# HumanGen: Generating Human Radiance Fields with Explicit Priors

Suyi Jiang[1]    Haoran Jiang[1]    Ziyu Wang[1]    Haimin Luo[1]    Wenzheng Chen[2]
Lan Xu[1,3]

[1]ShanghaiTech University     [2]University of Toronto     [3]Shanghai Engineering Research Center of Intelligent Vision and Imaging

## Abstract

*Recent years have witnessed the tremendous progress of 3D GANs for generating view-consistent radiance fields with photo-realism. Yet, high-quality generation of human radiance fields remains challenging, partially due to the limited human-related priors adopted in existing methods. We present HumanGen, a novel 3D human generation scheme with detailed geometry and 360° realistic free-view rendering. It explicitly marries the 3D human generation with various priors from the 2D generator and 3D reconstructor of humans through the design of "anchor image". We introduce a hybrid feature representation using the anchor image to bridge the latent space of HumanGen with the existing 2D generator. We then adopt a pronged design to disentangle the generation of geometry and appearance. With the aid of the anchor image, we adapt a 3D reconstructor for fine-grained details synthesis and propose a two-stage blending scheme to boost appearance generation. Extensive experiments demonstrate our effectiveness for state-of-the-art 3D human generation regarding geometry details, texture quality, and free-view performance. Notably, HumanGen can also incorporate various off-the-shelf 2D latent editing methods, seamlessly lifting them into 3D.*

## 1. Introduction

We are entering an era where the boundaries of real and virtually generated worlds are dismissing. An epitome of this revolution is the recent rise of 3D-aware and photo-realistic image synthesis in the past several years [5, 6, 10, 15, 51, 61, 88], which combine 2D Generative Adversarial Networks (GANs) with neural volume rendering, like neural radiance fields (NeRFs) [42]. But such 3D GANs mainly focus on 3D rigid contents like human/animal faces or CAD models. The further 3D generation of us humans with photo-realism is more attractive, with numerous applications in VR/AR or visual effects.

High-quality 3D human generative models should ideally generate 3D-aware humans with the following charac-



Figure 1. The proposed HumanGen can generate 3D humans with fine-detailed geometry and appearance while seamlessly lifting various 2D latent editing tools into 3D.

teristics: (1) detailed geometry, (2) photo-realistic appearance, and (3) even supporting 360° free-view rendering. Yet, it remains extremely challenging, mainly due to the significantly higher diversity of human apparel and skeletal pose. Only very recently, a few work [3, 21, 83] explore 3D GANs for human generation by using the parametric human model like SMPL [38] as motion priors. But such parametric human prior lacks sufficient geometry details, and the adopted neural rendering in these methods does not guarantee that meaningful 3D geometry can be generated, further leading to appearance artifacts. Besides, these 3D human generators are trained with limited human datasets that lack diversity [66] or suffer from imbalanced viewing angles (most are front views) [12, 37]. In a nutshell, existing methods fail to fulfill all the aforementioned three characteristics for 3D human generation.

We observe that 3D human generation can benefit from more explicit priors from other research domains of human modeling, except for the SMPL prior adopted in existing methods. Specifically, with the recent large-scale dataset SHHQ [12], the 2D human generators [28–30] achieve more decent synthesis results than the 3D ones. And various downstream 2D editing tools are available by disentangling the latent spaces [53, 56, 63, 76]. These abilities of 2D generation and subsequent can significantly benefit the 3D human generation if their latent spaces can be bridged. Be-

sides, recent advances in monocular 3D human reconstruction [2, 58] have achieved more fine-grained geometry details than the implicit geometry proxy in current 3D human generators. Yet, there lacks a well-designed mechanism to explicitly utilize the rich human priors from both 2D generator and 2D reconstructor for 3D human generation.

In this paper, we present *HumanGen* – a novel neural scheme to generate high-quality radiance fields for 3D humans from 2D images, as shown in Fig. 1. In stark contrast with existing methods that only use SMPL, our approach explicitly utilizes richer priors from the top-tier 2D generation and 3D reconstruction schemes. As a result, our approach not only enables photo-realistic human generation with detailed geometry and 360° free-view ability, but also maintains the compatibility to existing off-the-shelf 2D editing toolbox based on latent disentanglement.

Our key idea in HumanGen is to organically leverage a 2D human generator and a 3D human reconstructor as explicit priors into a 3D GAN-like framework. Specifically, we first introduce a hybrid feature representation of the generative 3D space, which consists of the tri-plane features from EG3D [5] as well as a 2D photo-real human image (denoted as "anchor image") generated through the pre-trained 2D human generator. Note that we adopt separated Style-GAN2 [30] architectures to generate both the tri-plane feature maps and the anchor image. But they share the same latent mapping network, so as to bridge and anchor the latent space of our 3D GAN to the pre-trained 2D human generator. Then, based on such hybrid representation, we design our 3D human generator into the pronged geometry and appearance branches. In the geometry branch, we explicitly utilize a pre-trained 3D reconstructor PIFuHD [58] to extract pixel-aligned features from the anchor image and provide extra fine-grained geometry supervision for our HumanGen. Note that the original PIFuHD encodes geometry as an implicit occupancy field. Thus, we propose a geometry adapting scheme to turn it into a generative version with signed distance field (SDF) output, so as to support efficient and high-resolution volume rendering with sphere tracing. For the appearance branch, we propose to learn an appearance field and a blending field from both the pixel-aligned and tri-plane features. Then, we adopt a two-stage blending scheme to make full use of the rich texture information in the anchor image. For our 3D GAN training procedure, we adopt a similar training strategy like EG3D [5] with pose conditioning, and introduce additional front-view and back-view consistency supervision to enhance the generated texture details.

Besides, we observe that existing 2D human generator StyleGAN2 [30] trained on the large-scale SHHQ [12] can potentially generate diverse human images including side-views and even back-views. Thus, we train our HumanGan using an augmented dataset from SHHQ by using the pre-trained 2D generator to cover 360° viewing angles. Once trained, our HumanGen enables high-quality 3D human generation. As an additional benefit, it shares the same latent mapping with the 2D generated anchor image. Thus, using the anchor image, we can seamlessly upgrade off-the-shelf 2D latent editing methods into our 3D setting. We showcase various 3D effects via convenient anchor image editing. To summarize, our main contributions include:

- We present a novel 3D-aware human generation scheme, with detailed geometry and 360° realistic free-view rendering, achieving significant superiority to state-of-the-arts.

