21] can achieve photo-realistic novel-view synthesis using only a limited set of images. Readers please refer to Fig. 7 for the overall framework of NeRF.

NeRF represents a 3D scene via an implicit function:

$$F_{\Theta}(\gamma(\mathbf{x})) \mapsto (\mathbf{c}, \sigma), \tag{22}$$

where $\gamma(\cdot)$ is a grid frequency encoder that lifts a 3D point $\mathbf{x}$ to a higher dimension, $\mathbf{c}$ is the RGB color, and $\sigma$ is the density value. Generally, the implicit function $F_{\Theta}(\cdot)$ is implemented as an MLP with trainable parameters $\Theta$. We omit $\Theta$ in the remaining sections for simplicity.

After getting the density and RGB color of a 3D point, NeRF employs a differentiable volume rendering module to render a 3D scene onto a 2D image. For each image pixel from a certain camera angle, the rendering involves casting a ray $\mathbf{r}$ from the pixel location into the 3D scene and sampling 3D points $\mu_i$ along the ray. The color $C$ of each image pixel is aggregated from the sampled points' color values $\mathbf{c}$:

$$C(\mathbf{r}) = \sum_i W_i \mathbf{c}_i, \quad W_i = \alpha_i \prod_{j<i} (1 - \alpha_j), \tag{23}$$

where $\alpha_i = 1 - e^{-\sigma_i \|\mu_i - \mu_{i+1}\|}$.

The loss for training NeRF is the total squared error between the rendered and true pixel colors:

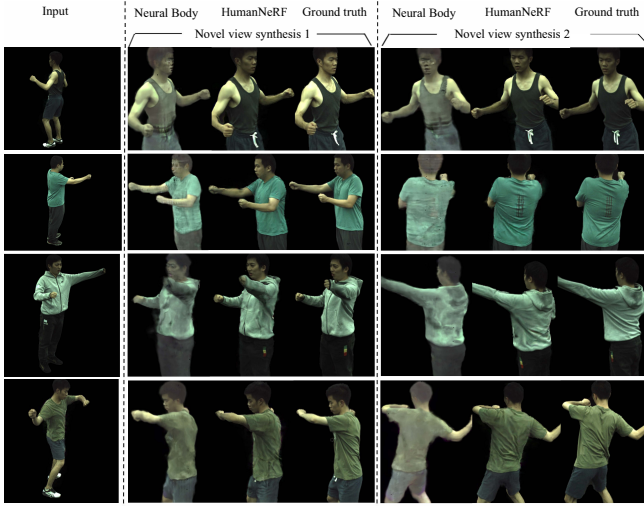$$\mathcal{L} = \sum_{\mathbf{r} \in \hat{R}} \|C_{\text{gt}}(\mathbf{r}) - C(\mathbf{r})\|_2^2, \tag{24}$$

**Figure 8:** *Input and output of NeRF-based reconstruction methods from static cameras. Figure obtained from [WCS\*22].*

where $\hat{R}$ is the set of rays in each batch.

## 4.2. NeRF-based Methods

Adopting the neural radiance field, many methods have been proposed for 3D human reconstruction and animation from unstructured photo sets under various scenarios.

### 4.2.1. 3D Human Reconstruction from Static Cameras

We first examine the scenario of 3D human reconstruction from static cameras, where all cameras are fixed in specific positions to capture images or videos of the human body.

Unlike PIFu which processes human subjects with only one pose during both training and inference, NeRF-based human reconstruction aims to optimize a 3D scene based on multiple images with various poses and viewpoints. To bridge the connections across different poses and angles, deformation fields based on SMPL(-X) [LMR\*15, PCG\*19] are usually applied to map a point $\mathbf{x}_o$ from an observation space to a corresponding point $\mathbf{x}_c$ in the canonical space. This process involves two key parts: (1) articulated deformation that applies the inverse transformation of SMPL linear blend skinning function $\mathtt{LBS}(\cdot)$ as in Equation (9) and Equation (11), and (2) non-rigid motion implemented as an MLP to learn the corrective offset:

$$\mathbf{x}_c = \mathbf{x}_c^{lbs} + \mathtt{MLP}_{\theta_{\mathrm{NR}}}\left(\gamma\left(\mathbf{x}_c^{lbs}\right)\right), \mathbf{x}_c^{lbs} = \mathcal{G}^{-1} \cdot \mathbf{x}_o, \qquad (25)$$

where $\gamma(\cdot)$ is the grid frequency encoder proposed in NeRF, and $\mathcal{G}$ annotated in Equation (10) is obtained from the observed SMPL vertex closet to $\mathbf{x}_o$.

To provide readers with a general understanding of the setups employed in NeRF-based reconstruction methods from static cameras, we present examples of their input and output in Fig. 8.

**HumanNeRF** HumanNeRF [WCS\*22] is one of the first to apply the deformation field to capture dynamic human models from monocular images under various viewpoints and motions. Specifically, it maps each point $\mathbf{x}_o$ sampled during volume rendering to its canonical counterpart $\mathbf{x}_c$ via Equation (25) before calculating the density and color value:

$$F(\gamma(\mathbf{x}_c)) \mapsto (\mathbf{c}, \sigma_c). \qquad (26)$$

Considering that the SMPL poses are separately estimated for each image (which may not be consistent), it further refines the body pose during training.

**Neural Body** Neural-Body [PZX\*21] advances HumanNeRF by introducing structured latent codes anchored to the SMPL model to provide regularization. These latent codes, $\mathcal{Z} = \{z_1, z_2, \ldots, z_{6890}\}$, correspond to SMPL's 6890 vertices. To address their sparsity in 3D space, Neural-Body uses SparseConvNet [GEVDM18] to process the latent codes and obtain 3D voxel-based feature volume. For any point $\mathbf{x}$, it is first transformed into the SMPL coordinate system, which aligns the point with the latent code volume. The latent code $\psi(\mathbf{x}, \mathcal{Z}, S_t)$ is then computed by tri-linear interpolation, and then input to two MLPs $F_\sigma$ and $F_\mathbf{c}$ to estimate the density and color:

$$F_\sigma(\psi(\mathbf{x}, \mathcal{Z}, S_t)) \mapsto \sigma, \qquad (27a)$$

$$F_\mathbf{c}(\psi(\mathbf{x}, \mathcal{Z}, S_t), \gamma_\mathbf{d}(\mathbf{d}), \gamma_\mathbf{x}(\mathbf{x}), \ell_t) \mapsto \mathbf{c}, \qquad (27b)$$

where $S_t$ is SMPL parameters at frame $t$, $\gamma_\mathbf{d}$ and $\gamma_\mathbf{x}$ are positional encoding functions for viewing direction $\mathbf{d}$ and spatial location $\mathbf{x}$, and $l_t$ is the latent embedding.

**Neural Human Performer** Unlike previous methods, Neural Human Performer [KKCF21] proposes a novel strategy for capturing information directly in the observation space. To achieve this, it first builds a skeletal feature bank by mapping the vertices of each SMPL model to its corresponding image and indexing the pixel-aligned image features. A temporal transformer then fuses these 2D features from different times (poses) and constructs a time-augmented skeletal feature bank. During the training and inference stages, Neural Human Performer (1) indexes skeletal features for a specific 3D point $\mathbf{x}$ from the feature bank using tri-linear interpolation, and (2) projects $\mathbf{x}$ to each image and acquires pixel-aligned image features via bi-linear interpolation. These features are then fused using a multi-view transformer to predict the density and color.

Besides the above methods: considering that multiple points in 3D space can be projected onto the same surface point on the mesh, Xu *et al.* [XFM22] utilizes barycentric interpolation with vertex normals to project points onto the mesh surface; I-M-Avatar [ZAB\*22] presents a morphing-based implicit model tailored for head avatars and utilizes expression and pose deformations for detailed geometry and appearance; Semantic-Facial-NeRF [GZX\*22] combines multi-level voxel fields with expression coefficients in the latent space to represent head avatars with complex facial attributes; Vid2Avatar [GJC\*23] proposes to jointly model the human and scene background via two neural radiance fields for a clean separation of the dynamic human and static background.

### 4.2.2. 3D Human Reconstruction from Dynamic Cameras

Besides static camera scenarios, researchers also propose to reconstruct humans under dynamic camera settings where the cameras
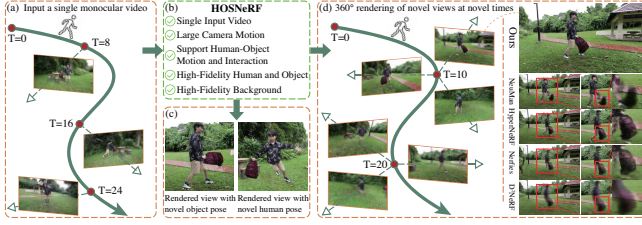
**Figure 9:** *Input and output of HOSNeRF. Figure obtained from [LCY\*23].*



**Figure 10:** *Input and output of NeRF-based animation methods. Figure obtained from [PDW\*21].*

are in motion during the capture process, which are more suitable for real-world applications. See Fig. 9 for an example of the input and output under such settings with dynamic cameras.

**HOSNeRF** Considering that the SMPL model represents only the rough human body, disregarding the clothing topology and accessories, HOSNeRF [LCY\*23] proposes object bones and state-conditional representations with learnable state embeddings for synthesizing 3D humans with items like bags. Specifically, object bones extend the human skeleton to better model deformations from human-object interactions. After defining the object bones, HOSNeRF also applies backward LBS and a non-rigid deformation module to map a 3D point from observation space to canonical space. In frame $i$, the dynamic 3D scene can be represented as:

$$F(\gamma(\mathbf{x}_c), \mathcal{O}_c^i) \mapsto (\mathbf{c}, \sigma), \tag{28}$$

where $\mathcal{O}_c^i$ is the learnable state embedding representing object states in the canonical space at frame $i$. As humans interact with objects at different times, state embeddings serve as conditions for learning representations of human-object interactions and the scene. Finally, HOSNeRF applies Mip-NeRF 360 [BMV\*22] to represent the background scene.

Similar to HOSNeRF, NeuMan [JYS\*22] trains separate NeRF models for humans and scenes to reconstruct both from the input video. To further handle the challenge of maintaining identities through occlusion events, 4DHumans [GPR\*23] proposes an HMR 2.0 network, preserving more details in multiple-person scenarios. PPR [YYZ\*23] combines differentiable physics simulation and differentiable rendering via coordinate descent, largely reducing reconstruction artifacts. RAC [YWRR23] introduces a reconstruction technique for both animals and humans by learning skeletons with constant bone lengths within a video.

#### 4.2.3. 3D Human Animation

3D human animation focuses on learning how the clothing topology changes with different poses based on multi-view and multi-pose images. This enables the synthesis of free-viewpoint animations under user-guided pose sequences that are out of training distribution from multi-view videos. An example of the input and output of these animation methods is illustrated in Fig. 10.

**Neural Actor** To achieve this, Neural Actor [LHR\*21] first adopts a deformation field that contains the articulate deformation and non-rigid motion (called residual deformation in Neural Actor), enabling the transfer of a 3D point $\mathbf{x}_o$ from the observation
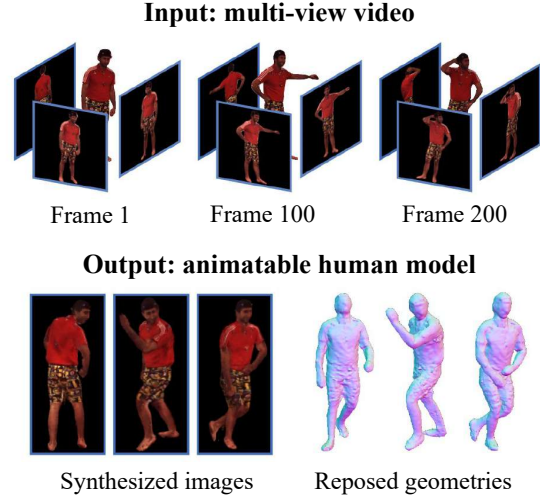
space to the canonical space $\mathbf{x}_c$. Additionally, it employs a feature extractor to acquire texture features from a 2D texture map, then indices the texture features of $\mathbf{x}_o$ based on its closest surface point $\mathbf{x}_s$:

$$F(\gamma_{\mathbf{x}_c}(\mathbf{x}_c), \gamma_{\mathbf{d}}(\mathbf{d}), \mathcal{Z}(\mathbf{x}_s)) \mapsto (\mathbf{c}, \sigma), \tag{29}$$

where $\mathcal{Z}(\cdot)$ denotes the extracted texture features. Specifically, during training, it obtains $\mathcal{Z}(\cdot)$ from the input multi-view images, while at the inference stage, it relies on a SMPL-based normal map to obtain $\mathcal{Z}(\cdot)$.

**Animatable NeRF** Considering that articulated deformation may not efficiently capture clothing topology, Peng *et al.* [PDW\*21] introduce a novel neural blend weight field to obtain better skinning weights $w_k$. To this end, they first define a per-frame latent code $l_i$ that encodes the human appearance in frame $i$. Given a 3D point $\mathbf{x}$, they transform it to the canonical space by:

$$w_i'(\mathbf{x}) = \text{norm}(F_{\delta_w}(\mathbf{x}, l_i) + w_i(\mathbf{x})), \tag{30a}$$

$$\mathbf{x}_c = \left( \sum_{k=1}^{K} w_k'(\mathbf{x}) \mathcal{G}_k \right)^{-1} \cdot \mathbf{x}, \tag{30b}$$

where $F_{\delta_w}$ is the neural blend weight field and the definition of other terms are the same as Equation (10). For each frame $i$, the color and density value are then given by:

$$F_{\sigma}(\gamma_{\mathbf{x}}(\mathbf{x}_c) \mapsto (\sigma_i(\mathbf{x}_c), \mathbf{z}_i(\mathbf{x}_c)), \tag{31a}$$

$$F_{\mathbf{c}}(\mathbf{z}_i(\mathbf{x}_c), \gamma_{\mathbf{d}}(\mathbf{d}), l_i) \mapsto \mathbf{c}_i(\mathbf{x}), \tag{31b}$$

where $\mathbf{z}_i(\cdot)$ is the intermediate shape feature.

**TAVA** Unlike previous methods that consider only one possible canonical point $\mathbf{x}_c$ after deformation, TAVA [LTV\*22] proposes to find the canonical candidates $\{\mathbf{x}_{c,1}, \mathbf{x}_{c,2}, ..., \mathbf{x}_{c,K}\}$ and predict the

color and density values for all the candidates:

$$F_{\mathbf{c},\sigma}(\mathbf{x}_{c,i}, \Sigma) \mapsto \mathbf{h} \mapsto (\mathbf{c}_i^*, \sigma_i), \tag{32a}$$

$$F_a(\mathbf{h}, \mathbf{d}) \mapsto a_i, \tag{32b}$$

where $\Sigma$ is the multivariate Gaussians applied in Mip-NeRF [BMV*22], $\mathbf{h}$ is the intermediate output, and $a_i$ is the ambient occlusion at point $\mathbf{x}_c$ under viewing direction $\mathbf{d}$. TAVA then chooses final values based on their density:

$$\mathbf{c} = \mathbf{c}_t^* \cdot a_t, \quad \sigma = \sigma_t, \quad \text{where} \quad t = \arg\max_i(\sigma_i). \tag{33}$$

**NARF** Evolving from TAVA, NARF [NSLH21] divides the articulated human into several rigid parts according to the skeleton of the SMPL model, and each rigid part represents a local coordinate system. For a 3D point $\mathbf{x}$, NARF first transforms it and the viewing direction $\mathbf{d}$ to each local system:

$$\mathbf{x}^i = \mathbf{R}^{i^{-1}}(\mathbf{x} - \mathbf{t}^i), \quad \mathbf{d}^i = \mathbf{R}^{i^{-1}}\mathbf{d}, \tag{34}$$

where $\mathbf{R}^i$ and $\mathbf{t}^i$ are the corresponding rotation matrix and translation vector for part $i$. To predict the exact part that point $\mathbf{x}$ belongs to, NARF includes a selector network $\mathcal{S}$ that consists of $P$ lightweight sub-networks for each rigid part:

$$F_{\mathcal{S}}(\gamma(\mathbf{x}^i, \gamma(\zeta)) \mapsto (s^i), \quad p^i = \frac{\exp(s^i)}{\sum_{k=1}^P \exp(s^k)}, \tag{35}$$

where $\zeta$ is the bone parameter, and $s^i$ is the occupancy value. The density and color can then be predicted via:

$$F_{\sigma}(\{\gamma_{\mathbf{x}}(\mathbf{x}^i) * p^i | i \in [1,P]\}, \gamma(\zeta)) \mapsto (\sigma, \mathbf{h}), \tag{36a}$$

$$F_{\mathbf{c}}(\mathbf{h}, \{(\gamma_{\mathbf{d}}(\mathbf{d}^i) * p^i, \gamma(\xi^i) * p^i | i \in [1,P]\}) \mapsto \mathbf{c}, \tag{36b}$$

where $\xi^i$ is the 6D transformation vector of part $i$.

**ARAH** Instead of applying the articulated deformation and non-rigid motion, ARAH [WSGT22] introduces a joint root-finding algorithm designed to find a canonical point and its depth along the viewing direction, which satisfies both the SDF iso-surface condition[‡] and the LBS condition[§]. Subsequently, the SDF value and color networks employed in ARAH can be formulated as:

$$F_{\text{SDF}}(\mathbf{x}_c, \theta, \beta, l) \mapsto \mathbf{h} \mapsto d, \tag{37a}$$

$$F_{\mathbf{c}}(\mathbf{x}_c, \mathbf{n}_o, \mathbf{d}, \mathbf{h}, l) \mapsto \mathbf{c}, \tag{37b}$$

where $\theta, \beta$ are the pose and shape parameters, $l$ denotes the latent code as in Neural Body, $\mathbf{n}_o$ is the normal vector in the observation space, and $\mathbf{h}$ is the intermediate output for calculating the SDF value $d$.

**PersonNeRF** Unlike previous methods that learn from videos and drive the human body following the given poses, Person-NeRF [WSCKS23] builds a neural radiance field spanned by camera view, body pose, and appearance based on an image collection of a specific person with different poses and clothes (see Fig. 11).
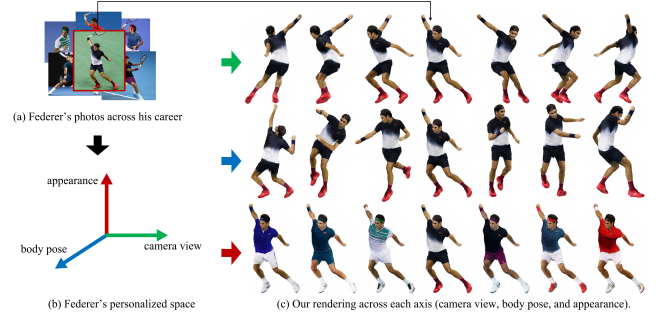
---

[‡] The canonical point should be on the iso-surface.

[§] After applying forward LBS to the canonical point, the transformed point should lie on the specified camera ray.



(a) Federer's photos across his career

appearance

body pose          camera view

(b) Federer's personalized space        (c) Our rendering across each axis (camera view, body pose, and appearance).