- We propose a hybrid feature representation using an anchor image with shared latent space to bridge our 3D GAN with the existing 2D generator.

- We propose a pronged design for appearance/geometry branches, and adapt a 3D reconstructor to aid the geometry branch for fine-grained details synthesis.

- We introduce an implicit blending field with two-stage blending strategy to generate high-quality appearance.

## 2. Related Work

**3D-aware GAN.** Early approaches mainly exploit explicit 3D representations for 3D-aware image synthesis, such as textured mesh [35, 50, 65], and voxels [13, 18, 19, 45, 46, 88]. These works commonly suffer from low model expressiveness or high memory footprint. The recent NeRF [42] has exerted tremendous momentum towards view-consistent 3D content generation [5–7, 10, 15, 16, 25, 48, 51, 61, 75, 77, 78, 87, 88]. However, low-resolution 3D volume generation is adopted [6, 61] and compensated with 2D upsampling layers [5, 15, 48, 51, 78], and limits the 3D view consistency and geometry detail. Hybrid representations, e.g., MPIs [10, 77, 87] partially address this issue but limit the visible range. Compared to these methods mainly focusing on the face or synthetic data [22], quite a few recent works [3, 5, 21, 83] explore 3D human generation using only 2D images. They tend to utilize SMPL [38] pose prior [3, 83] and coarse shape prior [21] to address the challenging human geometry diversity. However, such priors without geometry details cannot guarantee such generators generate faithful 3D geometry. Besides, the lack of view-balanced human datasets further limits their ability to around-view rendering. In contrast, our approach achieves generating high-quality freely-renderable human radiance fields with detailed geometry.

**2D Human Image Generation.** A large part of 2D image generation works fall in conditional image generation [24, 52, 73, 74, 89] and achieves a photo-realistic level. As to human image synthesis, a line of work focuses on conditioning the generator with semantic map [11, 27, 85],
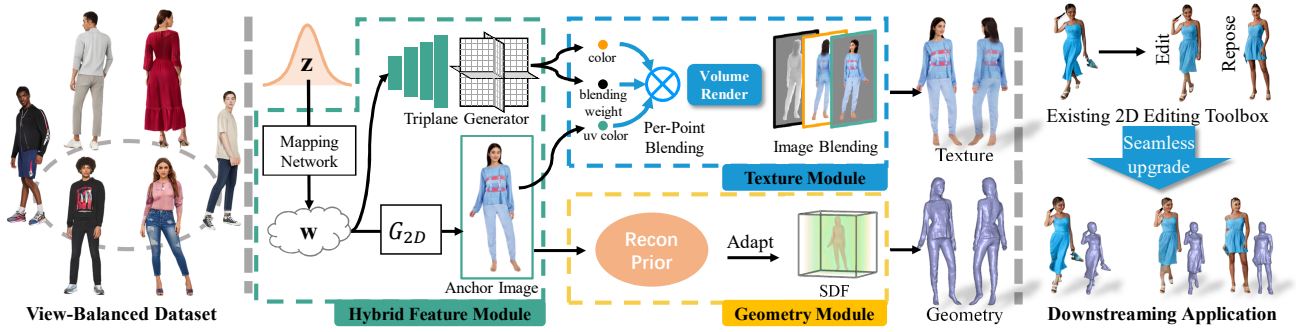
Figure 2. Our approach consists of three modules. The hybrid feature module includes anchor image and tri-plane feature generation. The geometry module includes reconstruction prior and SDF adaptation (Sec. 4.1). Texture module includes sphere tracing based volume rendering, texture and blending weight fields, and two-stage blending (Sec. 4.2)

pose [1, 60], texture [1, 14, 59] or even text [22, 27]. With the recent large-scale dataset SHHQ [12], the unconditional image generators [12, 28–30] achieve synthesizing appealing enough human images that are better than concurrent 3D ones. In the meanwhile, various downstream editing techniques over faces are then extended to full-scale human editing [53, 56, 62, 63, 76]. Rather than modeling a separate latent space for 3D human generator, we leverage the rich human appearance priors provided by 2D human generator by sharing the same latent space of a pre-trained 2D stylegan, and by the way, upgrade the editing toolchain to 3D era.

**Neural Human Modeling.** Rather than classical multi-view stereo methods that require complicated hardware, recent implicit occupancy field-based neural methods [17, 23, 57, 58] achieve reconstructing detailed human geometry with sparse or even one camera. The following works [26, 34, 64] further reveal the effectiveness of the neural occupancy field in the real-time textured dynamic human reconstruction, equipped with a neural texture blending scheme. Embracing the developing of NeRF techniques [8, 9, 39–41, 43, 44, 47, 67, 69, 71, 80, 84], the human shape prior augmented NeRFs achieve modeling realistic human bodies [31, 36, 49, 55, 86], learning animatable avatars [33, 54, 72] and generalizing across different persons [31, 68, 86] from temporal data. However, such techniques can only build human models from actually captured data, i.e., images and videos, and cannot generate novel individuals and appearances. In contrast, we learn a 3D human generator from only 2D human images and largely alleviate the cost of producing high-fidelity virtual humans.

## 3. Overview

By leveraging rich priors from 2D generation and 3D reconstruction models, HumanGen enables delicate 3D human generation with high-quality geometry details and photo-realistic textures. We achieve this by employing a 2D generator to synthesize an "anchor image" with exquisite textures, then lifting it to 3D space and utilizing a 3D reconstructor to enhance geometry details. We briefly discuss each step below, and provide detailed explanation in Sec. 4.

**Hybrid Feature.** As illustrated in Fig. 2, HumanGen first employs a pretrained 2D generator $G_{2D}$ to map gaussian noise $z$ to $w$ latent space and produce an $1024 \times 512$ anchor image. To lift the anchor image to 3D space, we further utilize another 3D generator (EG3D [5]), synthesizing a tri-plane from $w$ to complete the missing information. The tri-plane is composed of three feature planes $F_{xy}$, $F_{xz}$, $F_{yz}$ which align $xy$, $xz$ and $yz$ axes, respectively. We then align the anchor image with $F_{xy}$ to guide the tri-pane to synthesize consistent 3D information as the anchor image.