**Figure 11:** *Input and output of PersonNeRF. Figure obtained from [WSCKS23].*

To achieve this and enable traversing the space to explore unobserved combinations of human attributes, PersonNeRF incorporates (1) a pose embedding to correlate the estimated pose, and (2) an appearance embedding pre-trained on the image collection to ensure appearance consistency. It further proposes to learn a shared motion field that outputs joint angle residuals and is learned across different clothings. It constrains all body poses in the dataset regardless of their appearance, enhancing pose generalization.

Moreover, several other works also leverage deformation fields for free-viewpoint animations. MPS-NeRF [GYK*22] introduces two deformation fields to connect the canonical space with the observation space and target space, one for extracting image features from the input image for radiance prediction and another for rendering the output image. SHERF [HHP*23] utilizes 3D-aware global, point-level, and pixel-aligned features for effective encoding and a feature fusion transformer to predict the color and density. MonoHuman [YCL*23] proposes a Shared Bidirectional Deformation module to achieve generalizable consistent forward and backward deformation. ActorsNeRF [MSVW23] designs a 2-level canonical space (a category-level canonical space and an instance-level canonical space) for a coarse-to-fine strategy.

### 4.2.4. Methods Combining Mesh and Radiance Field

Following PIFu [SHN*19] and NeRF [MST*21], many methods have also been proposed to combine the merits of both mesh and radiance fields for high-fidelity 3D human reconstruction.

**NeuralHumanFVV** Among these methods, NeuralHuman-FVV [SJL*21] consists of a neural geometry reconstruction stage and a neural blending stage to produce live 4D renderings based on 6-view dynamic videos. In the neural geometry reconstruction stage, it (1) extracts a coarse geometry prior via Shape-from-Silhouette (SfS) [CBK03] algorithm, (2) obtains a finer geometry using a multi-view implicit function based on PIFu, and (3) utilizes a hierarchical sampling strategy to recover geometry details such as clothing folds. Specifically, in the hierarchical sampling strategy, a depth fine-tuning network takes the feature of the midpoint between two selected sample points $\mathbf{x}_1$ and $\mathbf{x}_2$ as input, and outputs the displacement of depth value, which is then used to refine the depth value. In the neural blending stage, it encodes the fine-detailed geometry and texture information from adjacent input views into a photo-realistic texture output.

1 of 5 views  PixelNeRF  PIFu  DoubleField   1 of 5 views  PixelNeRF  PIFu  DoubleField

**Figure 12:** *Comparison of methods combining mesh and radiance field with those using either alone (PIFu [SHN*19]: mesh only; PixelNeRF [YYTK21]: radiance field only). Figure obtained from [SZZ*22a].*

**Function4D** Function4D [YZG*21] achieves real-time reconstruction by integrating the proposed Dynamic Sliding Fusion (DSF) and deep implicit surface reconstruction. Unlike traditional volumetric fusion methods that attempt to complete surfaces by fusing all available temporal depth observations, DSF focuses on augmenting current observations to maintain consistency and reduce noise without relying on long-term tracking. This is done by confining tracking and fusion processes to a sliding window of the current, previous, and next frames. Having the depth map $\mathcal{D}$ extracted from DSF, the deep implicit surface reconstruction then includes (1) a GeoNet that uses truncated projective SDF (PSDF) values as a novel feature for preserving geometric details[¶]:

$$\text{PSDF}(\mathbf{x}) = \mathbf{T}\left(d - \mathcal{B}(\pi(\mathbf{x}), \mathcal{D})\right), \quad (38)$$

and (2) a ColorNet that utilizes a multi-head transformer network to aggregate features across multiple views.

**DoubleField** Different from the above methods, DoubleField [SZZ*22a] combines the surface and radiance field at the feature level in an implicit manner via a novel Network $F_{db}$. Given the query point $\mathbf{x}$, viewing direction $\mathbf{d}$ and 2D image features obtained from the input image $I$, $F_{db}$ learns a shared double embedding and predicts the occupancy $s$, the density $\sigma$ and the color $\mathbf{c}$ simultaneously. DoubleField network consists of a shared MLP (the Double MLP $F_{db}$) for learning the double embedding $e_{db}$ and two individual MLPs (the geometry MLP $F_g$ and the texture MLP $F_c$) for the surface and radiance fields prediction:

$$F_{db}(\gamma(\mathbf{x}), \mathcal{B}(f_{2D}, \pi(\mathbf{x}))) \mapsto e_{db}, \quad (39a)$$

$$F_g(e_{db}) \mapsto (s, \sigma), \quad F_c(e_{db}, d) \mapsto \mathbf{c}. \quad (39b)$$

Given that $s$ and $\sigma$ are both outputs from the same MLP layer, this approach inherently creates a robust link between the two fields, enabling their cooperation at the feature level.

Beyond the above-discussed methods, SelfRecon [JHBZ22] combines explicit and implicit geometry representations to obtain coherent geometry. DNA-Net [VPY*23] includes a Neural Articulations Prediction Network (NAP-Net) to project observed points into the canonical space for improved learning of the geometry (SDF) and color fields.

In order to showcase the advancements achieved by methods that integrate both mesh and neural radiance field, we compare Double-Field with PIFu and PixelNeRF [YYTK21] in Fig. 12. This comparison highlights the superior performance of DoubleFiled, as it effectively combines the strengths of both representations to enhance the modeling of geometry and texture.

## 5. 3D Gaussian-based Methods

### 5.1. 3D Gaussian Splatting

In contrast to NeRF [MST*21], which employs neural networks to synthesize novel views, 3D Gaussian Splatting (3DGS) [KKLD23] introduces a novel approach that directly optimizes the position and attributes of 3D Gaussians, i.e, 3D position, opacity $\alpha$, anisotropic covariance, and spherical harmonic (SH) [RH01] coefficients[∥]. This methodology enables us to efficiently represent and render intricate scenes at high resolutions, significantly reducing training time.

*(1) Definition.* 3D Gaussian Splatting takes a set of static scene images and corresponding camera parameters obtained from Structure-from-Motion (SfM) [SSS06] as inputs. For each SfM point $\mathbf{x}$, a 3D Gaussian is defined by a 3D covariance matrix $\Sigma$ centered at point (mean) $\mu$:

$$G(\mathbf{x}) = e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}. \quad (40)$$

*(2) Rendering.* To render the scene consisting of 3D Gaussians onto 2D image space, 3DGS incorporates the splatting rasterization.

a) Inspired by [LZ21], 3DGS designs a tile-based rasterizer. The screen is first divided into tiles (e.g., $16 \times 16$ pixels), and each Gaussian is instantiated according to the number of tiles they overlap and assigned a key that records view space depth and tile ID. Gaussians are then sorted based on their depth, allowing the rasterizer to correctly handle occlusions and overlapping geometry.

b) 3DGS introduces a point-based $\alpha$-blend rendering to obtain the RGB color $\mathbf{C}$. Specifically, it samples points along the ray with intervals $\delta_i$:

$$\mathbf{C}_{\text{color}} = \sum_{i \in N} \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1}(1 - \alpha_j), \quad (41)$$

with

$$\alpha_i = (1 - \exp(-\sigma_i \delta_i)), \quad (42)$$

where $\sigma_i$ and $\mathbf{c}_i$ are the density and color of each point along the ray.

---

[¶] $\mathbf{T}(\cdot)$ truncates the PSDF values in $[-\delta, \delta]$, with $\delta$ being a small positive threshold.

[∥] Spherical Harmonic (SH) is used for controlling the color of each Gaussian to accurately capture the view-dependent appearance of the scene.
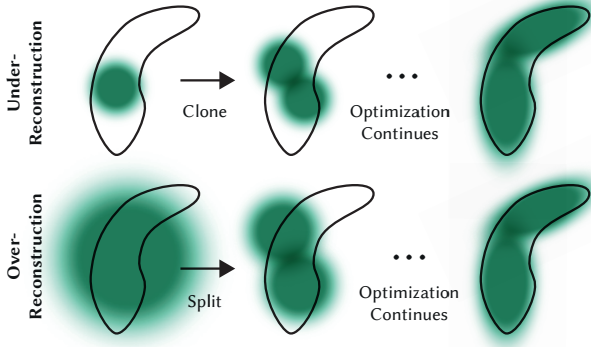
**Figure 13:** *The densification and culling scheme. Figure obtained from [KKLD23].*

*(3) Optimization.* Following [FKYT*22, SSC22], Stochastic Gradient Descent techniques are utilized for optimization. Specifically, 3DGS calculates the initial covariance matrix as an isotropic Gaussian with axes equal to the average distance to the closest three points. They use a standard exponential decay scheduling method similar to Plenoxels [FKYT*22] but for positions only. The optimization procedure is supervised by calculating the $\mathcal{L}_1$ loss and the D-SSIM term between the ground truth $g$ and the rendering $r$:

$$\mathcal{L} = (1-\lambda)\mathcal{L}_1 + \lambda\mathcal{L}_{\text{D}-\text{SSIM}}, \tag{43}$$

where

$$\mathcal{L}_1 = \frac{1}{N}\sum_{i=1}^{N}|g_i - r_i|, \tag{44a}$$

$$\mathcal{L}_{\text{D}-\text{SSIM}} = 1 - \frac{1}{M}\sum_{j=1}^{M}\frac{(2\mu_g\mu_r + c_1)(2\sigma_{gr} + c_2)}{(\mu_g^2 + \mu_r^2 + c_1)(\sigma_g^2 + \sigma_r^2 + c_2)}, \tag{44b}$$

where $\mu_g$ and $\mu_r$ are the average intensities of the ground truth and the rendered image within each of $M$ local windows[**], $\sigma_g^2$ and $\sigma_r^2$ are the variances of the ground truth and the rendered image, respectively, and $\sigma_{gr}$ is the covariance between them. Constants $c_1$ and $c_2$ help stabilize the division with small denominators.

*(4) Densification and Culling.* To ensure that enough details and accurate reconstruction of the scene can be optimized, 3DGS incorporates densification and culling during the optimization. Specifically, 3DGS densifies every 100 iterations and remove any Gaussians that are essentially transparent, i.e., with opacity $\alpha$ less than a threshold $\epsilon_\alpha$. Readers can see Fig. 13 for a detailed illustration of the densification scheme.

In order to provide readers with a clearer understanding of the differences between NeRF-based and Gaussian-based approaches, we have included a series of visualizations in Figure 14. By examining these visualizations, several key observations can be made: (1) NeRF-based methods employ MLPs to predict the density and color values while Gaussian-based methods directly optimize the

---

[**] Local windows refer to segmented areas of an image which are used to analyze statistical properties locally.
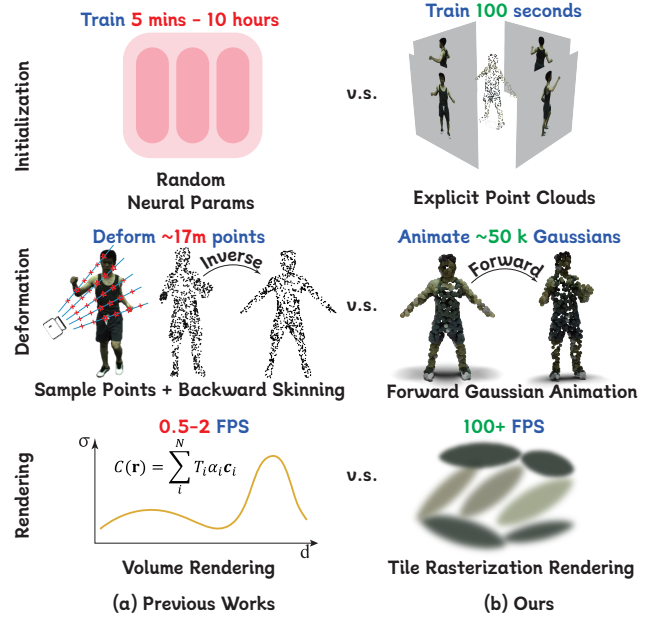


**Figure 14:** *Differences between Gaussian-based methods and NeRF-based methods. Figure obtained from [LTYY23].*

attributes of 3D Gaussians. This direct optimization leads to significantly enhanced training efficiency. (2) Due to the disparity in their 3D representations, NeRF-based methods encounter non-correlated points, resulting in a larger number of points and significantly increased computational requirements. (3) Gaussian-based methods have the advantage of rendering the 3D scene directly through tile-based rasterization. This enables them to achieve greatly improved real-time performance.

### 5.2. 3D Human Reconstruction

Leveraging 3DGS [KKLD23], several works have been introduced for human-related tasks.

**Animatable 3D Gaussian** Animatable-3D-Gaussian [LHQ*23] is one of the first methods that apply 3DGS to enhance training efficiency and reduce GPU requirements. It incorporates: (1) A deformation network that transforms the point $\mathbf{x}_c$, rotation matrix $\mathbf{R}_c$, and viewing direction $\mathbf{d}_c$ from the canonical space to the observation space:

$$\mathbf{x}_o = \sum_{i=1}^{n_b} w_i B_i^t \mathbf{x}_c,$$

$$\mathbf{R}_o = \sum_{i=1}^{n_b} w_i B_i^t \mathbf{R}_c, \quad \mathbf{d}_o = \left(\sum_{i=1}^{n_b} w_i B_i^t\right)^{-1}\mathbf{d}_c, \tag{45}$$

where $n_b$ is the number of bones, $w_i$ is the skinning weight as in Equation (10), and $B_i^t$ represents the transformation of the $i$-th bone in frame $t$, and (2) a time-dependent ambient occlusion that addresses the issue of dynamic shadows. Specifically, it predicts the ambient occlusion factor $a_o \in [0, 1]$ based on position $\mathbf{x}_o$ and hash-

encoded time $t$:

$$a_o = \text{MLP}(\mathbf{x}_o, \gamma(t)), \quad \mathbf{c} = a_o * \mathbf{c}_o. \tag{46}$$

**Drivable 3D Gaussian Avatars** Drivable-3D-Gaussian-Avatars (D3GA) [ZBS*23] leverages tetrahedral cages and cage-based deformation fields to model the body and individual garments, learning both the 3D human-related scenes and human segmentation maps. Given the canonical cage $\mathbf{v}$ and pose $\theta$, D3GA first employs two separate MLPs to model the deformation field of tetrahedron $i$:

$$\Psi_{\text{MLP}} : \{\theta, \gamma_{\mathbf{v}}(\mathbf{v})\} \to \Delta\mathbf{v}, \tag{47a}$$

$$\Pi_{\text{MLP}} : \{\theta, \mathbf{b}_i, \mathbf{r}_i, \mathbf{s}_i\} \to \{\Delta\mathbf{b}_i, \Delta\mathbf{r}_i, \Delta\mathbf{s}_i\}, \tag{47b}$$

where the cage node correction network $\Psi_{\text{MLP}}$ takes position-encoded canonical vertices $\gamma_{\mathbf{v}}(\mathbf{v})$ to predict offsets $\Delta\mathbf{v}$ for the cage node positions. The Gaussian correction network $\Pi_{\text{MLP}}$ uses the canonical Gaussian parameters (barycentric coordinates $\mathbf{b}_i \in \mathbb{R}^4$, rotation $\mathbf{r}_i \in \mathbb{R}^4$ and scale $\mathbf{s}_i \in \mathbb{R}^3$) to predict their corrections. A shading network $\Gamma_{\text{MLP}}$ is subsequently applied to learn the color $\mathbf{c}$ and opacity $o_i$:

$$\Gamma_{\text{MLP}} : \{\theta, \gamma_{\mathbf{d}}(\mathbf{d}_k), \mathbf{h}_i, \mathbf{f}_j\} \to \{\mathbf{c}_i, o_i\}, \tag{48}$$

where $\mathbf{h}_i$ is the auto-decoded [PFS*19] feature vector of the initial color, $\mathbf{f}_j$ is the embedding vector with the time frame of the current sample.

**Human101** Human101 [LTYY23] reconstructs and animates dynamic human avatars from single-view videos with real-time rendering. Given an input video, it first extracts multiple sets of images with various poses. Each set contains four images of four orientations (front, back, left, right), which are used to create corresponding meshes using ECON [XYC*23]. The meshes are then deformed to the canonical space using inverse LBS from Equation (9) and fused into a canonical point cloud. Following the original 3DGS, the point cloud is converted into canonical Gaussians for initialization. During animation, the Gaussians are deformed into the target pose by modifying their positions, rotations, and scales, and adjusting spherical harmonic coefficients. To address potential inconsistencies in human movement, Human101 employs an MLP to predict the residuals of position $\mathbf{x}$, rotation $\mathbf{r}$, and scale $\mathbf{s}$:

$$F(\gamma_{\mathbf{x}}(\mathbf{x}), \gamma_t(t), \theta, \beta) \mapsto (\Delta\mathbf{x}, \Delta\mathbf{r}, \Delta\mathbf{s}). \tag{49}$$

**Gaussian Head Avatar** Gaussian-Head-Avatar [XCL*23] proposes to reconstruct 3D heads based on 3DMM [BV99], 3D neutral landmarks, and triplane features. Specifically, it first construct a canonical neural Gaussian model with expression-independent attributes:

$$\{\mathbf{X}_o, \mathbf{F}_o, \mathbf{Q}_o, \mathbf{S}_o, \mathbf{A}_o\}, \tag{50}$$

where $\mathbf{X}_o \in \mathbb{R}^{N \times 3}$ is the position of the Gaussians in the canonical space, $\mathbf{F}_o \in \mathbb{R}^{N \times 128}$ denotes the point-wise feature vectors, $\mathbf{Q}_o \in \mathbb{R}^{N \times 4}, \mathbf{S}_o \in \mathbb{R}^{N \times 3}, \mathbf{A}_o \in \mathbb{R}^{N \times 1}$ represent the rotation, scale, and opacity respectively. To further improve the quality of geometry and texture, it proposes to condition the Gaussian attributes on the

expression coefficients $\varepsilon_{\text{3DMM}}$ and head pose $\theta_{\text{3DMM}}$:

$$\mathbf{X} = \mathbf{X}_o + \lambda_{\exp}(\mathbf{X}_o) \cdot \text{MLP}_{\mathbf{X}}^{\exp}(\mathbf{X}_o, \varepsilon_{\text{3DMM}})$$
$$+ \lambda_{\text{pose}}(\mathbf{X}_o) \cdot \text{MLP}_{\mathbf{X}}^{\text{pose}}(\mathbf{X}_o, \theta_{\text{3DMM}}), \tag{51a}$$

$$\mathbf{C} = \lambda_{\exp}(\mathbf{X}_o) \cdot \text{MLP}_{\mathbf{C}}^{\exp}(\mathbf{F}_o, \varepsilon_{\text{3DMM}})$$
$$+ \lambda_{\text{pose}}(\mathbf{X}_o) \cdot \text{MLP}_{\mathbf{C}}^{\text{pose}}(\mathbf{F}_o, \theta_{\text{3DMM}}), \tag{51b}$$

$$\{\mathbf{Q}, \mathbf{S}, \mathbf{A}\} = \{\mathbf{Q}_o, \mathbf{S}_o, \mathbf{A}_o\} + \lambda_{\exp}(\mathbf{X}_o) \cdot \text{MLP}_{att}^{\exp}(\mathbf{F}_o, \varepsilon_{\text{3DMM}})$$
$$+ \lambda_{\text{pose}}(\mathbf{X}_o) \cdot \text{MLP}_{att}^{\text{pose}}(\mathbf{F}_o, \theta_{\text{3DMM}}). \tag{51c}$$

By applying rigid rotations and translations to Gaussians, D3GA achieves dynamic 3D human head modeling in the observation space.

**SC-GS** Given an image sequence from a monocular dynamic video, SC-GS [HSY*23] utilize a set of sparse control points $\mathcal{P} = \{(p_i \in \mathbb{R}^3, \iota_i \in \mathbb{R}^+)\}, \quad i \in \{1, 2, ..., N_p\}^{\dagger\dagger}$ to reconstruct and drive 3D Gaussians for high-fidelity rendering. For each control point $p_i$, SC-GS first learns time-varying 6 DoF transformations $[\mathbf{R}_i^t | \mathbf{T}_i^t] \in \mathbf{SE}(3)$, which consists of a local frame rotation matrix $\mathbf{R}_i^t \in \mathbf{SO}(3)$ and a translation vector $\mathbf{T}_i^t \in \mathbb{R}^3$:

$$F(p_i, t) \mapsto (\mathbf{R}_i^t, \mathbf{T}_i^t). \tag{52}$$

Given the learned 6 DoF transformations for the control points and for each 3D Gaussian, SC-GS applies k-nearest neighbor (KNN) [CH67] to obtain its $K(K = 4)$ neighboring control points in the canonical space and calculate the interpolation weight:

$$w_{jk} = \frac{\hat{w}_{jk}}{\sum_{k \in K} \hat{w}_{jk}}, \quad \text{where} \quad \hat{w}_{jk} = \exp(-\frac{d_{jk}^2}{2o_k^2}), \tag{53}$$

where $d_{jk}$ is the distance between the center of Gaussian $G_j$ and the neighboring control point $p_k$. By applying LBS [SSP07], we can obtain the transformed Gaussian location $\mu_j^t$ and rotation matrix $R_j^t$ by:

$$\mu_j^t = \sum_{k \in K} (\mathbf{R}_k^t(\mu_j - p_k) + p_k + \mathbf{T}_k^t), \tag{54a}$$

$$\mathbf{R}_j^t = \left(\sum_{k \in K} w_{jk} r_k^t\right) \otimes \mathbf{R}_j, \tag{54b}$$

where $r_k^t$ is the quaternion of control point $k$ and $\otimes$ is the production process.

**GauHuman** After initializing the position $\mathbf{x}_c$ of 3D Gaussians from SMPL vertex points, GauHuman [HL23] incorporates an LBS weight field module and a pose refinement module to transform Gaussians from the canonical space to the posed space. In the LBS weight field module, for each 3D Gaussian, GauHuman adds the LBS weight $w_k$ of nearest SMPL vertex with the predicted offsets $\text{MLP}_{\Phi_{\text{lbs}}}(\gamma(\mathbf{x}_c))$ to predict the LBS weight coefficients $w_k'$:

$$w_k' = \frac{e^{\log(w_k + 10^{-8}) + \text{MLP}_{\text{lbs}}(\gamma(\mathbf{x}_c))[k]}}{\sum_{k=1}^{K} e^{\log(w_k + 10^{-8}) + \text{MLP}_{\text{lbs}}(\gamma(\mathbf{x}_c))[k]}}. \tag{55}$$

---

$\dagger\dagger$ $p_i$ denotes the learnable coordinate of control point in the canonical space, and $\iota_i$ is the learnable radius parameter of a radial-basis-function (RBF) kernel that controls the impact of a control point on a Gaussian.

In the pose refinement module, it updates the joint angles $\theta^{\mathrm{SMPL}}$ via another MLP:

$$\theta = \theta^{\mathrm{SMPL}} \otimes \mathtt{MLP}_{\mathrm{pose}}(\theta^{\mathrm{SMPL}}), \qquad (56)$$

where $\theta^{\mathrm{SMPL}}$ is the SMPL body pose parameter estimated from images. During optimization, GauHuman adopts the tile-based differentiable rasterizer from 3DGS for rapid rendering. They also introduce a strategy to adaptively control the number of Gaussians by employing human priors such as SMPL and Kullback-Leibler (KL) divergence to guide the split, clone, merge, and prune process.

Besides the methods mentioned above, Gaussian-Flow [LDZY23] proposes a novel Dual-Domain Deformation Model for 4D scene training and rendering, avoiding per-frame 3DGS optimization. 3DGS-Avatar [QWM*23] develops a non-rigid deformation network to reconstruct human avatars, further enhancing the reconstruction speed. To capture finer details, GaussianBody [LYX*24] employs explicit pose-guided deformation to reduce ambiguity between the observation and the canonical space. GVA [LWL*24] further introduces (1) a pose refinement method to enhance the alignment between the body and hand alignment, and (2) a surface-guided Gaussian re-initialization technique to address issues of unbalanced aggregation and initialization bias. There are also increasing efforts [GL23, SRV23, DWY*23, CLTS23] that utilize 3D Gaussian Splatting for object reconstruction. These papers can also provide valuable insights and inspiration for 3D human modeling.

## 6. GAN-based 3D Human Generation

Moving to 3D human generation, we first discuss methods utilizing Generative Adversarial Network (GAN) [GPAM*14]. Although StyleGAN [FLJ*22] and its follow-up works are conducted in the 2D image space, they provide the dataset and network basis for the recent 3D-GAN network. Therefore, this section mainly discusses StyleGAN-related techniques and 3D-GAN approaches that employ triplane [CLC*22] as 3D representation.

### 6.1. StyleGAN-related Techniques

StyleGAN [FLJ*22] extends GAN [GPAM*14] by introducing a "style" space to control human attributes. Its architecture contains two components: the mapping network and the generator. Given the latent code $\mathbf{z}$, the mapping network is applied to predict a control signal $\mathbf{w}$. The generator takes $\mathbf{w}$ as the input to generate the final image via a sequential model comprised of several layers that progressively increase the resolution of the image. To enhance the influence of the control signal, $\mathbf{w}$ is used in each convolution layer of the generator after an adaptive instance normalization (AdaIN).

After StyleGAN, StyleGAN2 [KLA*20] and Style-GAN3 [KAL*21] are proposed to further enhance the image quality. StyleGAN2 introduces improvements in data augmentation and architectural changes. Specifically, it discards the progressive increase of layers for a more stable training process and replaces AdaIN with weight demodulation to reduce artifacts. StyleGAN3 is specifically designed to preserve image quality and distinctive features regardless of transformations. Even when an
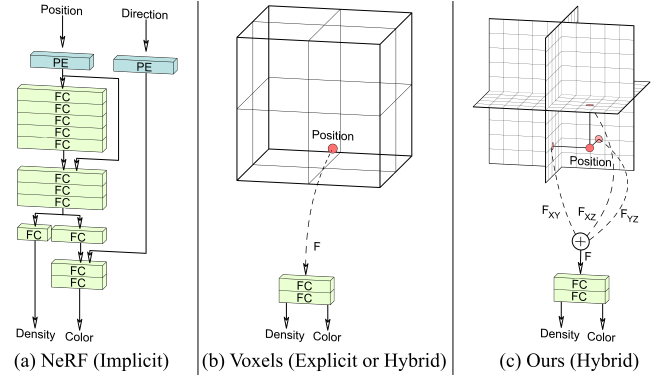


(a) NeRF (Implicit)  (b) Voxels (Explicit or Hybrid)  (c) Ours (Hybrid)

**Figure 15:** *Triplane representation architecture. Figure obtained from [CLC*22].*

image is subject to rotation or zoom, the resulting modifications are still coherent and predictable.

StyleGAN also collects and annotates the Stylish-Humans-HQ Dataset (SHHQ), which contains over 230000 high-quality images with various poses and textures. They introduce a "model zoo" that is trained on the SHHQ dataset using the StyleGAN2 framework and consists of six human-GAN models that are capable of generating full-body images with diverse poses and clothing textures.

Leveraging the SHHQ dataset and these three models, many methods have been proposed to improve the performance from various perspectives. One such method is StyleSDF [OELS*22], which elevates the model to the 3D space by employing SDF-based volume rendering and a 3D implicit network. Holo-GAN [NPLT*19] learns 3D features from a 4D constant tensor and separates pose, shape, and appearance for finer details. Block-GAN [NPRM*20] learns 3D scene representations directly from unlabelled 2D images, providing control over each object's 3D pose and identity. StylePeople [GII*21] goes a step further by integrating the deformable SMPL-X model and neural rendering techniques into the StyleGAN2 network.

### 6.2. 3D-GAN Approaches

As NeRF [MST*21] has advanced the field, EG3D [CLC*22] proposes a triplane representation and integrates it into the GAN framework. Remarkably, this approach represents one of the first instances where 3D model generation is accomplished using only 2D image data for training. See Fig. 15 for the visualization of triplane and its difference with NeRF and voxel representation.

Specifically, the process begins by projecting each query point $\mathbf{x} \in \mathbb{R}^3$ onto three feature planes, enabling the indexing of the corresponding feature vectors $(f_{xy}, f_{xz}, f_{yz})$ through bilinear interpolation. These features are then summed, and input to an MLP to predict the color $\mathbf{c}$ and density $\sigma$:

$$F(f_{xy}(\mathbf{x}) + f_{xz}(\mathbf{x}) + f_{yz}(\mathbf{x})) \mapsto (\sigma, \mathbf{c}). \qquad (57)$$

By incorporating volume rendering techniques, EG3D establishes a connection between the triplane representation and the Style-GAN2 [KLA*20] framework.

## 6.3. 3D Human Generation

Following EG3D, many methods have been proposed to handle 3D human generation.

**GNARF** To better present humans with dynamic motions, GNARF [BKY*22] proposes a Surface Field (SF) method to transform the query point **x** from the observation space to the canonical space. Its deformation function is written as:

$$D(\mathbf{x}) = t_{\mathbf{x}}^c \cdot [u, v, w]^\top + \left\langle \mathbf{x} - t_{\mathbf{x}}^o \cdot [u, v, w]^\top, \mathbf{n}_{t_{\mathbf{x}}}^o \right\rangle \mathbf{n}_{t_{\mathbf{x}}}^c, \quad (58)$$

where $t_{\mathbf{x}}^c, t_{\mathbf{x}}^o$ are the 3D point's nearest triangle on the canonical and observed SMPL respectively, $\mathbf{n}_{t_{\mathbf{x}}}^c, \mathbf{n}_{t_{\mathbf{x}}}^o$ are their normal values, and $[u, v, w]$ are the barycentric coordinates. Equation (57) can then be reformulated as:

$$F((f_{xy} \circ D)(\mathbf{x}) + (f_{xz} \circ D)(\mathbf{x}) + (f_{yz} \circ D)(\mathbf{x})) \mapsto (\sigma, \mathbf{c}). \quad (59)$$

**ENARF-GAN** Following NARF [NSLH21] (see Sec. 4.2.3) that transforms the 3D point **x** into different parts of the local coordinate system, ENARF-GAN [NSLH22] applies triplane features to (1) provide features for estimating the color and density, and (2) predict the probability of a query point **x** belonging to a specific body part:

$$F(\mathbf{f}) \mapsto (\sigma, \mathbf{c}), \quad \mathbf{f} = \sum_{k=1}^{K} p^k * \mathbf{f}^k, \quad (60)$$

where $\mathbf{f}^k = \sum_{ij \in (xy, yz, xz)} F_{ij}(\mathbf{x}_c^k)$, $p^k = p_{xy}^k p_{xz}^k p_{yz}^k$, and $K$ is the number of body parts. After obtaining the RGB image via volume rendering, it further applies two StyleGAN2 generators separately for generating the foreground and background images, ensuring high fidelity of both the articulated human and its surrounding environment.

**HumanGen** To achieve detailed geometry and realistic 360° free-view rendering, HumanGen [JJW*22] combines a 3D human reconstruction prior with a 3D-GAN network, utilizing a disentangled optimization for geometry and texture. The process begins with a pre-trained 2D generator $G_{2D}$ to generate an anchor image based on the latent code $z$. The anchor image serves a dual purpose within the framework: (1) For geometry reconstruction, a pre-trained PIFuHD [SSSJ20] reconstructs the geometry based on SDF estimation using the anchor image. (2) For the texture branch, a triplane generator takes the latent code $z$ as input to generate color and blending weight; Later on, HumanGen integrates color, blending weight, and UV color information from the anchor image for each 3D query point to estimate texture values.

**EVA3D** Improving from NARF [NSLH21], which predicts a point's probability of belonging to a specific body part, EVA3D [HCL*22] takes a step further by dividing the human avatar into 16 parts based on the SMPL [LMR*15] model. Each part corresponds to a dedicated NeRF [MST*21] network, enabling localized modeling for enhanced accuracy and detail. Specifically, for each body part $k$, EVA3D applies a sub-network $F_k$ to model the local bounding box $\{\mathbf{b}_{min}^k, \mathbf{b}_{max}^k\}$. For a 3D query point **x** within the $k$-th bounding box, the density and color are predicted by:

$$F_k(\mathbf{x}_k) \mapsto (\sigma^k, \mathbf{c}^k), \quad \text{where} \quad \mathbf{x}_k = \frac{2\mathbf{x} - (\mathbf{b}_{min}^k + \mathbf{b}_{max}^k)}{\mathbf{b}_{max}^k - \mathbf{b}_{min}^k}. \quad (61)$$



|              | a) EG3D | b) StyleSDF | c) EVA3D |

**Figure 16:** *Qualitative comparison of GAN-based methods. Figure obtained from [HCL*22].*

If **x** falls in multiple bounding boxes, EVA3D uses a window function [LSS*21] to linearly blend the predicted values:

$$(\sigma, \mathbf{c}) = \frac{1}{\sum \omega_k} \sum_{k \in \mathbb{K}} \{\sigma^k, \mathbf{c}^k\}, \quad (62a)$$

$$\omega_k = \exp(-m(\mathbf{x}_k(x)^n + \mathbf{x}_k(y)^n + \mathbf{x}_k(z)^n)), \quad (62b)$$

where $m, n$ are chosen empirically. Note that EVA3D also employs the deformation field to transform **x** from the observation space to the canonical space. Please refer to their paper for detailed illustrations.

Besides the above methods, GET3D [GSW*22] utilizes a differentiable explicit surface extraction method to directly optimize textured 3D meshes and a differentiable rendering technique, which can be directly used by 3D rendering engines. Next3D [SWW*23] advances 3D generation further by learning generative neural textures based on parametric mesh templates and mapping them onto three triplanes through rasterization to achieve both deformation accuracy and topological flexibility. 3DAvatarGAN [ALZ*23], for the first time, offers generation, editing, and animation of personalized avatars obtained from a single image via a domain-adaption framework. TriPlaneNet [BNS24] proposes to directly operate in the triplane space instead of the GAN parameter space by building upon a feed-forward convolutional encoder for the latent code and extending it with a fully convolutional predictor of triplane numerical offsets.

**Table 2:** *Quantitative comparison of GAN-based methods. Results obtained from [HCL*22].*

| Methods | DeepFashion [LLQ*16] | | | | SHHQ [FLJ*22] | | | |
|---|---|---|---|---|---|---|---|---|
| | FID ↓ | KID ↓ | PCK ↑ | Depth ↓ | FID ↓ | KID ↓ | PCK ↑ | Depth ↓ |
| EG3D [CLC*22] | 26.38 | 0.014 | - | 0.0779 | 32.96 | 0.033 | - | 0.0296 |
| StyleSDF [OELS*22] | 92.40 | 0.136 | - | 0.0359 | 14.12 | 0.010 | - | 0.0300 |
| EVA3D [HCL*22] | 15.91 | 0.011 | 87.50 | 0.0272 | 11.99 | 0.009 | 88.95 | 0.0177 |

## 6.4. Discussion and Limitations of GAN-based Methods

Upon analyzing the quantitative evaluations presented in Tab. 2, it becomes evident that incorporating triplane in recent 3D-aware GAN networks significantly enhances performance across both 2D and 3D metrics. This conclusion is further supported by the qualitative visualizations depicted in Fig. 16. Notably, even compared to StyleSDF, EG3D produces suboptimal outcomes, underscoring the importance of integrating 3D priors such as SMPL within 3D-aware GAN architectures.

Unfortunately, while 3D-aware GAN techniques achieve promising 3D human generation results by training the network only on 2D datasets, they face challenges including (1) the presence of artifacts like blurs, distortions, and noise in the generated images, and (2) a dependency on the training datasets, which limits the ability to generate content beyond the scope of the training datasets. We will then dive into more generalizable 3D human generation methods that leverage large vision-language models like Contrastive Language-Image Pretraining (CLIP) [RKH*21] and diffusion models [Sta22, SCS*22] in the next section.

## 7. 3D Human Generation via Large Language Models

To exploit the potential of 2D generative image models such as Contrastive Language-Image Pretraining (CLIP) [RKH*21] and diffusion model [Sta22, ZA23, LWVH*23] for generating 3D contents, methods have been proposed to optimize the 3D representation (e.g., mesh, point cloud, NeRF, 3D Gaussians) based solely on the text prompts.

### 7.1. CLIP-based 3D Human Generation

In this subsection, we first discuss the methods that utilize the pre-trained CLIP model. Typically, CLIP is capable of mapping images and text to the same feature spaces, allowing for comparison of their similarities.

DreamField is one of the first methods that integrate 3D representations with pre-trained language models by evaluating the similarity between renderings generated from 3D contents and text prompts. Specifically, the proposed self-optimization manner involves: (1) initializing the NeRF to a unit sphere; (2) rendering RGB images from randomly sampled camera direction during each iteration; (3) calculating the similarity distance between the rendered image and text via CLIP as the loss function; (4) backpropagating the loss to optimize the NeRF parameters. Through iterative optimization, DreamField effectively aligns NeRF with the text prompt, thus bridging the gap between natural language and 3D content generation.

Following DreamField, Text2Mesh [MBOL*22] proposes to simultaneously render multiple images from various viewpoints at each iteration to enhance the performance. CLIP-Mesh [MKXBP22] introduces a set of render augmentations and incorporates a text-to-image embedding prior. CLIPX-Plore [HHL*23] maps the encoded CLIP code to its associated shape code to ensure a coherent connection between CLIP and shape latent spaces. Unfortunately, considering the complexity of human subjects, these methods often fall short in generating topologically and structurally correct 3D human models.

**AvatarCLIP** To facilitate efficient 3D human generation, AvatarCLIP [HZP*22] introduces a novel approach that involves initializing the NeuS [WLL*21] model with a predefined SMPL model, instead of the typical use of a sphere unit in CLIP-based optimization. Alongside 3D human generation, AvatarCLIP expands its method to generate motion sequences for 3D human animation based on text prompts. This process leverages a pre-calculated codebook and a pre-trained motion VAE (Variational Autoencoder) model [KW13]. By employing these resources, the pre-generated 3D human model can be animated using the SMPL-based deformation described in Equation (10).

Following AvataCLIP, MotionCLIP [TGH*22] trains a transformer-based motion auto-encoder to reconstruct motion while being aligned to its text label's position in the CLIP space. T2M-GPT [ZZC*23] takes the CLIP text embedding as the language prior for text-motion model training. AttT2M [ZHZX23] proposes body-part attention to learn a discrete latent space and global-local motion-text attention to learn the sentence and word level motion-text cross-modal relationship. Wu *et al.* [WZH*23] introduce the descriptive code space as an intermediary for the mapping from the text embedding space to the 3D face parametric space.

Unfortunately, despite CLIP's advantages of avoiding expensive and hard-to-obtain 3D datasets, it still struggles with creating realistic 3D meshes and lacks broad generalization in human motion. This limitation mainly arises from CLIP's constrained capacity to fully comprehend complex human languages.