**Geometry Generation.** Prior works typically learn geometry in low-resolution 3D space [5, 15], which produce over-smoothed human shapes. In contrast, we utilize priors in 3D reconstruction models [58] to enhance more detailed shape generation. Specifically, HumanGen also represents geometry with signed distance field (SDF). We employ pixel-aligned global and local feature as well as occupancy field in PIFuHD [58] as guidance to regress SDF values and can therefore synthesize more detailed human geometry.

**Texture Generation and Blending.** HumanGen utilizes texture fields to decide RGB values at each 3D point and applies volume rendering to synthesize images. However, both reconstruction and generation priors help improve the texture quality. The former allows HumanGen to directly generate high-res images without any superresolution [5, 15]. Specifically, the SDF field adapted from reconstruction prior helps restrict the sampling regions, enabling a more efficient volume rendering to synthesize $512 \times 256$ images directly. We further incorporate the rich details from the anchor image to enhance the quality of rendered textures. We learn a blending weight field and propose a two-stage blending scheme to merge the anchor image and synthesized texture. We achieve view-consistent, photo-realistic texture generation, which we show later.

## 4. Method

### 4.1. Geometry Generation

Given a sampled anchor image, HumanGen first lifts it into 3D geometry. It incorporates reconstruction pri-

ors in the lifting to synthesize fine geometry details. We first discuss how to extract priors from the reconstruction model [58], then describe the way to utilize them in 3D generation.

**Reconstruction Priors.** We choose PIFuHD [58] to provide reconstruction prior, which can faithfully reconstruct fine human geometry with plausible details like hairs or wrinkles from single-view images. Given a 3D point $\mathbf{X} \in \mathbb{R}^3$, PIFuHD projects it on the image and applies a global and a local feature extractor to obtain the corresponding pixel-aligned features $(f^g, f^l)$. It further employs an implicit function $\mathbb{F}_{occ} : (\mathbf{X}, f^g, f^l) \mapsto o$ to map $\mathbf{X}$ to an occupancy $o \in [0, 1]$. We therefore collect all the features $f = (f^g, f^l)$ and the occupancy value $o$ as strong priors and apply them in the geometry generation.

**Prior-guided SDF Adaptation.**

We choose the SDF field to represent geometry, which has revealed better surface modelling [51, 70, 79]. Besides, SDF representation further allows us to derive an efficient sphere tracing to generate high-res images, which we detail in Sec. 4.2. We employ a four-layer MLP $\mathbb{F}_{sdf} : (\mathbf{X}, f) \mapsto s$ which predicts SDF value $s$ given a sample point $\mathbf{X}$ and PIFuHD feature $f$. During training, we sample surface points $\mathbf{X}_s$, where they should be 0.5 for occupancy denoted as: $o(\mathbf{X}_s) = 0.5$. We train $\mathbb{F}_{sdf}$ to predict correct SDF values ($s(\mathbf{X}_s) = 0$) for those surface points, where $\mathcal{L}_{3D\_SDF} = \|s(\mathbf{X}_s)\|_2^2$. In addition to the common Eikonal loss $\mathcal{L}_{eik} = \|\nabla s(X)\|_2 - 1.0$ for regularizing SDF gradients, we also add mask loss to make geometry converge. Given a random view $\mathbf{v}$, we get intersection mask $M_{o,\mathbf{v}}$ by ray marching from the occupancy field and alpha map $M_{s,\mathbf{v}}$ from the predicted SDF field. To get $M_{s,\mathbf{v}}$, we volume render following [51] where density is calculated by: $\sigma(\mathbf{X}) = \alpha^{-1} \text{sigmoid}(-s(\mathbf{X})/\alpha)$ and $\alpha$ here is a learnable parameter .The mask loss is to minimize their difference, where $\mathcal{L}_{mask,\mathbf{v}} = \|M_{o,\mathbf{v}} - M_{s,\mathbf{v}}\|_2^2$. The final geometry loss is defined as

$$\mathcal{L}_{geo} = \lambda_{mask}\mathcal{L}_{mask,v} + \lambda_{3D\_SDF}\mathcal{L}_{3D\_SDF} + \lambda_{eik}\mathcal{L}_{eik}. \tag{1}$$

Following PIFu [57], we further train a color module that will be further used in later texture generation. We first employ a CNN to extract image features. For a 3D point $\mathbf{X}$, we project it on the feature map and collect color feature $f^c$. We then train the color module $\mathbb{F}_{col.} : (f^c, f) \mapsto c \in [0, 1]^3$, and supervise the loss between predicted color and GT sampled color from textured meshes.

## 4.2. Texture Generation and Blending.

**Volume Render with Sphere Tracing.** As illustrated in Fig.4, geometry branch takes anchor image from $G_{2D}$ and generates the corresponding SDF field, which allows us to perform sphere tracing to find the surface where $s(\mathbf{X}) = 0$. Specifically, we choose an orthogonal camera model fitting
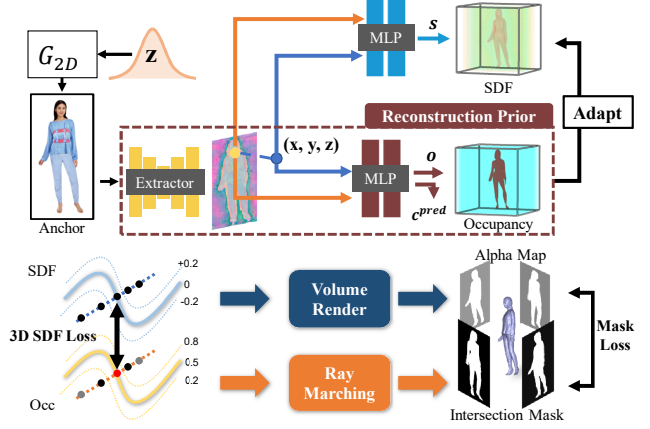


Figure 3. Illustration of reconstruction prior and our SDF adaptation scheme.(Sec. 4.1)

with the PIFuHD setting to render. For each ray, we first query $s(\mathbf{X}_0)$ at $\mathbf{X}_0 = \mathbf{o} + t_{start}\mathbf{d}$ where $\mathbf{o}$ and $\mathbf{d}$ is ray origin and direction, $t_{start}$ is the pre-defined starting step. Then it iterates to query $s(\mathbf{X}_n)$ at $\mathbf{X}_n = \mathbf{o} + (t_{start} + \sum_{i=0}^{n-1} s(\mathbf{X}_i))\mathbf{d}$, until $s(\mathbf{X}_n)$ converges to 0 or iteration exceeds $n_{max}$ times. We set it to be 12 empirically. With the intersection point of each ray with geometry, we only sample 6 points uniformly around the intersection point to efficiently apply volume rendering to synthesize images. The minimal sampling number allows us to generate $512 \times 256$ high-res images.