### 7.2. Diffusion Model-based 3D Human Generation and Editing

Unlike CLIP-based [RKH*21] methods that embed image and text into a shared latent space to compare their similarities and learn associations, diffusion models convert random Gaussian noise into structured data by a Markov process with a series of denoising steps:

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t), \tag{63}$$

where $p_\theta(\mathbf{x}_{0:T})$ is the joint distribution over all the states of the process, $p(\mathbf{x}_T)$ is the distribution at the final time step $T$, which is typically assumed to be a standard Gaussian distribution and represents the data in its most noisy form. $\prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$ is the sequential multiplication of conditional distributions, each representing the probability distribution of the state at time step $t-1$ given the state at time step $t$.

While CLIP may lead to less precise interpretations of complex instructions, diffusion models, on the other hand, bypass linguistic ambiguities inherent in language-image pairings and focus on the step-by-step refinement of visual data, and perform better in 3D human representations.

Benefiting from the development of text-guided diffusion models $\phi$ [Sta22, BNH*22], DreamFusion [PJBM22] proposes a novel Score Distillation Sampling (SDS) to enable the generation of a 3D scene $g(\theta)$. Let's denote the RGB rendering from NeRF as $I$ and the text embedding as $y$. SDS strategy firstly involves encoding $I$ to derive the latent features $z$ and introducing random noise $\epsilon$ to $z$ to generate a noisy latent variable $z_t$. A pre-trained denoising function $\epsilon_\phi(z_t; y, t)$ is then employed to predict the added noise. The SDS loss is defined as the difference between predicted and added noise, and its gradient is given by:

$$\nabla_\theta \mathcal{L}_{\text{SDS}}(\phi, g(\theta)) = \mathbb{E}_{t,\epsilon}\left[ w(t) \left( \epsilon_\phi(z_t; y, t) - \epsilon \right) \frac{\partial z}{\partial x} \frac{\partial x}{\partial \theta} \right], \quad (64)$$

where $w(t)$ weights the loss from noise level $t$. The SDS gradients will be back-propagated to optimize $g(\theta)$, generating expressive 3D content from the text prompt.

Following DreamFusion, many other methods have been developed to improve the performance from different perspectives. Among them, Latent-NeRF [MRP*22] employs a latent diffusion model to optimize NeRF in the latent space, largely increasing the training efficiency. Magic3D [LGT*23] utilizes a coarse-to-fine strategy that leverages latent diffusion model and DMTET [SGY*21] for high-resolution 3D content generation. Fantasia3D [CCJJ23] disentangles the generation process into geometry and texture generation to enhance the performance. TEX-Ture [RMA*23] innovates in texture generation, transfer, and editing by using a pre-trained depth-to-image diffusion model and applying an iterative scheme that paints a given 3D model from different viewpoints. ProlificDreamer [WLW*23] introduces a Variational Score Distillation (VSD) to improve the quality and diversity of the generated 3D contents. MVDream [SWY*23] fine-tunes a multi-view diffusion model to produce consistent multi-view 3D generations. However, these methods consistently face similar limitations in generating 3D human bodies. For example, they cannot control human motions, and the results often lack hands or feet, with inconsistencies in geometry and texture.

### 7.2.1. 3D Human Generation

**DreamAvatar** DreamAvatar [CCH*23a] designs a dual-observation space in which the canonical space and observation space are jointly optimized using a shared NeRF module. The connection between these two spaces is established through a deformation field (as in Equation (25)), encompassing articulated deformation and non-rigid motion. Additionally, DreamAvatar utilizes SMPL-derived density fields that allow the optimized NeRF to evolve from the density field derived from the SMPL model:

$$\bar{\sigma}_c = \max(0, \texttt{softplus}^{-1}(\frac{1}{a}\texttt{sigmoid}(-d_c/a))), \quad (65a)$$

$$\bar{\sigma}_o = \max(0, \texttt{softplus}^{-1}(\frac{1}{a}\texttt{sigmoid}(-d_o/a))), \quad (65b)$$

where $\texttt{sigmoid}(x) = 1/(1 + e^{-x})$, $\texttt{softplus}^{-1}(x) = \log(e^x - 1)$, $d_c, d_o$ are signed distance to the corresponding SMPL model, and $a$ is a predefined hyperparameter [XWC*22]. The density and color values of DreamAvatar's dual-observation space are then derived by:

$$F(\mathbf{x}_c, \bar{\sigma}_c) = F_\theta(\gamma(\mathbf{x}_c)) + (\bar{\sigma}_c, \mathbf{0}) \mapsto (\sigma_c, \mathbf{c}_c), \quad (66a)$$

$$F(\mathbf{x}_o, \bar{\sigma}_o) = F_\theta(\gamma(\hat{\mathbf{x}}_c)) + (\bar{\sigma}_o, \mathbf{0}) \mapsto (\sigma_o, \mathbf{c}_o), \quad (66b)$$

where $\hat{\mathbf{x}}_c$ is $\mathbf{x}_o$'s corresponding canonical point. DreamAvatar efficiently enables the distillation of the well-optimized texture and geometry from the canonical space to the observation space, achieving high-quality and controllable avatar generation under user-guided human pose.

**AvatarCraft** Following DreamFusion and AvatarCLIP, Avatar-Craft [JWZ*23] leverages NeuS [WLL*21] (initialized from SMPL) and a pre-trained diffusion model to generate 3D human avatars. Specifically, AvatarCraft proposes to divide the canonical avatar into face and body bounding boxes according to the SMPL model and separately render them in a coarse-to-fine manner to recover higher-resolution texture and geometry. Since the optimization process closely aligns with the SMPL model, the generated human models can be readily animated using Equation (10).

**DreamWaltz** Instead of initializing the optimization process from SMPL or constraining it to SMPL surfaces, DreamWaltz [HWZ*23] introduces an alternative approach by replacing the Stable Diffusion model [Sta22] with a pose-conditioned ControlNet [ZA23]. In addition to rendering RGB images, the network simultaneously generates pose skeleton images from the SMPL model at each iteration. These skeleton images, along with RGB renderings and text embeddings, are input into ControlNet, ensuring consistent 3D gradients for the SDS loss. Its further proposes occlusion culling [PT02] for the skeleton images to remove invisible parts, addressing the multi-face "Janus" problem.

**DreamHuman** In contrast, DreamHuman [KAZ*23] integrates imGHUM [AXS21] as its deformable NeRF to enforce constraints on the human body. DreamHuman introduces a semantic zooming strategy to enhance the generation quality. Specifically, during each iteration, the network identifies different body parts (e.g., hands, head, arms, legs, etc.) and performs zoomed-in rendering separately for each part to calculate the SDS. The renderings capture intricate details and therefore help refine the texture and geometry.

**ZeroAvatar** Unlike other methods that utilize text prompts to generate 3D avatars, ZeroAvatar [WWY23] takes a single human image as input. Given a single image, ZeroAvatar (1) estimates the body pose, shape, and UV map to initialize the density field, (2) optimizes the geometry via SDS loss, integrating depth from the posed body model for enhanced accuracy, and (3) applies the inferred UVs to parts of the body that are not visible in the original image for a complete appearance. Specifically, for texture completion, it first uses DensePose [GNK18] to regress the UV coordinates from the image, and samples RGB colors to fill in the visible region of the UV map. The symmetrical areas of the visible region are then predicted and used as a prior during optimization, serving as a stronger guide than SDS.

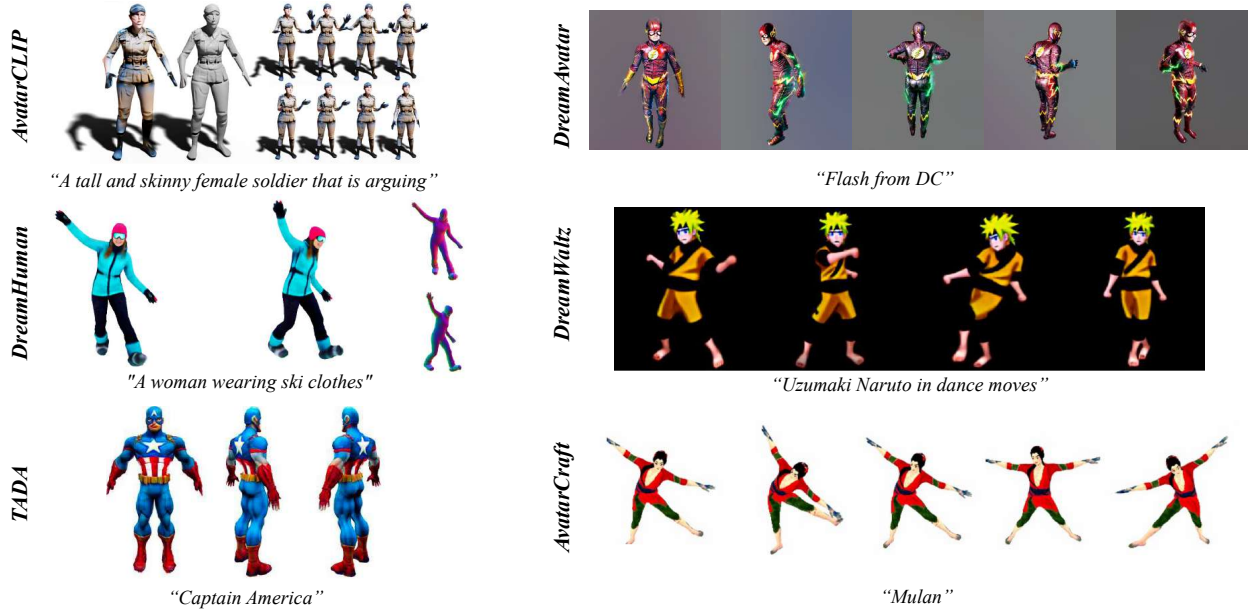**TADA** TADA [LYX*23] adopts 3D mesh representations and

**Figure 17:** *Results of various 3D human generation methods. Figure obtained from [HZP*22, CCH*23a, JWZ*23, HWZ*23, KAZ*23, LYX*23].*

directly initializes it from SMPL-X [PCG*19] surfaces, achieving robust 3D human generation and animation. After initializing from the SMPL-X surface with uniformly sampled points, TADA learns the SMPL-X parameters β, θ, ϕ and a displacement **D** which accounts for personalized details that are independent of pose, shape, and expression. Hence, the learnable 3D mesh can be written based on Equation (11):

$$M(\beta,\theta,\phi,\mathbf{D}) = \texttt{lbs}(\hat{T}(\beta,\theta,\phi,\mathbf{D}),J(\beta),\theta,\mathcal{W}), \quad (67a)$$

$$\hat{T}(\beta,\theta,\phi,\mathbf{D}) = \mathcal{S}(T(\beta,\theta,\phi)) + \mathbf{D}, \quad (67b)$$

where $\mathcal{S}(\cdot)$ is the mesh subdivision operation to add more details. To better align the geometry and texture for animation, it proposes to compute an additional SDS on the interpolation between normal and color image latents.

**AvatarBooth** To generate personalized 3D avatars, Avatar-Booth [ZLJ*23] introduces dual latent diffusion models to supervise the face and body generation separately. Specifically, Avatar-Booth needs a set of personalized images as input, and it first separates the input images into full body shots and headshots for fine-tuning two diffusion models. For the headshots, the pose-consistent constraint proposed in ControlNet [ZA23] is incorporated to ensure that the generated facial images maintain consistent identity across multiple views.

**HumanNorm** Considering that previous text-to-image diffusion models lack understanding of 3D structures, Human-Norm [HSZ*23] fine-tunes three diffusion models, i.e., text-to-normal diffusion model, text-to-depth diffusion model, and normal-aligned diffusion model, and adopts a disentangled optimization of the geometry and texture. It largely enhances the 2D perception of

the 3D geometry while ensuring the consistency between generated geometry and texture.

We provide several visualizations in Fig. 17 to offer further insights into the metrics employed by different methods for generating controllable 3D human avatars.

### 7.2.2. 3D Human Editing

**Control4D** Building upon Tensor4D [SZT*23], Control4D [SSP*23] introduces a 4D GAN architecture that can be edited via a ControlNet-based [ZA23] diffusion model, showcasing high-fidelity and consistent 4D editing based on 4D data and text prompts. It first employs Tensor4D to train the implicit representation of a 4D portrait scene from the 4D data, which is then rendered into latent features and RGB images via voxel rendering, serving as inputs to the generator. Meanwhile, the ControlNet takes original images and text prompts as inputs to produce edited images. The edited images serve as "real images" and the generator's outputs are "fake images", allowing for iterative 4D editing.

**HeadSculpt** HeadSculpt [HCH*23] adopts a coarse-to-fine strategy to generate high-resolution head avatars and perform fine-grained editing solely based on text prompts. To address the multi-face "Janus problem" in head generation, HeadSculpt proposes Prior-driven Score Distillation (PSD) that integrates 3D head priors and view-dependent textual inversion into the diffusion model. HeadSculpt further introduces Identity-aware Editing Score Distillation (IESD) that respects optimization gradients from both the original identity and instructive prompts, achieving fine-grained editing while maintaining its identity.

**HeadArtist** To address the inherent limitations of SDS (i.e.,

"A woman in a flowing sky blue sundress"  "An elderly woman in a cardigan and skirt"  "Lebron James"  "Joe Biden"  "Jason Statham"

"A boy with a beanie wearing a hoodie and joggers"  "A body builder wearing a tanktop"  "A viking"  "A black woman wearing sunglasses, a white t-shirt and jeans"  "A Texas ranger"

**Figure 18:** *Results of HumanNorm and HumanGaussian. Figure obtained from [HSZ*23, LZT*23].*

over-saturation and over-smoothing), HeadArtist [LWW*23] optimizes a parameterized 3D head model $R(\theta)$ under the supervision of the prior distillation itself, which is called Self Score Distillation (SSD):

$$\nabla_\theta \mathcal{L}_{\text{SSD}} = \mathbb{E}_{t,\epsilon} \left[ \omega(t) \left( \epsilon_\pi \left( x_t; y, t, c_L \right) - \hat{\epsilon}_\pi \left( x_t; y, t, c_L \right) \right) \frac{\partial x}{\partial \theta} \right], \quad (68)$$

where $c_L$ is the landmark that contains facial structure priors, $\epsilon_\pi$ and $\hat{\epsilon}_\pi$ are two pre-trained ControlNets with the same parameters, $x_t$ is the marginal distribution of the ControlNet given the text and landmarks.

**AvatarStudio** Evolving from Instruct-NeRF2NeRF [HTE*23], AvatarStudio [PET*23] achieves text-driven editing of dynamic head avatars by (1) fine-tuning pre-trained diffusion models on images with various viewpoints and time stamps, and (2) a novel view- and time-aware score distillation sampling (VT-SDS):

$$\nabla_x \mathcal{L}_{VT-SDS} = w(t) \left( \epsilon_t - \Psi \left( x_t, t, \mathbf{s}, \mathbf{s}_i \right) \right), \quad (69)$$

where $\mathbf{s}$ is the text embedding, $\mathbf{s}_i = \Gamma(\mathbf{P}_i)$ is a conditioning vector with $\mathbf{P}_i$ being the label assigned to the input image, and $\Psi(x_t, t, \mathbf{s}, \mathbf{s}_i)$ is the noise predicted by the diffusion model.

Besides, other works also achieve good results in 3D object and human editing. Progressive3D [CYW*23] decomposes the generation into a series of locally progressive editing steps to create precise 3D content for complex prompts. Vox-E [SFHAE23] introduces a novel volumetric regularization loss that operates directly in 3D space to maintain global coherence between the original and edited object. RODIN [WZZ*23] trains an image encoder to extract a semantic latent vector as the conditional input of the diffusion model, allowing semantic editing of generated results. DiffusionRig [DZX*23] learns generic and person-specific facial priors from extensive and individual datasets, respectively, for high-fidelity editing. TECA [ZFK*23] combines traditional 3D mesh models for the head, face, and upper body with NeRF for the modeling and editing of hair, clothing, and accessories.

## 7.3. 3D Gaussian-based 3D Human Generation and Editing

With the development of 3D Gaussian Splatting (3DGS), methods have been introduced to apply 3DGS for high-quality 3D generation with increased training efficiency.

Among them, DreamGaussian [TRZ*23] designs a mesh extraction algorithm from 3D Gaussians and a UV-space texture refinement for both efficiency and quality. GSGEN [CWL23] adopts a coarse-to-fine strategy for more delicate details and accurate geometry. GaussianDreamer [YFW*23] utilizes a 3D diffusion model to provide priors for initialization and a 2D diffusion model to enrich the geometry and appearance. LucidDreamer [CLN*23] leverages Stable Diffusion [Sta22], depth estimation, and explicit 3D representation for a domain-free high-quality 3D scene generation. GALA3D [ZRX*24] introduces a layout-guided Gaussian model and a compositional optimization mechanism to ensure geometry and texture consistency and accurate object interactions.

### 7.3.1. 3D Human Generation

**HumanGaussian** HumanGaussian [LZT*23] is an efficient and effective framework that combines Structure-Aware SDS and Annealed Negative Prompt Guidance for high-quality 3D human generation. In Structure-Aware SDS, HumanGaussian (1) utilizes SMPL-X [PCG*19] as a prior to densely sample Gaussians on the human mesh surface, (2) train a Texture-Structure Joint Model to simultaneously denoise the image and depth conditioned on the posed skeleton, and (3) design a dual-branch SDS to jointly optimize the appearance and geometry. In Annealed Negative Prompt Guidance, HumanGaussian uses the cleaner classifier score with an annealed negative score to regularize the stochastic SDS gradient of high variance. The floating artifacts are further eliminated based on Gaussian size in a prune-only phase to enhance generation smoothness.

Different from methods [HZP*22, KAZ*23, CCH*23a, HWZ*23, LYX*23, JWZ*23] which aim to control the generation of avatars in complex poses (see Fig. 17), HumanNorm and HumanGaussian focus primarily on pre-training diffusion models to enhance the alignment between geometry and texture. The generated results of their approach can be observed in Fig. 18.

**HeadStudio** HeadStudio [ZMFY24] achieves high-quality and animated head avatars differently by integrating FLAME [LBB*17] into 3D Gaussian splatting and SDS. It first proposes to deform 3D Gaussian points with facial expressions via FLAME-based 3D Gaussian Splatting (F-3DGS), where each 3D point is linked to a FLAME mesh and then rotated, scaled, and translated by the mesh deformation. Subsequently, it utilizes FLAME-based Score Distillation Sampling (F-SDS) which employs FLAME-based fine-grained control signals to guide the score distillation process. Finally, it enhances generation quality by applying uniform super-resolution and mesh regularization in F-3DGS.

We show several methods in Fig. 19 that are specifically designed for head avatars, which can efficiently generate or edit human heads.

**Figure 19:** *Generation or editing results of methods specifically designed for head avatars. Figure obtained from [HCH\*23, LWW\*23, ZMFY24].*

### 7.3.2. 3D Human Editing

3D editing methods based on implicit 3D representations like NeRF [MST\*21] suffer from slow processing speeds and limited control over complex scenes, while 3D Gaussian Splatting [KKLD23] offers an opportunity to efficiently locate user-guided semantics with enhanced efficiency.