**Texture and Blending Weight Field.** With tri-plane from hybrid feature generation, for any queried point $\mathbf{X}$ in 3D space, it is projected onto each feature plane to get feature $f_{xy}$, $f_{xz}$ and $f_{yz}$. We model a two-layer MLP following [5] as an implicit decoder to decode color $c \in [0, 1]^3$ and blending weight $b \in [0, 1]$. To make texture generation branch geometry-aware, we also apply the PIFuHD feature $f$ from Sec.4.1: $f_{decoder} : (E(f_{xy}, f_{xz}, f_{yz}), f) \mapsto (c, b)$. $E$ denotes the mean operation. For convenience, we use $c(\mathbf{X})$ and $b(\mathbf{X})$ to denote color and blending weight at $\mathbf{X}$.

**Two-Stage Blending.** Our goal is to synthesize a high-detailed texture map. While texture field tends to produce under-detailed results, we further blend it with anchor image to enhance details. Before volume rendering, each sample point on rays will query its SDF $s$, RGB $c$ and blending weight $b$. To incorporate information from anchor image, sample points are also projected to anchor image to fetch pixel-aligned RGB $c_{uv}(\mathbf{X})$. $c_{uv}(\mathbf{X})$ and $c(\mathbf{X})$ are then blended through $c_b(\mathbf{X}) = c(\mathbf{X}) \cdot (1 - b(\mathbf{X})) + c_{uv}(\mathbf{X}) \cdot b(\mathbf{X})$.

Then we render $512 \times 256$ raw RGB image $I_{raw}$ and blending map $I_{map}$ with volume rendering. For each ray, we get integral color $C(\mathbf{r})$ and blending weight $B(\mathbf{r})$ of ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ by sampling 6 sample points in an interval surrounding its intersection point with geometry. The sampling interval $[t_{start}, t_{end}]$ is empirically determined by
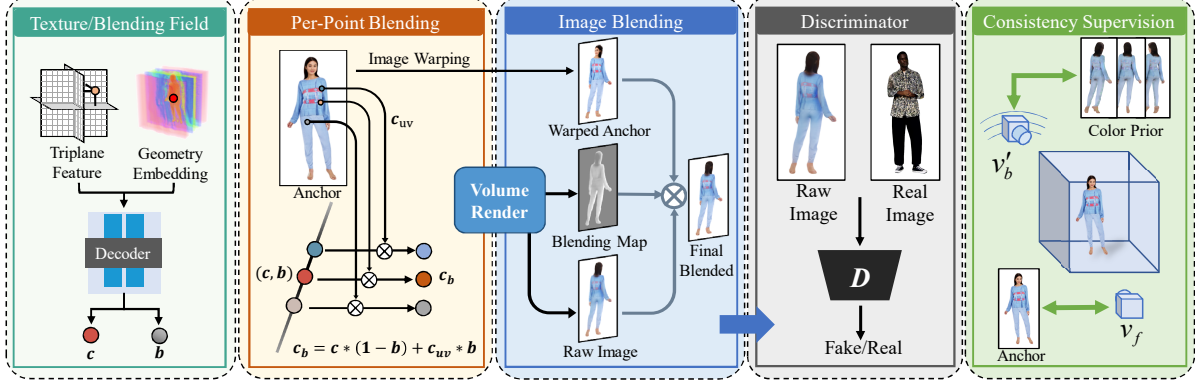
4

Figure 4. Illustration of texture and blending weight fields, two-stage blending and consistency supervision.(Sec. 4.2)

convergence of $\alpha$ in Sec.4.1

$$C(\mathbf{r}) = \int_{t_{st}}^{t_{ed}} T(t)\sigma(\mathbf{r}(t))c_b(\mathbf{r}(t))dt,$$

$$B(\mathbf{r}) = \int_{t_{st}}^{t_{ed}} T(t)\sigma(\mathbf{r}(t))b(\mathbf{r}(t))dt, \quad (2)$$

$$T(t) = \exp\left(-\int_{t_{st}}^{t} \sigma(\mathbf{r}(s))ds\right).$$

Since points with same $xy$ coordinates blend same color from 2D image, from side view, the raw rendering result tends to have "stretching" artifact. To alleviate such artifact, we further utilize calculated blending map to blend rendered image $I_v$ with warped anchor image as post-processing: $I_v = I_{raw} \cdot (1 - I_{map}) + \hat{I}_{anchor} \cdot I_{map}$ where $I_v$ denotes final image from specific rendering view $v$, $\hat{I}_{anchor}$ denotes warped anchor image from frontal view to render view.

**Consistency Supervision.** To enhance 3D consistency with the 2D anchor image, we design two consistency losses $\mathcal{L}_{front}^{CS}$ and $\mathcal{L}_{back}^{CS}$. Firstly, the anchor image $I$ from $G_{2D}$ is aligned with $xy$ axes, so it can naturally be used to supervise rendered image $I_{v_f}$ from frontal camera $v_f$ viewing through z-axis with 2D photometric loss $\mathcal{L}_{2D}$. Besides, from the frontal camera $v_f$, we calculate the ray intersection points $\mathbf{X}_{inter}$ with sphere tracing. The decoded color $c(\mathbf{X}_{inter})$ is supervised to be close to their pixel-aligned RGB color $c_{uv}(\mathbf{X}_{inter})$ from anchor image by 3D RGB loss. Finally, we add a regularization term $\mathcal{L}_{reg}$ on the blending weight of $\mathbf{X}_{inter}$ to enhance the blending effect on the frontal side. The final loss is as $\mathcal{L}_{front}^{CS} = \lambda_{2D\_front}\mathcal{L}_{2D,v_f} + \lambda_{3D\_RGB}\mathcal{L}_{3D\_RGB} + \lambda_b\mathcal{L}_{reg,b}$, where

$$\mathcal{L}_{2D,v_f} = \|I_{v_f} - I\|_2^2,$$

$$\mathcal{L}_{3D\_RGB} = \|c(\mathbf{X}_{inter}) - c_{uv}(\mathbf{X}_{inter})\|_2^2, \quad (3)$$

$$\mathcal{L}_{reg,b} = \|1 - b(\mathbf{X}_{inter})\|_2^2.$$

To supervise consistency on the back view, we use the pre-trained reconstruction color prior as mentioned in Sec. 4.1 to calculate consistency loss $\mathcal{L}_{back}^{CS}$. For rendering view $v_b$

opposite to $v_f$, a standard gaussian noise is added to its spherical coordinates to get random view $v_b'$ and 2D photometric loss is calculated between volume rendered image $I_{v_b'}$ and predicted color image $I_{v_b'}^{pred}$ which is calculated by query $f_c(\mathbf{X}_{inter})$ on intersection points. $\mathcal{L}_{back}^{CS}$ is as:

$$\mathcal{L}_{back}^{CS} = \lambda_{2D\_back}\mathcal{L}_{2D,v_b'},$$

$$\mathcal{L}_{2D,v_b'} = \|I_{v_b'} - I_{v_b'}^{pred}\|_2^2. \quad (4)$$

### 4.3. Training.

**Traning Set.** Current 2D human image collections typically have view distribution bias [12,81], as they focus more on taking photos in front views. Previous method [83] use sampling trick to alleviate this. We notice that $G_{2D}$ also generates view-biased images and produces more diverse side view and back view at a lower frequency. So we apply a human pose calibration [82] on $G_{2D}$ to filter images and generate a relatively view-balanced training set of 230k images.

**GAN Training.** We further use adversarial loss to refine whole texture generation. For previous triplane-based methods, they all assume that the learned object is in a canonical state and thus condition discriminator on live state parameters like camera pose and human skeleton pose. In our setting, we assume that images from $G_{2D}$ are aligned with $xy$ axes. However, human in image is probably not facing to front, so we condition discriminator on relative pose of human. In order to render with same view distribution with dataset, for image $I$ given by $G_{2D}$, we calibrate it with PyMAF [82] to get relative human root transformation $M_{human}$ and skeleton pose. For sampled camera view $v$ from dataset, it is transformed with $v_r = M_{human}^{-1}v$, where $v_r$ denotes actual rendering view. The GAN loss is

$$\mathcal{L}_{adv}(\theta_G, \theta_D) = \mathbf{E}_{z\sim p_z, v\sim p_{dst}}[\log(D(G(z,v;\theta_G)))]$$
$$+\mathbf{E}_{I_r\sim p_{real}}[\log(-D(I_r;\theta_D)) + \lambda\|\nabla D(I_r;\theta_D)\|_2^2]. \quad (5)$$

HumanGen training is separated into first stage training for geometry branch and second stage training for texture branch. More details can be found in the supplementary.

5

Figure 5. The geometry and texture generation results of our HumanGen on various identities.

# 5. Experimental Results

In this section, we evaluate our HumanGen on quality of generated texture and geometry. Metrics used are Frechet Inception Distance (FID) [20] and Kernel Inception Distance (KID) [4]. 50k generated images are used to compute scores. However, FID, KID cannot properly assess geometry quality, so we further evaluate geometry by comparing generated depth with aligned depth predicted by MiDaS [32] in masked region on 5k generated samples.

## 5.1. Comparison

We compare our method with state-of-the-art 3D-aware generation methods, EG3D [5], StyleSDF [51]. Another compared baseline is the combination of Stylegan-human [12] and PIFu [57]. Besides, we compare GNARF [3] implemented by ourselves with super-resolution module to achieve same resolution. As illustrated in Fig. 6, StyleSDF fails to generate full-body and diverse texture and performs poor on generated geometry. EG3D generates wrong geometry and its texture appearance lacks view consistency, especially on the head region. GNARF generates geometry with artifacts on body parts because of self-intersection problem caused by imperfect pose calibration from single-view human image. 2D Generator with PIFu generates correct full-body geometry but lacks fine geometry details and its texture tends to be blurred. While our HumanGen achieves better detailed full-body geometry as well as photo-realistic texture generation. We compare above methods on FID and KID for texture and depth for geometry. The quantitative results in Tab. 1 demonstrate that our method achieves the best FID and depth.

## 5.2. Ablation Study

**Adversarial Training.** We evaluate texture generation based on the same fixed geometry branch. In Fig. 7(b), we train a tri-plane generator to decode texture only with ad-
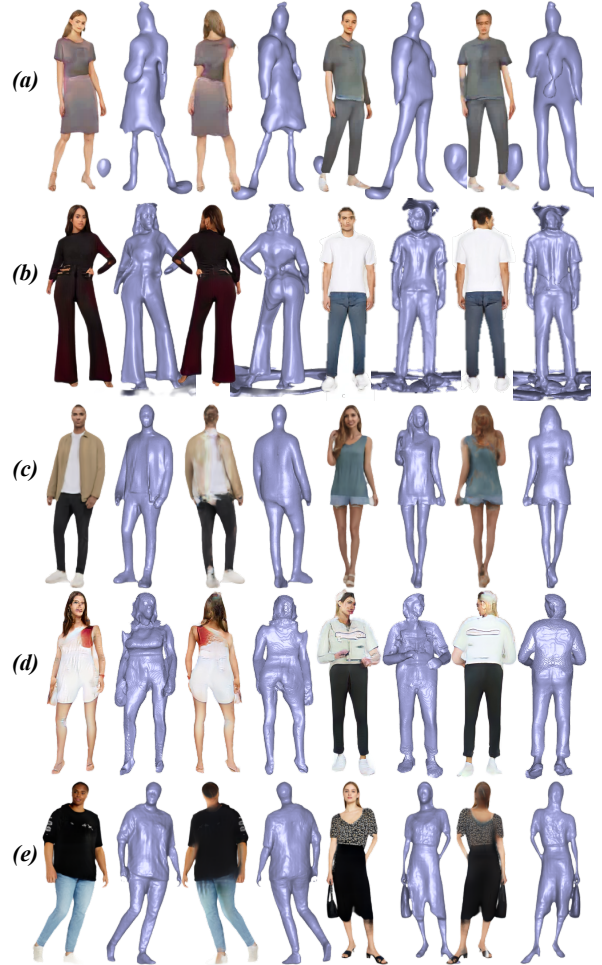