**GaussianEditor** GaussianEditor [CCZ\*23] achieves precise local editing via Hierarchical Gaussian Splatting (HGS) and Gaussian Semantic Tracing. By providing a 2D segmentation mask with the editing area, Gaussian Semantic Tracing first identifies the corresponding parts in the 3D scene via back-projection, followed by assigning semantic tags to the affiliated Gaussians in these areas. HGS then records the semantic attributes of all the initialized Gaussians. By inheriting this attribute to child Gaussians during the densification and pruning, HGS enables more detailed and effective local editing. Besides precise local editing, GaussianEditor can also remove objects and integrate new objects via 3D impainting: 1) When an object needs to be removed, it identifies and isolates the object using Gaussian Semantic Tracing. 2) When an object needs to be added, it creates a 3D representation of the new object, which is then converted into the Gaussians that are compatible with the HGS system.

**SC-GS** Besides reconstruction illustrated in Section 5.2, SC-GS [HSY\*23] also allows for efficient motion editing by manipulating the learned control points. For each control point $p_j$, SC-GS calculates its trajectory $p_i^{\text{traj}}$ that includes its locations across $N_t (= 8)$ randomly sampled time steps as:

$$p_i^{\text{traj}} = \frac{1}{N_t} p_i^{t_1} \oplus p_i^{t_2} \oplus \cdots \oplus p_i^{t_{N_t}}, \tag{70}$$

where $\oplus$ denotes the vector concatenation operation. With the trained control points and the deformation MLP, SC-GS constructs



**Figure 20:** *Editing results of SC-GS. Figure obtained from [HSY\*23].*

a control point graph $\mathcal{P}'$ that connects control points based on their trajectories. An ARAP [SA07] deformation module is applied to the control graph to maintain local rigidity. Specifically, given a set of user-defined points $\left\{ h_l \in \mathbb{R}^3 \mid l \in \mathcal{H} \subset \{1, 2, \cdots, N_p\} \right\}$ ($N_p$ is the number of the control points), the control graph $\mathcal{P}'$ is deformed by minimizing the ARAP energy, which can be formulated as:

$$E\left(\mathcal{P}'\right) = \sum_{i=1}^{N_p} \sum_{j \in \mathcal{N}_i} w_{ij} \left\| \left(p_i' - p_j'\right) - \hat{R}_i \left(p_i - p_j\right) \right\|^2, \tag{71}$$

with fixed position condition $p_l' = h_l$ for $l \in \mathcal{H}$. $w_{ij}$ is the interpolation weights for control point $p_j$ and Gaussian $G_i$. $\hat{R}_i$ is the rigid local rotation defined on each control point. Examples of editing results are shown in Fig. 20.

**TIP-Editor** To compensate for the lack of accurate control over the specified appearance and location of the editing result in existing methods, TIP-Editor [ZKC\*24] uses a 3D bounding box to specify the editing region for 3D scene editing. It begins with a 2D stepwise personalization that incorporates (1) a scene personalization step with a localization loss to enhance the interaction between the existing 3D scene and the edited content, and (2) a content personalization step utilizing LoRA [HSW\*21]. TIP-Editor proposes a coarse editing stage via SDS, and a pixel-level reconstruction loss to refine the texture of the 3D scene following SDEdit [MHS\*21].

We present several visualizations in Fig. 21 to demonstrate the distinct metrics utilized by GaussianEditor and TIP-Editor for achieving locally controllable editing.

In addition to the above methods, Point'n Move [HY23] devises a two-stage self-prompting mask propagation process that produces 3D semantic segmentation masks from 2D image prompts for interactive scene object editing. Texture-GS [XHL\*24] separates appearance from geometry by mapping 2D texture onto the 3D surface. GaussCtrl [WBL\*24] proposes depth-conditioned editing and attention-based latent code alignment for multi-view consistent editing. VcEdit [WYW\*24] also aims to maintain consistency via a Cross-attention Consistency Module and an Editing Consistency Module.
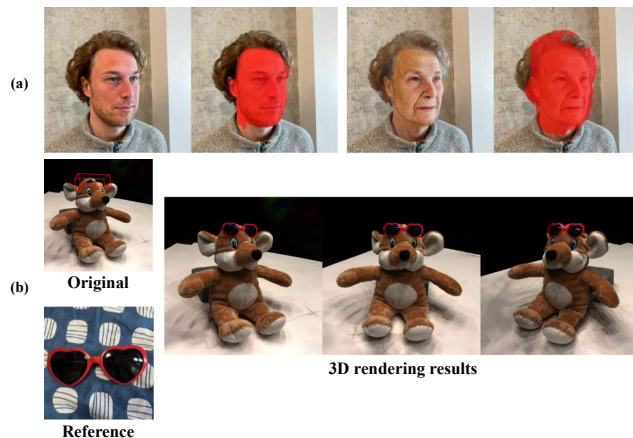
**Figure 21:** *Comparison of (a) GuassianEditor and (b) TIP-Editor. Figures obtained from [CCZ\*23, ZKC\*24].*

## 8. Reflection

Despite recent advancements in 3D human reconstruction and generation, there are still many areas in existing methods that demand further study. In this section, we explore some prominent perspectives with the aim of offering inspiration for advancing the field. While nowadays "reconstruction" and "generation" are always integrated to achieve good results, we categorize the discussion into optimization-based techniques and feed-forward pipelines.

### 8.1. Optimization-based Methods

#### 8.1.1. Integration of Image References and Text Prompts

As previously discussed, we can perceive NeRF [MST\*21] / 3DGS [KKLD23] and DreamFusion [PJBM22] as two extremes. NeRF and 3DGS rely on a multitude of images from diverse viewpoints to reconstruct the 3D scene, while DreamFusion [PJBM22] bridges the connection between images and text prompts via pretrained large language models and generates 3D content solely from text prompts. Consequently, NeRF and 3DGS excel in robust 3D scene reconstruction but are limited in their ability to generalize beyond trained specific subjects. On the other hand, DreamFusion and subsequent methods exhibit promising potential in generating previously unseen content distributions, but they face challenges related to quality, training efficiency, and robustness. Therefore, we aim to explore opportunities for integrating the advantages of both image references and text prompts.

We divide the discussion into three parts based on the input type and number.

**(1) Text prompt only.** Although AvatarCLIP [HZP\*22], DreamAvatar [CCH\*23a], AvatarCraft [JWZ\*23], DreamWaltz [HWZ\*23], TADA [LYX\*23], and following works [HSZ\*23, LZT\*23, ZLJ\*23] introduce several metrics to control the 3D avatars and ensure structural and topological accuracy, there remain challenges in:

*(I) Training efficiency and stability.* Existing methods always require hours to optimize a single 3D model, which is impractical

for real-world applications. One direction to improve the situation would be reducing the stochastic feature of the SDS. Methods like ProlificDreamer [WLW\*23], Consistent3D [WZY\*24], and SteinDreamer [WFX\*23] seek to minimize the variations of SDS, enhancing optimization robustness but not fully addressing training deficiency and stability issues. With the release of 3D large datasets, e.g, Objaverse [DSS\*23], Objaverse-XL [DLW\*24], and OmniObject3D [WZF\*23], recent methods [SWY\*23, SCZ\*23] have proposed to fine-tune the diffusion model for generating multi-view consistent images. Unfortunately, initialization is still the key factor that constrains their training efficiency and stability. Another direction for obvious improvement would be the introduction of retrieval techniques that can retrieve the related model from 3D datasets based on text prompts, and the retrieved 3D model will serve as the basis for further optimization.

*(II) Geometric quality.* Since SDS is more stochastic than reconstruction-based loss functions, existing methods always exhibit noisy geometry. The problem becomes more severe in Gaussian-based generation, as SDS gradients will directly impact the Gaussian position and attributes instead of the MLP in NeRF. One may consider improving the 3D representation based on existing approaches like Mosaic-SDF [YPN\*23], Flexi-Cube [SMH\*23], Hierarchical Gaussian Splatting [CCZ\*23], and SuGaR [GL23]. We believe the disentanglement of the position and texture attributes of 3D Gaussian Splatting will also provide a promising pathway.

*(III) Reliance on SMPL prior.* As discussed before, current methods rely heavily on SMPL [LMR\*15] and SMPL-X [PCG\*19], which, despite providing useful 3D priors, bring certain limitations. Specifically, the inherent characteristics of SMPL, e.g., containing only minimal clothing topology, largely constrain avatar customization, especially for the generation and animation of loose clothing. Learning from PIFu [SHN\*19], an alternative approach could involve pre-training an encoder that can perceive depth, silhouette, and semantic information.

**(2) Single image and text prompt.** With only a single image as input, the situation is similar to the single-view reconstruction in PIFu and related works, where the limited information in only one image poses issues, especially for human subjects where self-occlusion and depth-ambiguity often occur. To mitigate this, incorporating text prompts and pre-trained diffusion models could provide supplementary guidance for unseen views. In 3D object generation, methods like Zero-1-to-3 [LWVH\*23], One-2-3-45 [LXJ\*23], Magic123 [QMH\*23], and others [WCMB\*22, ZT23, CNC\*23, LSC\*23, MKRLV23, SJK\*23, LHG\*23, YWL\*23, LLZ\*23] estimate the image from other camera views via diffusion model, creating more robust 3D models from a single image and text prompts. TeCH [HYX\*23] proposes to apply the garment parsing model (SegFormer [XWY\*21]) and BLIP [LLXH22] to extract text prompts from the input image for further optimization. However, capturing all visual attributes to accurately reconstruct unseen areas and obtaining the correct pose remains a challenge. Thus, exploring better uses of diffusion models, e.g., adding more conditions and leveraging geometric correspondences, would be a promising future direction.

**(3) Multi-view images and text prompt.** Multi-view images

can provide more information in the 3D space but demand efficient feature fusion methods, a long-standing problem in 3D human reconstruction. Meanwhile, ensuring consistency between multi-view images and text prompts remains an open challenge. With the development of diffusion models, Guide3D [CCH*23b] generates multi-view images via ControlNet [ZA23] and textual inversion [GAA*23], while facing issues like inconsistent poses, texture, and orientations. Although Guide3D introduces a joint optimization of multi-resolution DMTET [SGY*21] grids, the geometric quality still lacks essential details. MVDream [SWY*23] proposes to train a diffusion model for generating multi-view images from text prompts but still falls short with human subjects. Moreover, both Guide3D and MVDream deal with sparse-view images. Thus, designing a diffusion model that is capable of generating consistent and dense multi-view images needs to be studied. We believe recent video diffusion models [BDK*23, BTCT*24, GZH*23, GZL*24] would provide inspiration for better spatial and temporal consistency.

### 8.1.2. 3D Human Editing

Given user preferences for privacy in virtual world applications like VR / AR, 3D human editing becomes a useful area of research. Existing methods for 3D human editing mainly detail under a global style, which are effective strategies to apply diffusion models but inevitably change remaining parts or the environment. TIP-Editor [ZKC*24] achieves local editing via the usage of a predefined bounding box. Yet, achieving automatic local editing represents a more intuitive and effective future direction.

### 8.1.3. 3D Human Animation

While 3D human generation can yield attractive results, the static 3D model cannot be readily applied in films, gaming, etc. Therefore, generating realistic human motion sequences and animating 3D human avatars based on large language models are promising research directions. Existing forms of human-motion generation [TRG*22, ZCP*24, YSI*23, ZGP*23, PBV23, JWF*24, ZMR*23] are all effective ways to make diffusion models more controllable in generating motion sequences from text prompts. However, these methods still fall short of real-world demands and realism. Meanwhile, generating human-object interaction motion sequences remains a hard open problem. On the other hand, by leveraging video diffusion models, DreamGaussian4D [RPT*23] and 4D-fy [BSR*23] attempt multi-stage animations. However, their scale of animation is still minimal, and they struggle to represent complex poses or movements accurately. Therefore, exploring user-guided and controllable control points as the anchor would be an encouraging direction in addressing these problems.

### 8.2. Feed-forward Methods

Feed-forward methods aim to pre-train a model on 3D datasets, offering text-based or image-based inference with much less time compared with optimization-based methods. While only a limited number of recent feed-forward methods specifically cater for human subjects, this section focuses on discussing existing methods and assessing their potential for 3D human modeling.

### 8.2.1. 3D Foundation Model

Recently, the advent of extensive 3D datasets such as Objaverse [DSS*23, DLW*24] and OmniObject3D [WZF*23] has fueled progress in 3D modeling techniques. Point-E [NJD*22] and Shap-E [JN23] train a 3D foundation model which can generate text-guided point clouds within minutes. MVDream [SWY*23], Zero-1-to-3 [LWVH*23], and Zero123++ [SCZ*23] put forth the idea of training diffusion models for generating consistent multi-view images using either a text prompt or a single image as input. Other methods, including SyncDreamer [LLZ*23], Wonder3D [LGL*23], One-2-3-45 [LXJ*23], UniDream [LLL*23], HexaGen3D [MNR*24], Sculpt3D [CYY*24], MVDiffusion++ [TCW*24], Make-Your-3D [LWC*24], and more [LCCT23, LSC*23, WS23, QMH*23, LXZ*23, KDJ*23, SWC*23], focus on achieving 3D-consistent generation with intricate details in geometry and texture from generated multi-view images. Additionally, the advancements in video diffusion models have led to the development of V3D [CWW*24], which leverages the temporal consistency characteristic of such models and fine-tunes pre-trained stable video diffusion models to generate dense multi-view images from single-view inputs. Yet, building a 3D feature space from the multi-view generated images to enhance the level of detail remains an unresolved challenge. Furthermore, accurately capturing the pose and clothing topologies of humans poses challenges due to the intricate nature of the human body. Hence, developing robust algorithms that can deal with both general objects and human subjects becomes crucial.

### 8.2.2. Large Reconstruction Model

Building on the success of transformers [VSP*17] and DINO [CTM*21, ODM*23], LRM (Large Reconstruction Model) [HZG*23] proposes a large transformer-based architecture to decode 3D triplane representation from DINO-encoded image features. Taking a single image as input, LRM reduces 3D modeling time to about 5 seconds. Following LRM, PF-LRM [WTB*23] achieves joint pose and shape prediction in 3D object reconstruction from a few unposed images. DMV-3D [XTL*23] also achieves fast 3D generation from text or a single image via a multi-view 2D image diffusion model and an LRM-based multi-view denoiser that reconstructs noise-free triplane NeRFs. Instant3D [LTZ*23] employs a two-stage approach that generates multi-view images from text prompts and reconstructs the 3D model via the large transformer-based module. More recently, TripoSR [TPL*24] enhances network quality and efficiency, and 3DTopia [HTC*24] introduces hybrid diffusion priors that include both text-conditioned triplane latent diffusion model and 2D diffusion priors to generate high-quality 3D objects. LGM [TCC*24], CRM [WWC*24], and GRM [YZW*24] explore various 3D representations, e.g., 3DGS [KKLD23], FlexiCube [SMH*23], for improved performance. However, similar problems in capturing the details and dynamics of human avatars as 3D foundation models can be observed. Although HumanLRM [WLT*24] proposes to distill multi-view reconstruction into single-view via a conditional triplane diffusion model for human subjects, its performance does not surpass traditional 3D implicit function-based human reconstruction methods.

### 8.2.3. Combination of Explicit Representaion and 3D Gaussian Splatting

While 3DGS [KKLD23] has demonstrated notable performance, it possesses a non-structural nature. The optimization process prioritizes updating the appearance of 3D Gaussians rather than directly moving them to the desired 3D locations. To enhance geometric optimization, researchers have been exploring methods that separately optimize the geometry (i.e., 3D point positions) and Gaussian attributes (texture). One such approach is Triplane-Meets-Gaussian-Splatting [ZYG*23], which employs two transformer-based networks, a point decoder and a triplane decoder. The point decoder generates point clouds from a single image, acting as the explicit 3D prior for the triplane decoder to predict Gaussian features for each point. Another method, AGG [XYM*24], decomposes the generation of 3D Gaussian locations and appearance attributes by first generating a coarse 3D representation and subsequently upsampling it via a 3D Gaussian super-resolution module. It is a promising direction to improve the performance of 3D human modeling while incorporating better 3D representation like Mosaic-SDF [YPN*23] would further improve the quality.

### 8.2.4. Real-time Gaussian-based Generation

Inspired by PIFu [SHN*19] that uses an image encoder and MLP for occupancy prediction based on pixel-aligned image features, GPS-Gaussian [ZZS*23] proposes to learn a 2D Gaussian parameter map from sparse-view RGB images of human-centered scenes to reconstruct free-viewpoint renderings. Notably, GPS-Gaussian achieves real-time rendering of dynamic scenes by efficiently querying the 2D Gaussian parameter map. However, directly applying this method to a 3D foundation model or large-scale reconstruction model is challenging due to the complexity of training the corresponding encoder and decoder. Finding solutions to incorporate 2D Gaussian parameter maps into generative pipelines remains a valuable open problem, given its potential for real-time 3D generation.

## References

[ALZ*23]  ABDAL R., LEE H.-Y., ZHU P., CHAI M., SIAROHIN A., WONKA P., TULYAKOV S.: 3davatargan: Bridging domains for personalized editable avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 4552–4562. 15

[APMTM19]  ALLDIECK T., PONS-MOLL G., THEOBALT C., MAGNOR M.: Tex2shape: Detailed full human body geometry from a single image. In *International Conference on Computer Vision* (2019). 1

[ASK*05]  ANGUELOV D., SRINIVASAN P., KOLLER D., THRUN S., RODGERS J., DAVIS J.: Scape: shape completion and animation of people. In *ACM SIGGRAPH 2005 Papers*. 2005, pp. 408–416. 1

[AXS21]  ALLDIECK T., XU H., SMINCHISESCU C.: imghum: Implicit generative models of 3d human shape and articulated pose. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 5461–5470. 17

[BDK*23]  BLATTMANN A., DOCKHORN T., KULAL S., MENDELEVITCH D., KILIAN M., LORENZ D., LEVI Y., ENGLISH Z., VOLETI V., LETTS A., ET AL.: Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127* (2023). 22

[BKY*22]  BERGMAN A., KELLNHOFER P., YIFAN W., CHAN E., LINDELL D., WETZSTEIN G.: Generative neural articulated radiance fields.