Figure 6. Qualitative comparison. (a) StyleSDF (b) EG3D (c) Stylegan-human + PIFu (d) GNARF (e) Ours

versarial loss (Eqn. 5) and discriminator without condition (**base**). We find it hard for the generator to maintain color or identity consistency with the anchor image and generate high-quality texture of human part. In Fig. 7(d), with-

Table 1. Quantitative comparison of generation results.

| Method | FID↓ | KID↓ | Depth↓ |
|--------|------|------|--------|
| EG3D [5] | 21.33 | **0.0110** | 0.0395 |
| StyleSDF [51] | 36.69 | 0.0309 | 0.0493 |
| 2D-G [12]+PIFu [57] | 39.20 | 0.0351 | 0.0379 |
| GNARF [3] | 24.61 | 0.0169 | 0.0408 |
| Ours | **20.97** | 0.0157 | **0.0201** |

Table 2. Quantitative evaluation of texturing generation.

| Method | FID↓ | KID↓ |
|--------|------|------|
| base | 31.07 | 0.0251 |
| w/o both CS | 34.52 | 0.0246 |
| w/o frontal CS | 25.75 | 0.0193 |
| w/o back CS | 21.69 | 0.0137 |
| w/o both blending scheme | 48.37 | 0.0436 |
| w/o image blending | 25.94 | 0.0193 |
| full (view-biased dataset) | 21.81 | 0.0161 |
| full | **20.97** | 0.0157 |

Table 3. Quantitative evaluation of geometry generation.

| Method | Depth Diff ↓ |
|--------|--------------|
| w/o PIFuHD feat. | 0.0383 |
| w/o SDF | 0.0183 |
| Full | **0.0145** |

out relative camera pose conditioned on discriminator, the generator tends to be confused about correct body part position. While our method maintains high consistency with given anchor image and learns to generate full-body texture. Quantitative result in Tab. 2 demonstrates that our method achieves better score.

**Consistency Supervision.** Let CS denote consistency supervision described in Sec. 4.2. In Fig. 8 (a), **without both CS** from Eqn. 3 and Eqn. 4, the generator is prone to change cloth type or color and is not consistent with anchor image. In (c) and (e), **without frontal CS** or **without back CS**, the texture generation on unsupervised side is prone to disobey color or identity consistency with supervised side. Our full method maintains consistency with anchor image while achieving the best self-consistency. As in Tab. 2, our full method achieves better score.

**Two-Stage Blending.** In Fig. 9 (a), **without both blending scheme**, the texture generation depends completely on decoded RGB, which is prone to be blurred and cannot recover high-fidelity texture consistent with anchor image. The second-stage image blending depends on learned first-stage per-sample-point blending. As illustrated in (c), **without image blending** as second-stage post-processing, the texture will have "stretching" artifacts because sample points sharing the same $xy$-coordinates are blended with the same pixel color from anchor image. In our full method, per-sample-point blending enables recovering high-frequency details and image blending alleviates artifacts of first-stage blending. Quantitative results can refer to Tab. 2.

**View Distribution of Training Set.** As for evaluation of view distribution of training set, we further train our full method on another synthesised view-biased training set without data from back views. As shown in Fig. 10, without adversarial loss on back $180°$ region, the texture generation on back tends to degrade to predicted color prior from reconstruction model which is prone to be blurred and has obvious color difference with frontal texture. Quantitative results can refer to Tab. 2.

**Analysis on Geometry Adaptation.** We conduct an analysis on geometry adaptation with reconstruction prior. In Fig. 11 (b), without PIFuHD [58] feature, we directly adapt a tri-plane generator with losses from Eqn. 1, but the geometry fails to converge. In (c)&(d), we respectively adapt an MLP to output density and SDF. Geometry in SDF representation is more smooth and has clear surface level set.
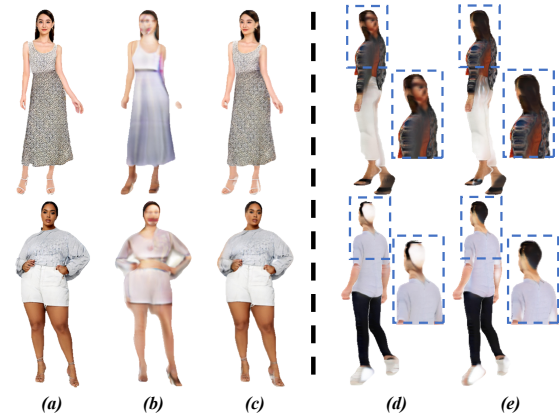


Figure 7. Qualitative evaluation of texture generation. (a) Anchor image. (b) base; (c)&(e) Ours (d) w/o relative cam pose condition.



Figure 8. Qualitative evaluation of texture generation. (a) w/o both CS. (b)&(d)&(f) Ours ; (c) w/o frontal CS; (e) w/o back CS.

As shown in Tab. 3, with better representation of SDF, our full method achieves the lowest difference with depth from PIFuHD.

**Application.** Our method is naturally compatible with existing 2D editing methods. As shown in Fig. 12(a), with style mixing on 2D, we are able to achieve 3D results with given pose. In Fig. 12(b), real images can be inverted to latent space to generate 3D results. In Fig. 12(c), by editing in
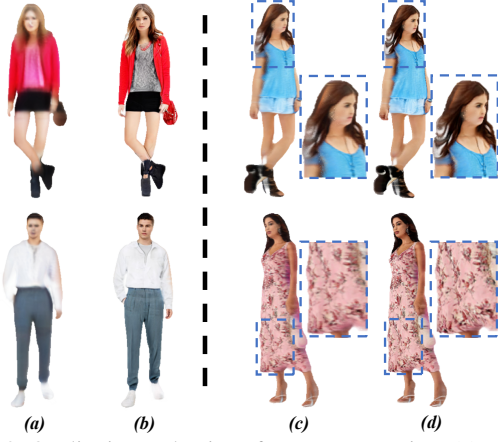
Figure 9. Qualitative evaluation of texture generation. (a) w/o two-stage blending scheme. (b)&(d) Ours; (c) w/o image blending.
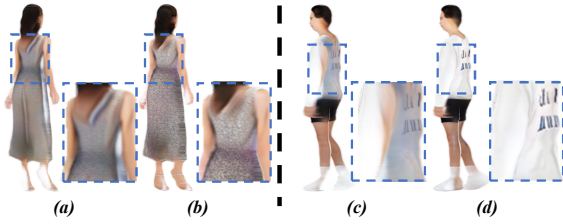


Figure 10. Qualitative evaluation of view-distribution of training set. (a)&(c) results on view-biased training set; (b)&(d) Ours.
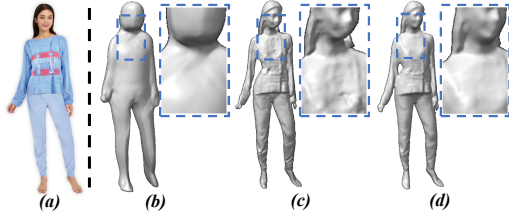


Figure 11. Qualitative evaluation of human geometry generation. (a) Anchor images. (b) w/o PIFuHD feature; (c) w/o SDF; (d) Full.

latent space, the length change of upper or lower clothes can be seamlessly upgraded to 3D. In Fig. 12(d), 2D text-guided generation can be extended to 3D.

### 5.3. Limitation and Discussion

Although HumanGen achieves generating human with detailed geometry and 360° realistic rendering, it still has some limitations. First, the adopted view-balanced training set is essentially generated from a view-biased 2D generator, so the bias problem still potentially exists on uneven generation quality among different poses and views. It is still meaningful to have a real dataset with balanced view distribution and data of same identities. Second, with rich explicit priors in our method, it limits the out-of-domain generation ability of our method. Furthermore, our method does not support continuous skeletal pose control. It's promising to include motion prior into current framework for high-quality deformable generation, but requiring much more diverse training data.



(a) Pose editing with 2D style mixing

(b) Real image Inverse

(c) Editing with latent space

A lady is wearing a blue blouse and black pants. She has black hair.
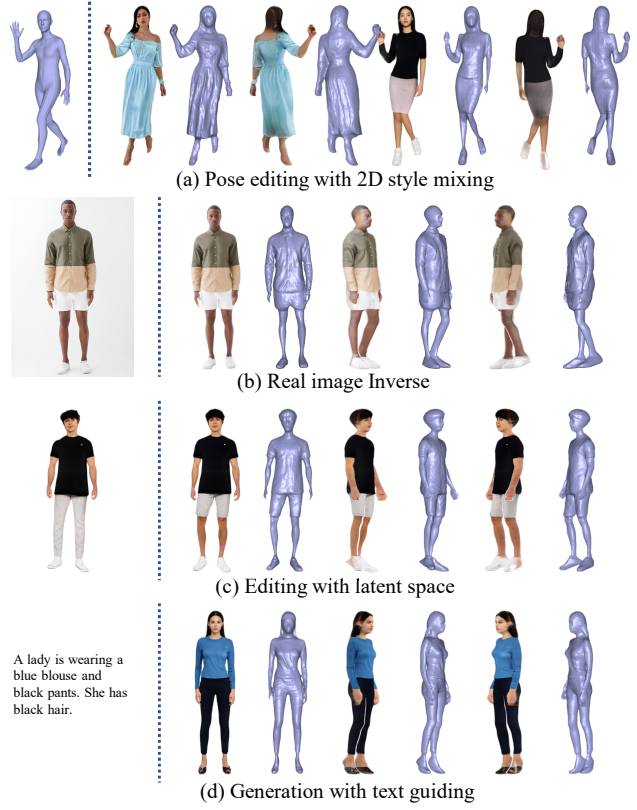
(d) Generation with text guiding

Figure 12. Various 3D applications of HumanGen.

Though human eyes can still distinguish results of HumanGen, we care much about the ethical issue behind it. HumanGen should not be utilized to create any fake result of the real person and deceive people unfamiliar with this domain. Meanwhile, all our results are carefully chosen to ensure impartiality.

### 6. Conclusion

We have presented a novel 3D human generation scheme with detailed geometry and 360° realistic free-view rendering. Our key idea is to introduce the concept of "anchor image" to aid the human generation using various human priors explicitly. Our hybrid feature representation efficiently bridges the latent space of HumanGen with the existing 2D generator. Our geometry adapting scheme enables fine-grained details synthesis from 3D human reconstruction prior, while our two-stage blending scheme further encodes the rich texture information in the anchor image for appearance generation. Our experimental results demonstrate the effectiveness of HumanGen for state-of-the-art 3D human generation. Various 3D applications of HumanGen further demonstrate its compatibility to existing off-the-shelf 2D editing toolbox based on latent disentanglement. With the above unique ability, we believe that our approach is a critical step for high-quality 3D human generation, with various potential applications in VR/AR.

# References

[1] Badour Albahar, Jingwan Lu, Jimei Yang, Zhixin Shu, Eli Shechtman, and Jia-Bin Huang. Pose with style: Detail-preserving pose-guided image synthesis with conditional stylegan. *ACM Transactions on Graphics (TOG)*, 40(6):1–11, 2021. 3

[2] Thiemo Alldieck, Mihai Zanfir, and Cristian Sminchisescu. Photorealistic monocular 3d reconstruction of humans wearing clothing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1506–1515, 2022. 2

[3] Alexander W. Bergman, Petr Kellnhofer, Wang Yifan, Eric R. Chan, David B. Lindell, and Gordon Wetzstein. Generative neural articulated radiance fields. In *NeurIPS*, 2022. 1, 2, 6, 7

[4] Mikolaj Binkowski, Danica J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans, 2018. 6

[5] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. 1, 2, 3, 4, 6, 7

[6] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5799–5809, 2021. 1, 2

[7] Anpei Chen, Ruiyang Liu, Ling Xie, Zhang Chen, Hao Su, and Jingyi Yu. Sofgan: A portrait image generator with dynamic styling. *ACM Transactions on Graphics (TOG)*, 41(1):1–26, 2022. 2