*Advances in Neural Information Processing Systems 35* (2022), 19900–19916. 2, 15

[BMV*22]  BARRON J. T., MILDENHALL B., VERBIN D., SRINIVASAN P. P., HEDMAN P.: Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *IEEE Conference on Computer Vision and Pattern Recognition* (2022). 9, 10

[BNH*22]  BALAJI Y., NAH S., HUANG X., VAHDAT A., SONG J., KREIS K., AITTALA M., AILA T., LAINE S., CATANZARO B., ET AL.: ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324* (2022). 17

[BNS24]  BHATTARAI A. R., NIESSNER M., SEVASTOPOLSKY A.: Triplanenet: An encoder for eg3d inversion. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision* (2024), pp. 3055–3065. 15

[BSB*07]  BALAN A. O., SIGAL L., BLACK M. J., DAVIS J. E., HAUSSECKER H. W.: Detailed human shape and pose from images. In *2007 IEEE Conference on Computer Vision and Pattern Recognition* (2007), IEEE, pp. 1–8. 1

[BSR*23]  BAHMANI S., SKOROKHODOV I., RONG V., WETZSTEIN G., GUIBAS L., WONKA P., TULYAKOV S., PARK J. J., TAGLIASACCHI A., LINDELL D. B.: 4d-fy: Text-to-4d generation using hybrid score distillation sampling. *arXiv preprint arXiv:2311.17984* (2023). 22

[BSTPM20]  BHATNAGAR B. L., SMINCHISESCU C., THEOBALT C., PONS-MOLL G.: Combining implicit function learning and parametric models for 3d human reconstruction. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16* (2020), Springer, pp. 311–329. 6

[BTCT*24]  BAR-TAL O., CHEFER H., TOV O., HERRMANN C., PAISS R., ZADA S., EPHRAT A., HUR J., LI Y., MICHAELI T., ET AL.: Lumiere: A space-time diffusion model for video generation. *arXiv preprint arXiv:2401.12945* (2024). 22

[BTTPM19]  BHATNAGAR B. L., TIWARI G., THEOBALT C., PONS-MOLL G.: Multi-garment net: Learning to dress 3d people from images. In *Proceedings of the IEEE/CVF international conference on computer vision* (2019), pp. 5420–5430. 1

[BV99]  BLANZ V., VETTER T.: A morphable model for the synthesis of 3d faces. In *Computer graphics and interactive techniques* (1999). 5, 13

[CBK03]  CHEUNG K., BAKER S., KANADE T.: Shape-from-silhouette of articulated objects and its use for human body kinematics estimation and motion capture. In *2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003. Proceedings.* (2003), vol. 1, IEEE, pp. I–I. 10

[CCH*22]  CAO Y., CHEN G., HAN K., YANG W., WONG K.-Y. K.: Jiff: Jointly-aligned implicit face function for high quality single view clothed human reconstruction. In *IEEE Conference on Computer Vision and Pattern Recognition* (2022). 1, 5

[CCH*23a]  CAO Y., CAO Y.-P., HAN K., SHAN Y., WONG K.-Y. K.: Dreamavatar: Text-and-shape guided 3d human avatar generation via diffusion models. *arXiv preprint arXiv:2304.00916* (2023). 2, 17, 18, 19, 21

[CCH*23b]  CAO Y., CAO Y.-P., HAN K., SHAN Y., WONG K.-Y. K.: Guide3d: Create 3d avatars from text and image guidance. *arXiv preprint arXiv:2308.09705* (2023). 22

[CCJJ23]  CHEN R., CHEN Y., JIAO N., JIA K.: Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. *arXiv preprint arXiv:2303.13873* (2023). 17

[CCS*15]  COLLET A., CHUANG M., SWEENEY P., GILLETT D., EVSEEV D., CALABRESE D., HOPPE H., KIRK A., SULLIVAN S.: High-quality streamable free-viewpoint video. *ACM Transactions on Graphics (ToG) 34*, 4 (2015), 1–13. 1

[CCZ*23]  CHEN Y., CHEN Z., ZHANG C., WANG F., YANG X., WANG Y., CAI Z., YANG L., LIU H., LIN G.: Gaussianeditor: Swift and controllable 3d editing with gaussian splatting. *arXiv preprint arXiv:2311.14521* (2023). 20, 21

[CH67] COVER T., HART P.: Nearest neighbor pattern classification. *IEEE transactions on information theory 13*, 1 (1967), 21–27. 13

[CHIS23] CROITORU F.-A., HONDRU V., IONESCU R. T., SHAH M.: Diffusion models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023). 3

[CHW23] CAO Y., HAN K., WONG K.-Y. K.: Sesdf: Self-evolved signed distance field for implicit 3d clothed human reconstruction. *arXiv preprint arXiv:2304.00359* (2023). 1, 6, 7

[CLC*22] CHAN E. R., LIN C. Z., CHAN M. A., NAGANO K., PAN B., MELLO S. D., GALLO O., GUIBAS L., TREMBLAY J., KHAMIS S., KARRAS T., WETZSTEIN G.: Efficient geometry-aware 3D generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition* (2022). 2, 14, 16

[CLN*23] CHUNG J., LEE S., NAM H., LEE J., LEE K. M.: Luciddreamer: Domain-free generation of 3d gaussian splatting scenes. *arXiv preprint arXiv:2311.13384* (2023). 19

[CLTS23] CHARATAN D., LI S., TAGLIASACCHI A., SITZMANN V.: pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. *arXiv preprint arXiv:2312.12337* (2023). 14

[CLZL22a] CHAN K., LIN G., ZHAO H., LIN W.: S-pifu: Integrating parametric human models with pifu for single-view clothed human reconstruction. *Advances in Neural Information Processing Systems 35* (2022), 17373–17385. 6

[CLZL22b] CHAN K. Y., LIN G., ZHAO H., LIN W.: Integratedpifu: Integrated pixel aligned implicit function for single-view human reconstruction. In *European conference on computer vision* (2022), Springer, pp. 328–344. 6

[CNC*23] CHAN E. R., NAGANO K., CHAN M. A., BERGMAN A. W., PARK J. J., LEVY A., AITTALA M., DE MELLO S., KARRAS T., WETZSTEIN G.: Genvs: Generative novel view synthesis with 3d-aware diffusion models, 2023. 21

[CPM20] CHIBANE J., PONS-MOLL G.: Implicit feature networks for texture completion from partial 3d data. In *European Conference on Computer Vision* (2020). 1, 6

[CSS*22] CAO X., SANTO H., SHI B., OKURA F., MATSUSHITA Y.: Bilateral normal integration. In *European Conference on Computer Vision* (2022), Springer. 6

[CTG*24] CAO H., TAN C., GAO Z., XU Y., CHEN G., HENG P.-A., LI S. Z.: A survey on generative diffusion models. *IEEE Transactions on Knowledge and Data Engineering* (2024). 3

[CTM*21] CARON M., TOUVRON H., MISRA I., JÉGOU H., MAIRAL J., BOJANOWSKI P., JOULIN A.: Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision* (2021), pp. 9650–9660. 22

[CW24] CHEN G., WANG W.: A survey on 3d gaussian splatting. *arXiv preprint arXiv:2401.03890* (2024). 3

[CWL23] CHEN Z., WANG F., LIU H.: Text-to-3d using gaussian splatting. *arXiv preprint arXiv:2309.16585* (2023). 19

[CWW*24] CHEN Z., WANG Y., WANG F., WANG Z., LIU H.: V3d: Video diffusion models are effective 3d generators. *arXiv preprint arXiv:2403.06738* (2024). 22

[CYW*23] CHENG X., YANG T., WANG J., LI Y., ZHANG L., ZHANG J., YUAN L.: Progressive3d: Progressively local editing for text-to-3d content creation with complex semantic prompts. *arXiv preprint arXiv:2310.11784* (2023). 19

[CYY*24] CHEN C., YANG X., YANG F., FENG C., FU Z., FOO C.-S., LIN G., LIU F.: Sculpt3d: Multi-view consistent text-to-3d generation with sparse 3d prior. *arXiv preprint arXiv:2403.09140* (2024). 22

[CZ19] CHEN Z., ZHANG H.: Learning implicit fields for generative shape modeling. In *IEEE Conference on Computer Vision and Pattern Recognition* (2019). 1

[DLJ*20] DENG B., LEWIS J. P., JERUZALSKI T., PONS-MOLL G., HINTON G., NOROUZI M., TAGLIASACCHI A.: Nasa neural articulated shape approximation. In *European Conference on Computer Vision* (2020). 1

[DLW*24] DEITKE M., LIU R., WALLINGFORD M., NGO H., MICHEL O., KUSUPATI A., FAN A., LAFORTE C., VOLETI V., GADRE S. Y., ET AL.: Objaverse-xl: A universe of 10m+ 3d objects. *Advances in Neural Information Processing Systems 36* (2024). 21, 22

[DSS*23] DEITKE M., SCHWENK D., SALVADOR J., WEIHS L., MICHEL O., VANDERBILT E., SCHMIDT L., EHSANI K., KEMBHAVI A., FARHADI A.: Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 13142–13153. 21, 22

[DWY*23] DAS D., WEWER C., YUNUS R., ILG E., LENSSEN J. E.: Neural parametric gaussians for monocular non-rigid object reconstruction. *arXiv preprint arXiv:2312.01196* (2023). 14

[DZX*23] DING Z., ZHANG X., XIA Z., JEBE L., TU Z., ZHANG X.: Diffusionrig: Learning personalized priors for facial appearance editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 12736–12746. 19

[FKYT*22] FRIDOVICH-KEIL S., YU A., TANCIK M., CHEN Q., RECHT B., KANAZAWA A.: Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 5501–5510. 12

[FLJ*22] FU J., LI S., JIANG Y., LIN K.-Y., QIAN C., LOY C. C., WU W., LIU Z.: Stylegan-human: A data-centric odyssey of human generation. In *European Conference on Computer Vision* (2022). 2, 14, 16

[FXZ*24] FEI B., XU J., ZHANG R., ZHOU Q., YANG W., HE Y.: 3d gaussian as a new vision era: A survey. *arXiv preprint arXiv:2402.07181* (2024). 3

[GAA*23] GAL R., ALALUF Y., ATZMON Y., PATASHNIK O., BERMANO A. H., CHECHIK G., COHEN-OR D.: An image is worth one word: Personalizing text-to-image generation using textual inversion. In *ICLR* (2023). 22

[GEVDM18] GRAHAM B., ENGELCKE M., VAN DER MAATEN L.: 3d semantic segmentation with submanifold sparse convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2018), pp. 9224–9232. 8

[GII*21] GRIGOREV A., ISKAKOV K., IANINA A., BASHIROV R., ZAKHARKIN I., VAKHITOV A., LEMPITSKY V.: Stylepeople: A generative model of fullbody human avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 5151–5160. 14

[GJC*23] GUO C., JIANG T., CHEN X., SONG J., HILLIGES O.: Vid2avatar: 3d avatar reconstruction from videos in the wild via self-supervised scene decomposition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 12858–12868. 8

[GL23] GUÉDON A., LEPETIT V.: Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering. *arXiv preprint arXiv:2311.12775* (2023). 14, 21

[GNK18] GÜLER R. A., NEVEROVA N., KOKKINOS I.: Densepose: Dense human pose estimation in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition* (2018). 17

[GPAM*14] GOODFELLOW I., POUGET-ABADIE J., MIRZA M., XU B., WARDE-FARLEY D., OZAIR S., COURVILLE A., BENGIO Y.: Generative adversarial nets. *Advances in neural information processing systems 27* (2014). 2, 14

[GPR*23] GOEL S., PAVLAKOS G., RAJASEGARAN J., KANAZAWA* A., MALIK* J.: Humans in 4D: Reconstructing and tracking humans with transformers. In *International Conference on Computer Vision (ICCV)* (2023). 9

[GSW*22] GAO J., SHEN T., WANG Z., CHEN W., YIN K., LI D., LITANY O., GOJCIC Z., FIDLER S.: Get3d: A generative model of high quality 3d textured shapes learned from images. In *Advances in Neural Information Processing Systems* (2022). 15

[GWH*20] GUO Y., WANG H., HU Q., LIU H., LIU L., BENNAMOUN M.: Deep learning for 3d point clouds: A survey. *IEEE transactions on pattern analysis and machine intelligence 43*, 12 (2020), 4338–4364. 3

[GYK*22] GAO X., YANG J., KIM J., PENG S., LIU Z., TONG X.: Mps-nerf: Generalizable 3d human rendering from multiview images. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022). 10

[GZH*23] GUO X., ZHENG M., HOU L., GAO Y., DENG Y., MA C., HU W., ZHA Z., HUANG H., WAN P., ET AL.: I2v-adapter: A general image-to-video adapter for video diffusion models. *arXiv preprint arXiv:2312.16693* (2023). 22

[GZL*24] GONG L., ZHU Y., LI W., KANG X., WANG B., GE T., ZHENG B.: Atomovideo: High fidelity image-to-video generation. *arXiv preprint arXiv:2403.01800* (2024). 22

[GZX*22] GAO X., ZHONG C., XIANG J., HONG Y., GUO Y., ZHANG J.: Reconstructing personalized semantic facial nerf models from monocular video. *ACM Transactions on Graphics (TOG) 41*, 6 (2022), 1–12. 8

[HCH*23] HAN X., CAO Y., HAN K., ZHU X., DENG J., SONG Y.-Z., XIANG T., WONG K.-Y. K.: Headsculpt: Crafting 3d head avatars with text. *arXiv preprint arXiv:2306.03038* (2023). 2, 18, 20

[HCJS20] HE T., COLLOMOSSE J., JIN H., SOATTO S.: Geo-pifu: Geometry and pixel aligned implicit functions for single-view human reconstruction. In *Advances on Neural Information Processing Systems* (2020). 1, 4

[HCL*22] HONG F., CHEN Z., LAN Y., PAN L., LIU Z.: Eva3d: Compositional 3d human generation from 2d image collections. *arXiv preprint arXiv:2210.04888* (2022). 2, 15, 16

[HHL*23] HU J., HUI K.-H., LIU Z., ZHANG H., FU C.-W.: Clipxplore: Coupled clip and shape spaces for 3d shape exploration. In *SIGGRAPH Asia 2023 Conference Papers* (2023), pp. 1–12. 16

[HHP*23] HU S., HONG F., PAN L., MEI H., YANG L., LIU Z.: Sherf: Generalizable human nerf from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 9352–9364. 10

[HL23] HU S., LIU Z.: Gauhuman: Articulated gaussian splatting from monocular human videos. *arXiv preprint arXiv:2312.02973* (2023). 2, 13

[HMZA21] HUANG X., MEI G., ZHANG J., ABBAS R.: A comprehensive survey on point cloud registration. *arXiv preprint arXiv:2103.02690* (2021). 3

[HSW*21] HU E. J., SHEN Y., WALLIS P., ALLEN-ZHU Z., LI Y., WANG S., WANG L., CHEN W.: Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685* (2021). 20

[HSY*23] HUANG Y.-H., SUN Y.-T., YANG Z., LYU X., CAO Y.-P., QI X.: Sc-gs: Sparse-controlled gaussian splatting for editable dynamic scenes. *arXiv preprint arXiv:2312.14937* (2023). 2, 13, 20

[HSZ*23] HUANG X., SHAO R., ZHANG Q., ZHANG H., FENG Y., LIU Y., WANG Q.: Humannorm: Learning normal diffusion model for high-quality and realistic 3d human generation. *arXiv preprint arXiv:2310.01406* (2023). 2, 18, 19, 21

[HTC*24] HONG F., TANG J., CAO Z., SHI M., WU T., CHEN Z., WANG T., PAN L., LIN D., LIU Z.: 3dtopia: Large text-to-3d generation model with hybrid diffusion priors. *arXiv preprint arXiv:2403.02234* (2024). 22

[HTE*23] HAQUE A., TANCIK M., EFROS A. A., HOLYNSKI A., KANAZAWA A.: Instruct-nerf2nerf: Editing 3d scenes with instructions. *arXiv preprint arXiv:2303.12789* (2023). 19

[HWZ*23] HUANG Y., WANG J., ZENG A., CAO H., QI X., SHI Y., ZHA Z.-J., ZHANG L.: Dreamwaltz: Make a scene with complex 3d animatable avatars. *arXiv preprint arXiv:2305.12529* (2023). 2, 17, 18, 19, 21

[HXL*20] HUANG Z., XU Y., LASSNER C., LI H., TUNG T.: Arch: Animatable reconstruction of clothed humans. In *IEEE Conference on Computer Vision and Pattern Recognition* (2020). 1, 5

[HXS*21] HE T., XU Y., SAITO S., SOATTO S., TUNG T.: Arch++: Animation-ready clothed human reconstruction revisited. In *International Conference on Computer Vision* (2021). 1, 5

[HY23] HUANG J., YU H.: Point'n move: Interactive scene object manipulation on gaussian splatting radiance fields. *arXiv preprint arXiv:2311.16737* (2023). 20

[HYX*23] HUANG Y., YI H., XIU Y., LIAO T., TANG J., CAI D., THIES J.: Tech: Text-guided reconstruction of lifelike clothed humans. *arXiv preprint arXiv:2308.08545* (2023). 21

[HZG*23] HONG Y., ZHANG K., GU J., BI S., ZHOU Y., LIU D., LIU F., SUNKAVALLI K., BUI T., TAN H.: Lrm: Large reconstruction model for single image to 3d. *arXiv preprint arXiv:2311.04400* (2023). 22

[HZJ*21] HONG Y., ZHANG J., JIANG B., GUO Y., LIU L., BAO H.: Stereopifu: Depth aware clothed human digitization via stereo vision. In *IEEE Conference on Computer Vision and Pattern Recognition* (2021). 1, 4, 6

[HZP*22] HONG F., ZHANG M., PAN L., CAI Z., YANG L., LIU Z.: Avatarclip: zero-shot text-driven generation and animation of 3d avatars. *ACM Transactions on Graphics (TOG)* (2022). 2, 16, 18, 19, 21

[JBAT17] JACKSON A. S., BULAT A., ARGYRIOU V., TZIMIROPOULOS G.: Large pose 3d face reconstruction from a single image via direct volumetric cnn regression. In *IEEE Conference on Computer Vision and Pattern Recognition* (2017). 4

[JHBZ22] JIANG B., HONG Y., BAO H., ZHANG J.: Selfrecon: Self reconstruction your digital avatar from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 5605–5615. 11

[JJW*22] JIANG S., JIANG H., WANG Z., LUO H., CHEN W., XU L.: Humangen: Generating human radiance fields with explicit priors. *arXiv preprint arXiv:2212.05321* (2022). 2, 15

[JMB*22] JAIN A., MILDENHALL B., BARRON J. T., ABBEEL P., POOLE B.: Zero-shot text-guided object generation with dream fields. In *IEEE Conference on Computer Vision and Pattern Recognition* (2022). 2

[JN23] JUN H., NICHOL A.: Shap-e: Generating conditional 3d implicit functions. *arXiv preprint arXiv:2305.02463* (2023). 22

[JWF*24] JIN P., WU Y., FAN Y., SUN Z., YANG W., YUAN L.: Act as you wish: Fine-grained control of motion diffusion model with hierarchical semantic graphs. *Advances in Neural Information Processing Systems 36* (2024). 22

[JWZ*23] JIANG R., WANG C., ZHANG J., CHAI M., HE M., CHEN D., LIAO J.: Avatarcraft: Transforming text into neural human avatars with parameterized shape and pose control. *arXiv preprint arXiv:2303.17606* (2023). 2, 17, 18, 19, 21

[JYQ*22] JIANG Y., YANG S., QIU H., WU W., LOY C. C., LIU Z.: Text2human: Text-driven controllable human image generation. *ACM TOG* (2022). 2

[JYS*22] JIANG W., YI K. M., SAMEI G., TUZEL O., RANJAN A.: Neuman: Neural human radiance field from a single video. In *European Conference on Computer Vision* (2022), Springer, pp. 402–418. 9

[KAL*21] KARRAS T., AITTALA M., LAINE S., HÄRKÖNEN E., HELLSTEN J., LEHTINEN J., AILA T.: Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems 34* (2021), 852–863. 14

[KAZ*23] KOLOTOUROS N., ALLDIECK T., ZANFIR A., BAZAVAN E. G., FIERARU M., SMINCHISESCU C.: Dreamhuman: Animatable

3d avatars from text. *arXiv preprint arXiv:2306.09329* (2023). 2, 17, 18, 19

[KDJ*23] KWAK J.-G., DONG E., JIN Y., KO H., MAHAJAN S., YI K. M.: Vivid-1-to-3: Novel view synthesis with video diffusion models. *arXiv preprint arXiv:2312.01305* (2023). 22

[KH13] KAZHDAN M., HOPPE H.: Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG) 32*, 3 (2013), 1–13. 6

[KKCF21] KWON Y., KIM D., CEYLAN D., FUCHS H.: Neural human performer: Learning generalizable radiance fields for human performance rendering. *Advances on Neural Information Processing Systems* (2021). 1, 8

[KKLD23] KERBL B., KOPANAS G., LEIMKÜHLER T., DRETTAKIS G.: 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (ToG) 42*, 4 (2023), 1–14. 2, 11, 12, 20, 21, 22, 23

[KLA*20] KARRAS T., LAINE S., AITTALA M., HELLSTEN J., LEHTINEN J., AILA T.: Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2020), pp. 8110–8119. 14

[KPD19] KOLOTOUROS N., PAVLAKOS G., DANIILIDIS K.: Convolutional mesh regression for single-image human shape reconstruction. In *IEEE Conference on Computer Vision and Pattern Recognition* (2019). 5

[KW13] KINGMA D. P., WELLING M.: Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013). 16

[LBB*17] LI T., BOLKART T., BLACK M. J., LI H., ROMERO J.: Learning a model of facial shape and expression from 4d scans. *ACM Transactions on Graphics (TOG)* (2017). 19

[LC98] LORENSEN W. E., CLINE H. E.: Marching cubes: A high resolution 3d surface construction algorithm. In *Seminal graphics: pioneering efforts that shaped the field*. 1998. 4

[LCCT23] LI W., CHEN R., CHEN X., TAN P.: Sweetdreamer: Aligning geometric priors in 2d diffusion for consistent text-to-3d. *arXiv preprint arXiv:2310.02596* (2023). 22

[LCY*23] LIU J.-W., CAO Y.-P., YANG T., XU E. Z., KEPPO J., SHAN Y., QIE X., SHOU M. Z.: Hosnerf: Dynamic human-object-scene neural radiance fields from a single video. *arXiv preprint arXiv:2304.12281* (2023). 1, 9

[LDZY23] LIN Y., DAI Z., ZHU S., YAO Y.: Gaussian-flow: 4d reconstruction with dynamic 3d gaussian particle. *arXiv preprint arXiv:2312.03431* (2023). 14

[LGL*23] LONG X., GUO Y.-C., LIN C., LIU Y., DOU Z., LIU L., MA Y., ZHANG S.-H., HABERMANN M., THEOBALT C., ET AL.: Wonder3d: Single image to 3d using cross-domain diffusion. *arXiv preprint arXiv:2310.15008* (2023). 22

[LGT*23] LIN C.-H., GAO J., TANG L., TAKIKAWA T., ZENG X., HUANG X., KREIS K., FIDLER S., LIU M.-Y., LIN T.-Y.: Magic3d: High-resolution text-to-3d content creation. In *IEEE Conference on Computer Vision and Pattern Recognition* (2023). 17

[LHG*23] LIN Y., HAN H., GONG C., XU Z., ZHANG Y., LI X.: Consistent123: One image to highly consistent 3d asset using case-aware diffusion priors. *arXiv preprint arXiv:2309.17261* (2023). 21

[LHQ*23] LIU Y., HUANG X., QIN M., LIN Q., WANG H.: Animatable 3d gaussian: Fast and high-quality reconstruction of multiple human avatars. *arXiv preprint arXiv:2311.16482* (2023). 2, 12

[LHR*21] LIU L., HABERMANN M., RUDNEV V., SARKAR K., GU J., THEOBALT C.: Neural actor: Neural free-view synthesis of human actors with pose control. *ACM SIGGRAPH Asia* (2021). 1, 9

[LLL*23] LIU Z., LI Y., LIN Y., YU X., PENG S., CAO Y.-P., QI X., HUANG X., LIANG D., OUYANG W.: Unidream: Unifying diffusion priors for relightable text-to-3d generation. *arXiv preprint arXiv:2312.08754* (2023). 22

[LLQ*16] LIU Z., LUO P., QIU S., WANG X., TANG X.: Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *IEEE Conference on Computer Vision and Pattern Recognition* (2016). 16

[LLXH22] LI J., LI D., XIONG C., HOI S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML* (2022). 21

[LLZ*23] LIU Y., LIN C., ZENG Z., LONG X., LIU L., KOMURA T., WANG W.: Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453* (2023). 21, 22

[LMR*15] LOPER M., MAHMOOD N., ROMERO J., PONS-MOLL G., BLACK M. J.: SMPL: A skinned multi-person linear model. *ACM Trans. Graphics, Asia* (2015). 1, 4, 8, 15, 21

[LSC*23] LIU M., SHI R., CHEN L., ZHANG Z., XU C., WEI X., CHEN H., ZENG C., GU J., SU H.: One-2-3-45++: Fast single image to 3d objects with consistent multi-view generation and 3d diffusion. *arXiv preprint arXiv:2311.07885* (2023). 21, 22

[LSS*21] LOMBARDI S., SIMON T., SCHWARTZ G., ZOLLHOEFER M., SHEIKH Y., SARAGIH J.: Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (ToG) 40*, 4 (2021), 1–13. 15

[LTV*22] LI R., TANKE J., VO M., ZOLLHÖFER M., GALL J., KANAZAWA A., LASSNER C.: Tava: Template-free animatable volumetric actors. In *European Conference on Computer Vision* (2022), Springer, pp. 419–436. 1, 9

[LTYY23] LI M., TAO J., YANG Z., YANG Y.: Human101: Training 100+ fps human gaussians in 100s from 1 view. *arXiv preprint arXiv:2312.15258* (2023). 2, 12, 13

[LTZ*23] LI J., TAN H., ZHANG K., XU Z., LUAN F., XU Y., HONG Y., SUNKAVALLI K., SHAKHNAROVICH G., BI S.: Instant3d: Fast text-to-3d with sparse-view generation and large reconstruction model. *arXiv preprint arXiv:2311.06214* (2023). 22

[LWC*24] LIU F., WANG H., CHEN W., SUN H., DUAN Y.: Make-your-3d: Fast and consistent subject-driven 3d content generation. *arXiv preprint arXiv:2403.09625* (2024). 22

[LWL*24] LIU X., WU C., LIU X., LIU J., WU J., ZHAO C., FENG H., DING E., WANG J.: Gea: Reconstructing expressive 3d gaussian avatar from monocular video. *arXiv preprint arXiv:2402.16607* (2024). 14

[LWVH*23] LIU R., WU R., VAN HOORICK B., TOKMAKOV P., ZAKHAROV S., VONDRICK C.: Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 9298–9309. 16, 21, 22

[LWW*23] LIU H., WANG X., WAN Z., SHEN Y., SONG Y., LIAO J., CHEN Q.: Headartist: Text-conditioned 3d head generation with self score distillation. *arXiv preprint arXiv:2312.07539* (2023). 2, 19, 20

[LXJ*23] LIU M., XU C., JIN H., CHEN L., XU Z., SU H., ET AL.: One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. *arXiv preprint arXiv:2306.16928* (2023). 21, 22

[LXZ*23] LORRAINE J., XIE K., ZENG X., LIN C.-H., TAKIKAWA T., SHARP N., LIN T.-Y., LIU M.-Y., FIDLER S., LUCAS J.: Att3d: Amortized text-to-3d object synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 17946–17956. 22

[LYX*23] LIAO T., YI H., XIU Y., TANG J., HUANG Y., THIES J., BLACK M. J.: Tada! text to animatable digital avatars. *arXiv preprint arXiv:2308.10899* (2023). 2, 17, 18, 19, 21

[LYX*24] LI M., YAO S., XIE Z., CHEN K., JIANG Y.-G.: Gaussianbody: Clothed human reconstruction via 3d gaussian splatting. *arXiv preprint arXiv:2401.09720* (2024). 14

[LZ21] LASSNER C., ZOLLHOFER M.: Pulsar: Efficient sphere-based neural rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 1440–1449. 11

[LZK*24] LI X., ZHANG Q., KANG D., CHENG W., GAO Y., ZHANG J., LIANG Z., LIAO J., CAO Y.-P., SHAN Y.: Advances in 3d generation: A survey. *arXiv preprint arXiv:2401.17807* (2024). 3

[LZT*23] LIU X., ZHAN X., TANG J., SHAN Y., ZENG G., LIN D., LIU X., LIU Z.: Humangaussian: Text-driven 3d human generation with gaussian splatting. *arXiv preprint arXiv:2311.17061* (2023). 19, 21

[LZW*23] LI C., ZHANG C., WAGHWASE A., LEE L.-H., RAMEAU F., YANG Y., BAE S.-H., HONG C. S.: Generative ai meets 3d: A survey on text-to-3d in aigc era. *arXiv preprint arXiv:2305.06131* (2023). 3

[MBOL*22] MICHEL O., BAR-ON R., LIU R., BENAIM S., HANOCKA R.: Text2mesh: Text-driven neural stylization for meshes. In *IEEE Conference on Computer Vision and Pattern Recognition* (2022). 2, 16

[MHS*21] MENG C., HE Y., SONG Y., SONG J., WU J., ZHU J.-Y., ERMON S.: Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073* (2021). 20

[MKRLV23] MELAS-KYRIAZI L., RUPPRECHT C., LAINA I., VEDALDI A.: Realfusion: 360 {\deg} reconstruction of any object from a single image. *arXiv preprint arXiv:2302.10663* (2023). 21

[MKXBP22] MOHAMMAD KHALID N., XIE T., BELILOVSKY E., POPA T.: Clip-mesh: Generating textured meshes from text using pre-trained image-text models. In *SIGGRAPH Asia 2022 conference papers* (2022), pp. 1–8. 16

[MNR*24] MERCIER A., NAKHLI R., REDDY M., YASARLA R., CAI H., PORIKLI F., BERGER G.: Hexagen3d: Stablediffusion is just one step away from fast and diverse text-to-3d generation. *arXiv preprint arXiv:2401.07727* (2024). 22

[MRP*22] METZER G., RICHARDSON E., PATASHNIK O., GIRYES R., COHEN-OR D.: Latent-nerf for shape-guided generation of 3d shapes and textures. *arXiv preprint arXiv:2211.07600* (2022). 17

[MST*21] MILDENHALL B., SRINIVASAN P. P., TANCIK M., BARRON J. T., RAMAMOORTHI R., NG R.: Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM* (2021). 1, 7, 10, 11, 14, 15, 20, 21

[MSVW23] MU J., SANG S., VASCONCELOS N., WANG X.: Actorsnerf: Animatable few-shot human rendering with generalizable nerfs. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 18391–18401. 10

[MYR*20] MA Q., YANG J., RANJAN A., PUJADES S., PONS-MOLL G., TANG S., BLACK M. J.: Learning to dress 3d people in generative clothing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2020), pp. 6469–6478. 6, 7

[NJD*22] NICHOL A., JUN H., DHARIWAL P., MISHKIN P., CHEN M.: Point-e: A system for generating 3d point clouds from complex prompts. *arXiv preprint arXiv:2212.08751* (2022). 22

[NPLT*19] NGUYEN-PHUOC T., LI C., THEIS L., RICHARDT C., YANG Y.-L.: Hologan: Unsupervised learning of 3d representations from natural images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2019), pp. 7588–7597. 14

[NPRM*20] NGUYEN-PHUOC T. H., RICHARDT C., MAI L., YANG Y., MITRA N.: Blockgan: Learning 3d object-aware scene representations from unlabelled images. *Advances in neural information processing systems 33* (2020), 6767–6778. 14

[NSLH21] NOGUCHI A., SUN X., LIN S., HARADA T.: Neural articulated radiance field. In *International Conference on Computer Vision* (2021). 1, 10, 15

[NSLH22] NOGUCHI A., SUN X., LIN S., HARADA T.: Unsupervised learning of efficient geometry-aware neural articulated representations. In *European Conference on Computer Vision* (2022). 2, 15

[NYD16] NEWELL A., YANG K., DENG J.: Stacked hourglass networks for human pose estimation. In *European Conference on Computer Vision* (2016). 4

[ODM*23] OQUAB M., DARCET T., MOUTAKANNI T., VO H., SZAFRANIEC M., KHALIDOV V., FERNANDEZ P., HAZIZA D., MASSA F., EL-NOUBY A., ET AL.: Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193* (2023). 22

[OELS*22] OR-EL R., LUO X., SHAN M., SHECHTMAN E., PARK J. J., KEMELMACHER-SHLIZERMAN I.: Stylesdf: High-resolution 3d-consistent image and geometry generation. In *IEEE Conference on Computer Vision and Pattern Recognition* (2022). 14, 16

[PBV23] PETROVICH M., BLACK M. J., VAROL G.: Tmr: Text-to-motion retrieval using contrastive 3d human motion synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 9488–9497. 22

[PCG*19] PAVLAKOS G., CHOUTAS V., GHORBANI N., BOLKART T., OSMAN A. A. A., TZIONAS D., BLACK M. J.: Expressive body capture: 3d hands, face, and body from a single image. In *IEEE Conference on Computer Vision and Pattern Recognition* (2019). 4, 8, 18, 19, 21

[PDW*21] PENG S., DONG J., WANG Q., ZHANG S., SHUAI Q., ZHOU X., BAO H.: Animatable neural radiance fields for modeling dynamic human bodies. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 14314–14323. 1, 9

[PET*23] PAN M. M., ELGHARIB M., TEOTIA K., TEWARI A., GOLYANIK V., KORTYLEWSKI A., THEOBALT C., ET AL.: Avatarstudio: Text-driven editing of 3d dynamic human head avatars. *arXiv preprint arXiv:2306.00547* (2023). 2, 19

[PFS*19] PARK J. J., FLORENCE P., STRAUB J., NEWCOMBE R., LOVEGROVE S.: Deepsdf: Learning continuous signed distance functions for shape representation. In *IEEE Conference on Computer Vision and Pattern Recognition* (2019). 13

[PJBM22] POOLE B., JAIN A., BARRON J. T., MILDENHALL B.: Dreamfusion: Text-to-3d using 2d diffusion. In *International Conference on Learning Representations* (2022). 2, 17, 21

[PNM*20] PENG S., NIEMEYER M., MESCHEDER L., POLLEFEYS M., GEIGER A.: Convolutional occupancy networks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16* (2020), Springer, pp. 523–540. 1

[PT02] PANTAZOPOULOS I., TZAFESTAS S.: Occlusion culling algorithms: A comprehensive survey. *Journal of Intelligent and Robotic Systems 35* (2002), 123–156. 17

[PYG*23] PO R., YIFAN W., GOLYANIK V., ABERMAN K., BARRON J. T., BERMANO A. H., CHAN E. R., DEKEL T., HOLYNSKI A., KANAZAWA A., ET AL.: State of the art on diffusion models for visual computing. *arXiv preprint arXiv:2310.07204* (2023). 3

[PZX*21] PENG S., ZHANG Y., XU Y., WANG Q., SHUAI Q., BAO H., ZHOU X.: Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *IEEE Conference on Computer Vision and Pattern Recognition* (2021). 1, 8

[QMH*23] QIAN G., MAI J., HAMDI A., REN J., SIAROHIN A., LI B., LEE H.-Y., SKOROKHODOV I., WONKA P., TULYAKOV S., ET AL.: Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. *arXiv preprint arXiv:2306.17843* (2023). 21, 22

[QSMG17] QI C. R., SU H., MO K., GUIBAS L. J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition* (2017). 5

[QWM*23] QIAN Z., WANG S., MIHAJLOVIC M., GEIGER A., TANG S.: 3dgs-avatar: Animatable avatars via deformable 3d gaussian splatting. *arXiv preprint arXiv:2312.09228* (2023). 14

[QYSG17] QI C. R., YI L., SU H., GUIBAS L. J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv preprint arXiv:1706.02413* (2017). 5

[Ren] https://renderpeople.com. 6, 7

[RH01] RAMAMOORTHI R., HANRAHAN P.: An efficient representation for irradiance environment maps. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques* (2001), pp. 497–500. 11

[RKH*21] RADFORD A., KIM J. W., HALLACY C., RAMESH A., GOH G., AGARWAL S., SASTRY G., ASKELL A., MISHKIN P., CLARK J.,

ET AL.: Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning* (2021). 2, 16

[RMA*23] RICHARDSON E., METZER G., ALALUF Y., GIRYES R., COHEN-OR D.: Texture: Text-guided texturing of 3d shapes. *arXiv preprint arXiv:2302.01721* (2023). 17

[RPT*23] REN J., PAN L., TANG J., ZHANG C., CAO A., ZENG G., LIU Z.: Dreamgaussian4d: Generative 4d gaussian splatting. *arXiv preprint arXiv:2312.17142* (2023). 22

[RRN*20] RAVI N., REIZENSTEIN J., NOVOTNY D., GORDON T., LO W.-Y., JOHNSON J., GKIOXARI G.: Accelerating 3d deep learning with pytorch3d. *arXiv preprint arXiv:2007.08501* (2020). 5

[SA07] SORKINE O., ALEXA M.: As-rigid-as-possible surface modeling. In *Symposium on Geometry processing* (2007), vol. 4, pp. 109–116. 20

[SCS*22] SAHARIA C., CHAN W., SAXENA S., LI L., WHANG J., DENTON E. L., GHASEMIPOUR K., GONTIJO LOPES R., KARAGOL AYAN B., SALIMANS T., ET AL.: Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems 35* (2022), 36479–36494. 16

[SCZ*23] SHI R., CHEN H., ZHANG Z., LIU M., XU C., WEI X., CHEN L., ZENG C., SU H.: Zero123++: a single image to consistent multi-view diffusion base model. *arXiv preprint arXiv:2310.15110* (2023). 21, 22

[SFHAE23] SELLA E., FIEBELMAN G., HEDMAN P., AVERBUCH-ELOR H.: Vox-e: Text-guided voxel editing of 3d objects. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 430–440. 19

[SGY*21] SHEN T., GAO J., YIN K., LIU M.-Y., FIDLER S.: Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis. *Advances in Neural Information Processing Systems 34* (2021), 6087–6101. 2, 17, 22

[SHN*19] SAITO S., HUANG Z., NATSUME R., MORISHIMA S., KANAZAWA A., LI H.: Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *International Conference on Computer Vision* (2019). 1, 2, 4, 6, 7, 10, 11, 21, 23

[SJK*23] SEO J., JANG W., KWAK M.-S., KIM H., KO J., KIM J., KIM J.-H., LEE J., KIM S.: Let 2d diffusion model know 3d-consistency for robust text-to-3d generation. *arXiv preprint arXiv:2303.07937* (2023). 21

[SJL*21] SUO X., JIANG Y., LIN P., ZHANG Y., WU M., GUO K., XU L.: Neuralhumanfvv: Real-time neural volumetric human performance rendering using rgb cameras. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2021), pp. 6226–6237. 1, 10

[SMH*23] SHEN T., MUNKBERG J., HASSELGREN J., YIN K., WANG Z., CHEN W., GOJCIC Z., FIDLER S., SHARP N., GAO J.: Flexible isosurface extraction for gradient-based mesh optimization. *ACM Transactions on Graphics (TOG) 42*, 4 (2023), 1–16. 21, 22

[SRV23] SZYMANOWICZ S., RUPPRECHT C., VEDALDI A.: Splatter image: Ultra-fast single-view 3d reconstruction. *arXiv preprint arXiv:2312.13150* (2023). 14

[SSC22] SUN C., SUN M., CHEN H.-T.: Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 5459–5469. 12

[SSP07] SUMNER R. W., SCHMID J., PAULY M.: Embedded deformation for shape manipulation. In *ACM siggraph 2007 papers*. 2007, pp. 80–es. 13

[SSP*23] SHAO R., SUN J., PENG C., ZHENG Z., ZHOU B., ZHANG H., LIU Y.: Control4d: Dynamic portrait editing by learning 4d gan from 2d diffusion-based editor. *arXiv preprint arXiv:2305.20082* (2023). 2, 18

[SSS06] SNAVELY N., SEITZ S. M., SZELISKI R.: Photo tourism: exploring photo collections in 3d. In *ACM siggraph 2006 papers*. 2006, pp. 835–846. 11

[SSSJ20] SAITO S., SIMON T., SARAGIH J., JOO H.: Pifuhd: Multilevel pixel-aligned implicit function for high-resolution 3d human digitization. In *IEEE Conference on Computer Vision and Pattern Recognition* (2020). 1, 4, 6, 7, 15

[Sta22] STABILITY.AI: Stable diffusion. https://stability.ai/blog/stable-diffusion-public-release, 2022. 16, 17, 19

[SWC*23] SHI Y., WANG J., CAO H., TANG B., QI X., YANG T., HUANG Y., LIU S., ZHANG L., SHUM H.-Y.: Toss: High-quality text-guided novel view synthesis from a single image. *arXiv preprint arXiv:2310.10644* (2023). 22

[SWW*23] SUN J., WANG X., WANG L., LI X., ZHANG Y., ZHANG H., LIU Y.: Next3d: Generative neural texture rasterization for 3d-aware head avatars. In *IEEE Conference on Computer Vision and Pattern Recognition* (2023). 15

[SWY*23] SHI Y., WANG P., YE J., MAI L., LI K., YANG X.: Mvdream: Multi-view diffusion for 3d generation. *arXiv:2308.16512* (2023). 17, 21, 22

[SZT*23] SHAO R., ZHENG Z., TU H., LIU B., ZHANG H., LIU Y.: Tensor4d: Efficient neural 4d decomposition for high-fidelity dynamic reconstruction and rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 16632–16642. 18

[SZZ*22a] SHAO R., ZHANG H., ZHANG H., CHEN M., CAO Y.-P., YU T., LIU Y.: Doublefield: Bridging the neural surface and radiance fields for high-fidelity human reconstruction and rendering. In *IEEE Conference on Computer Vision and Pattern Recognition* (2022). 1, 11

[SZZ*22b] SHAO R., ZHENG Z., ZHANG H., SUN J., LIU Y.: Diffustereo: High quality human reconstruction via diffusion-based stereo using sparse cameras. In *European Conference on Computer Vision* (2022). 1, 7

[TCC*24] TANG J., CHEN Z., CHEN X., WANG T., ZENG G., LIU Z.: Lgm: Large multi-view gaussian model for high-resolution 3d content creation. *arXiv preprint arXiv:2402.05054* (2024). 22

[TCW*24] TANG S., CHEN J., WANG D., TANG C., ZHANG F., FAN Y., CHANDRA V., FURUKAWA Y., RANJAN R.: Mvdiffusion++: A dense high-resolution multi-view diffusion model for single or sparse-view 3d object reconstruction. *arXiv preprint arXiv:2402.12712* (2024). 22

[TFT*20] TEWARI A., FRIED O., THIES J., SITZMANN V., LOMBARDI S., SUNKAVALLI K., MARTIN-BRUALLA R., SIMON T., SARAGIH J., NIESSNER M., ET AL.: State of the art on neural rendering. In *Computer Graphics Forum* (2020), vol. 39, Wiley Online Library, pp. 701–727. 3

[TGH*22] TEVET G., GORDON B., HERTZ A., BERMANO A. H., COHEN-OR D.: Motionclip: Exposing human motion generation to clip space. In *European Conference on Computer Vision* (2022), Springer, pp. 358–374. 16

[TPL*24] TOCHILKIN D., PANKRATZ D., LIU Z., HUANG Z., LETTS A., LI Y., LIANG D., LAFORTE C., JAMPANI V., CAO Y.-P.: Triposr: Fast 3d object reconstruction from a single image. *arXiv preprint arXiv:2403.02151* (2024). 22

[TRG*22] TEVET G., RAAB S., GORDON B., SHAFIR Y., COHEN-OR D., BERMANO A. H.: Human motion diffusion model. *arXiv preprint arXiv:2209.14916* (2022). 22

[TRZ*23] TANG J., REN J., ZHOU H., LIU Z., ZENG G.: Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653* (2023). 19

[TTM*22] TEWARI A., THIES J., MILDENHALL B., SRINIVASAN P., TRETSCHK E., YIFAN W., LASSNER C., SITZMANN V., MARTIN-BRUALLA R., LOMBARDI S., ET AL.: Advances in neural rendering. In *Computer Graphics Forum* (2022), vol. 41, Wiley Online Library, pp. 703–735. 3

[TZLW23] TIAN Y., ZHANG H., LIU Y., WANG L.: Recovering 3d human mesh from monocular images: A survey. *IEEE transactions on pattern analysis and machine intelligence* (2023). 3

[VPY*23] VO K., PHAM T.-T., YAMAZAKI K., TRAN M., LE N.: Dna: Deformable neural articulations network for template-free dynamic 3d human reconstruction from monocular rgb-d video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 3675–3684. 11

[VSP*17] VASWANI A., SHAZEER N., PARMAR N., USZKOREIT J., JONES L., GOMEZ A. N., KAISER Ł., POLOSUKHIN I.: Attention is all you need. *Advances in neural information processing systems 30* (2017). 22

[WBL*24] WU J., BIAN J.-W., LI X., WANG G., REID I., TORR P., PRISACARIU V. A.: Gaussctrl: Multi-view consistent text-driven 3d gaussian splatting editing. *arXiv preprint arXiv:2403.08733* (2024). 20

[WCMB*22] WATSON D., CHAN W., MARTIN-BRUALLA R., HO J., TAGLIASACCHI A., NOROUZI M.: Novel view synthesis with diffusion models. *arXiv preprint arXiv:2210.04628* (2022). 21

[WCS*22] WENG C.-Y., CURLESS B., SRINIVASAN P. P., BARRON J. T., KEMELMACHER-SHLIZERMAN I.: Humannerf: Free-viewpoint rendering of moving people from monocular video. In *IEEE Conference on Computer Vision and Pattern Recognition* (2022). 1, 2, 8

[WFX*23] WANG P., FAN Z., XU D., WANG D., MOHAN S., IANDOLA F., RANJAN R., LI Y., LIU Q., WANG Z., ET AL.: Steindreamer: Variance reduction for text-to-3d score distillation via stein identity. *arXiv preprint arXiv:2401.00604* (2023). 21

[WLL*21] WANG P., LIU L., LIU Y., THEOBALT C., KOMURA T., WANG W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *Advances on Neural Information Processing Systems* (2021). 16, 17

[WLT*24] WENG Z., LIU J., TAN H., XU Z., ZHOU Y., YEUNG-LEVY S., YANG J.: Single-view 3d human digitalization with large reconstruction models. *arXiv preprint arXiv:2401.12175* (2024). 22

[WLW*23] WANG Z., LU C., WANG Y., BAO F., LI C., SU H., ZHU J.: Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *arXiv preprint arXiv:2305.16213* (2023). 17, 21

[WS23] WANG P., SHI Y.: Imagedream: Image-prompt multi-view diffusion for 3d generation. *arXiv preprint arXiv:2312.02201* (2023). 22

[WSCKS23] WENG C.-Y., SRINIVASAN P. P., CURLESS B., KEMELMACHER-SHLIZERMAN I.: Personnerf: Personalized reconstruction from photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 524–533. 10

[WSGT22] WANG S., SCHWARZ K., GEIGER A., TANG S.: Arah: Animatable volume rendering of articulated human sdfs. In *European conference on computer vision* (2022), Springer, pp. 1–19. 1, 10

[WTB*23] WANG P., TAN H., BI S., XU Y., LUAN F., SUNKAVALLI K., WANG W., XU Z., ZHANG K.: Pf-lrm: Pose-free large reconstruction model for joint pose and shape prediction. *arXiv preprint arXiv:2311.12024* (2023). 22

[WTZ*21] WANG J., TAN S., ZHEN X., XU S., ZHENG F., HE Z., SHAO L.: Deep 3d human pose estimation: A review. *Computer Vision and Image Understanding 210* (2021), 103225. 3

[WWC*24] WANG Z., WANG Y., CHEN Y., XIANG C., CHEN S., YU D., LI C., SU H., ZHU J.: Crm: Single image to 3d textured mesh with convolutional reconstruction model. *arXiv preprint arXiv:2403.05034* (2024). 22

[WWY23] WENG Z., WANG Z., YEUNG S.: Zeroavatar: Zero-shot 3d avatar generation from a single image. *arXiv preprint arXiv:2305.16411* (2023). 2, 17

[WYW*24] WANG Y., YI X., WU Z., ZHAO N., CHEN L., ZHANG H.: View-consistent 3d editing with gaussian splatting. *arXiv preprint arXiv:2403.11868* (2024). 20

[WYZ*24] WU T., YUAN Y.-J., ZHANG L.-X., YANG J., CAO Y.-P., YAN L.-Q., GAO L.: Recent advances in 3d gaussian splatting. *arXiv preprint arXiv:2403.11134* (2024). 3

[WZF*23] WU T., ZHANG J., FU X., WANG Y., REN J., PAN L., WU W., YANG L., WANG J., QIAN C., ET AL.: Omniobject3d: Large-vocabulary 3d object dataset for realistic perception, reconstruction and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 803–814. 21, 22

[WZH*23] WU M., ZHU H., HUANG L., ZHUANG Y., LU Y., CAO X.: High-fidelity 3d face generation from natural language descriptions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 4521–4530. 16

[WZY*24] WU Z., ZHOU P., YI X., YUAN X., ZHANG H.: Consistent3d: Towards consistent high-fidelity text-to-3d generation with deterministic sampling prior. *arXiv preprint arXiv:2401.09050* (2024). 21

[WZZ*23] WANG T., ZHANG B., ZHANG T., GU S., BAO J., BALTRUSAITIS T., SHEN J., CHEN D., WEN F., CHEN Q., ET AL.: Rodin: A generative model for sculpting 3d digital avatars using diffusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2023), pp. 4563–4573. 19

[XCL*23] XU Y., CHEN B., LI Z., ZHANG H., WANG L., ZHENG Z., LIU Y.: Gaussian head avatar: Ultra high-fidelity head avatar via dynamic gaussians. *arXiv preprint arXiv:2312.03029* (2023). 2, 13

[XFC*23] XING Z., FENG Q., CHEN H., DAI Q., HU H., XU H., WU Z., JIANG Y.-G.: A survey on video diffusion models. *arXiv preprint arXiv:2310.10647* (2023). 3

[XFM22] XU T., FUJITA Y., MATSUMOTO E.: Surface-aligned neural radiance fields for controllable 3d human synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 15883–15892. 8

[XHL*24] XU T.-X., HU W., LAI Y.-K., SHAN Y., ZHANG S.-H.: Texture-gs: Disentangling the geometry and texture for 3d gaussian splatting editing. *arXiv preprint arXiv:2403.10050* (2024). 20

[XTL*23] XU Y., TAN H., LUAN F., BI S., WANG P., LI J., SHI Z., SUNKAVALLI K., WETZSTEIN G., XU Z., ET AL.: Dmv3d: Denoising multi-view diffusion using 3d large reconstruction model. *arXiv preprint arXiv:2311.09217* (2023). 22

[XWC*22] XU J., WANG X., CHENG W., CAO Y.-P., SHAN Y., QIE X., GAO S.: Dream3d: Zero-shot text-to-3d synthesis using 3d shape prior and text-to-image diffusion models. *arXiv preprint arXiv:2212.14704* (2022). 17

[XWY*21] XIE E., WANG W., YU Z., ANANDKUMAR A., ALVAREZ J. M., LUO P.: Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems 34* (2021), 12077–12090. 21

[XYC*23] XIU Y., YANG J., CAO X., TZIONAS D., BLACK M. J.: ECON: Explicit Clothed humans Obtained from Normals. In *IEEE Learning Transferable Visual Models From Natural Language Supervision Conference on Computer Vision and Pattern Recognition* (2023). 1, 6, 7, 13

[XYM*24] XU D., YUAN Y., MARDANI M., LIU S., SONG J., WANG Z., VAHDAT A.: Agg: Amortized generative 3d gaussians for single image to 3d. *arXiv preprint arXiv:2401.04099* (2024). 23

[XYTB22] XIU Y., YANG J., TZIONAS D., BLACK M. J.: ICON: Implicit Clothed humans Obtained from Normals. In *IEEE Conference on Computer Vision and Pattern Recognition* (2022). 1, 5, 6, 7

[YCL*23] YU Z., CHENG W., LIU X., WU W., LIN K.-Y.: Monohuman: Animatable human neural field from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 16943–16953. 10

[YFW*23] YI T., FANG J., WU G., XIE L., ZHANG X., LIU W., TIAN Q., WANG X.: Gaussiandreamer: Fast generation from text to 3d gaussian splatting with point cloud priors. *arXiv preprint arXiv:2310.08529* (2023). 19

[YPN*23]  YARIV L., PUNY O., NEVEROVA N., GAFNI O., LIPMAN Y.: Mosaic-sdf for 3d generative models. *arXiv preprint arXiv:2312.09222* (2023). 21, 23

[YSI*23]  YUAN Y., SONG J., IQBAL U., VAHDAT A., KAUTZ J.: Physdiff: Physics-guided human motion diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 16010–16021. 22

[YWL*23]  YE J., WANG P., LI K., SHI Y., WANG H.: Consistent-1-to-3: Consistent image to 3d view synthesis via geometry-aware diffusion models. *arXiv preprint arXiv:2310.03020* (2023). 21

[YWRR23]  YANG G., WANG C., REDDY N. D., RAMANAN D.: Reconstructing animatable categories from videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 16995–17005. 9

[YYTK21]  YU A., YE V., TANCIK M., KANAZAWA A.: pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 4578–4587. 11

[YYZ*23]  YANG G., YANG S., ZHANG J. Z., MANCHESTER Z., RAMANAN D.: Ppr: Physically plausible reconstruction from monocular videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 3914–3924. 9

[YZG*21]  YU T., ZHENG Z., GUO K., LIU P., DAI Q., LIU Y.: Function4d: Real-time human volumetric capture from very sparse consumer rgbd sensors. In *IEEE Conference on Computer Vision and Pattern Recognition* (2021). 1, 11

[YZS*23]  YANG L., ZHANG Z., SONG Y., HONG S., XU R., ZHAO Y., ZHANG W., CUI B., YANG M.-H.: Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys 56*, 4 (2023), 1–39. 3

[YZW*24]  YINGHAO X., ZIFAN S., WANG Y., HANSHENG C., CEYUAN Y., SIDA P., YUJUN S., GORDON W.: Grm: Large gaussian reconstruction model for efficient 3d reconstruction and generation, 2024. `arXiv:2403.14621`. 22

[ZA23]  ZHANG L., AGRAWALA M.: Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543* (2023). 16, 17, 18, 22

[ZAB*22]  ZHENG Y., ABREVAYA V. F., BÜHLER M. C., CHEN X., BLACK M. J., HILLIGES O.: Im avatar: Implicit morphable head avatars from videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 13545–13555. 8

[ZBS*23]  ZIELONKA W., BAGAUTDINOV T., SAITO S., ZOLLHÖFER M., THIES J., ROMERO J.: Drivable 3d gaussian avatars. *arXiv preprint arXiv:2311.08581* (2023). 2, 13

[ZCP*24]  ZHANG M., CAI Z., PAN L., HONG F., GUO X., YANG L., LIU Z.: Motiondiffuse: Text-driven human motion generation with diffusion model. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2024). 22

[ZFK*23]  ZHANG H., FENG Y., KULITS P., WEN Y., THIES J., BLACK M. J.: Text-guided generation and editing of compositional 3d avatars. *arXiv preprint arXiv:2309.07125* (2023). 19

[ZGP*23]  ZHANG M., GUO X., PAN L., CAI Z., HONG F., LI H., YANG L., LIU Z.: Remodiffuse: Retrieval-augmented motion diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 364–373. 22

[ZHZX23]  ZHONG C., HU L., ZHANG Z., XIA S.: Attt2m: Text-driven human motion generation with multi-perspective attention mechanism. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 509–519. 16

[ZJY*22]  ZHANG J., JIANG Z., YANG D., XU H., SHI Y., SONG G., XU Z., WANG X., FENG J.: Avatargen: A 3d generative model for animatable human avatars. In *Arxiv* (2022). 2

[ZKC*24]  ZHUANG J., KANG D., CAO Y.-P., LI G., LIN L., SHAN Y.: Tip-editor: An accurate 3d editor following both text-prompts and image-prompts. *arXiv preprint arXiv:2401.14828* (2024). 20, 21, 22

[ZLJ*23]  ZENG Y., LU Y., JI X., YAO Y., ZHU H., CAO X.: Avatarbooth: High-quality and customizable 3d human avatar generation. *arXiv preprint arXiv:2306.09864* (2023). 2, 18, 21

[ZMFY24]  ZHOU Z., MA F., FAN H., YANG Y.: Headstudio: Text to animatable head avatars with 3d gaussian splatting. *arXiv preprint arXiv:2402.06149* (2024). 19, 20

[ZMR*23]  ZHU W., MA X., RO D., CI H., ZHANG J., SHI J., GAO F., TIAN Q., WANG Y.: Human motion generation: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023). 22

[ZRX*24]  ZHOU X., RAN X., XIONG Y., HE J., LIN Z., WANG Y., SUN D., YANG M.-H.: Gala3d: Towards text-to-3d complex scene generation via layout-guided generative gaussian splatting. *arXiv preprint arXiv:2402.07207* (2024). 19

[ZSZ*21]  ZHENG Y., SHAO R., ZHANG Y., YU T., ZHENG Z., DAI Q., LIU Y.: Deepmulticap: Performance capture of multiple characters using sparse multiview cameras. *arXiv preprint arXiv:2105.00261* (2021). 1, 6

[ZT23]  ZHOU Z., TULSIANI S.: Sparsefusion: Distilling view-conditioned diffusion for 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 12588–12597. 21

[ZWC*23]  ZHENG C., WU W., CHEN C., YANG T., ZHU S., SHEN J., KEHTARNAVAZ N., SHAH M.: Deep learning-based human pose estimation: A survey. *ACM Computing Surveys 56*, 1 (2023), 1–37. 3

[ZYG*23]  ZOU Z.-X., YU Z., GUO Y.-C., LI Y., LIANG D., CAO Y.-P., ZHANG S.-H.: Triplane meets gaussian splatting: Fast and generalizable single-view 3d reconstruction with transformers. *arXiv preprint arXiv:2312.09147* (2023). 23

[ZYLD21]  ZHENG Z., YU T., LIU Y., DAI Q.: Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021). 1, 5, 6, 7

[ZZC*23]  ZHANG J., ZHANG Y., CUN X., HUANG S., ZHANG Y., ZHAO H., LU H., SHEN X.: T2m-gpt: Generating human motion from textual descriptions with discrete representations. *arXiv preprint arXiv:2301.06052* (2023). 16

[ZZS*23]  ZHENG S., ZHOU B., SHAO R., LIU B., ZHANG S., NIE L., LIU Y.: Gps-gaussian: Generalizable pixel-wise 3d gaussian splatting for real-time human novel view synthesis. *arXiv preprint arXiv:2312.02155* (2023). 23