[8] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision (ECCV)*, 2022. 3

[9] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–14133, 2021. 3

[10] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. Gram: Generative radiance manifolds for 3d-aware image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10673–10683, 2022. 1, 2

[11] Patrick Esser, Ekaterina Sutter, and Björn Ommer. A variational u-net for conditional appearance and shape generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8857–8866, 2018. 2

[12] Jianglin Fu, Shikai Li, Yuming Jiang, Kwan-Yee Lin, Chen Qian, Chen Change Loy, Wayne Wu, and Ziwei Liu. Stylegan-human: A data-centric odyssey of human generation. In *European Conference on Computer Vision*, pages 1–19. Springer, 2022. 1, 2, 3, 5, 6, 7, 13

[13] Matheus Gadelha, Subhransu Maji, and Rui Wang. 3d shape induction from 2d views of multiple objects. In *2017 International Conference on 3D Vision (3DV)*, pages 402–411. IEEE, 2017. 2

[14] Artur Grigorev, Karim Iskakov, Anastasia Ianina, Renat Bashirov, Ilya Zakharkin, Alexander Vakhitov, and Victor Lempitsky. Stylepeople: A generative model of fullbody human avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5151–5160, 2021. 3

[15] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d aware generator for high-resolution image synthesis. In *International Conference on Learning Representations*, 2022. 1, 2, 3

[16] Zekun Hao, Arun Mallya, Serge Belongie, and Ming-Yu Liu. Gancraft: Unsupervised 3d neural rendering of minecraft worlds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14072–14082, 2021. 2

[17] Tong He, Yuanlu Xu, Shunsuke Saito, Stefano Soatto, and Tony Tung. Arch++: Animation-ready clothed human reconstruction revisited. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11046–11056, 2021. 3

[18] Philipp Henzler, Niloy J. Mitra, and Tobias Ritschel. Escaping plato's cave: 3d shape from adversarial rendering. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 2

[19] Philipp Henzler, Niloy J Mitra, and Tobias Ritschel. Escaping plato's cave: 3d shape from adversarial rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9984–9993, 2019. 2

[20] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Günter Klambauer, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a nash equilibrium. *CoRR*, abs/1706.08500, 2017. 6

[21] Fangzhou Hong, Zhaoxi Chen, Yushi Lan, Liang Pan, and Ziwei Liu. Eva3d: Compositional 3d human generation from 2d image collections. *arXiv preprint arXiv:2210.04888*, 2022. 1, 2

[22] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu. Avatarclip: Zero-shot text-driven generation and animation of 3d avatars. *ACM Transactions on Graphics (TOG)*, 41(4):1–19, 2022. 2, 3

[23] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. Arch: Animatable reconstruction of clothed humans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3093–3102, 2020. 3

[24] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. *CVPR*, 2017. 2

[25] Ajay Jain, Ben Mildenhall, Jonathan T. Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided object generation with dream fields. *CVPR*, 2022. 2

[26] Yuheng Jiang, Suyi Jiang, Guoxing Sun, Zhuo Su, Kaiwen Guo, Minye Wu, Jingyi Yu, and Lan Xu. Neuralhofusion: Neural volumetric rendering under human-object interactions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6155–6165, 2022. 3

[27] Yuming Jiang, Shuai Yang, Haonan Qju, Wayne Wu, Chen Change Loy, and Ziwei Liu. Text2human: Text-driven controllable human image generation. *ACM Transactions on Graphics (TOG)*, 41(4):1–11, 2022. 2, 3

[28] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*, 34:852–863, 2021. 1, 3

[29] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 1, 3

[30] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 1, 2, 3, 13

[31] Youngjoong Kwon, Dahun Kim, Duygu Ceylan, and Henry Fuchs. Neural human performer: Learning generalizable radiance fields for human performance rendering. *Advances in Neural Information Processing Systems*, 34:24741–24752, 2021. 3

[32] Katrin Lasinger, René Ranftl, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *CoRR*, abs/1907.01341, 2019. 6

[33] Ruilong Li, Julian Tanke, Minh Vo, Michael Zollhofer, Jurgen Gall, Angjoo Kanazawa, and Christoph Lassner. Tava: Template-free animatable volumetric actors. In *European Conference on Computer Vision (ECCV)*, 2022. 3

[34] Ruilong Li, Yuliang Xiu, Shunsuke Saito, Zeng Huang, Kyle Olszewski, and Hao Li. Monocular real-time volumetric performance capture. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, pages 49–67, Cham, 2020. Springer International Publishing. 3

[35] Yiyi Liao, Katja Schwarz, Lars Mescheder, and Andreas Geiger. Towards unsupervised learning of generative models for 3d controllable image synthesis. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[36] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. *ACM Trans. Graph.(ACM SIGGRAPH Asia)*, 2021. 3

[37] Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1096–1104, 2016. 1

[38] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Smpl: A skinned multi-person linear model. *ACM Trans. Graph.*, 34(6):248:1–248:16, Oct. 2015. 1, 2

[39] H. Luo, A. Chen, Q. Zhang, B. Pang, M. Wu, L. Xu, and J. Yu. Convolutional neural opacity radiance fields. In *2021 IEEE International Conference on Computational Photog-*

[40] Haimin Luo, Teng Xu, Yuheng Jiang, Chenglin Zhou, Qiwei Qiu, Yingliang Zhang, Wei Yang, Lan Xu, and Jingyi Yu. Artemis: Articulated neural pets with appearance and motion synthesis. *ACM Trans. Graph.*, 41(4), jul 2022. 3

[41] Quan Meng, Anpei Chen, Haimin Luo, Minye Wu, Hao Su, Lan Xu, Xuming He, and Jingyi Yu. Gnerf: Gan-based neural radiance field without posed camera. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6351–6361, 2021. 3

[42] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, pages 405–421, Cham, 2020. Springer International Publishing. 1, 2

[43] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020. 3

[44] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022. 3

[45] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. Hologan: Unsupervised learning of 3d representations from natural images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7588–7597, 2019. 2

[46] Thu H Nguyen-Phuoc, Christian Richardt, Long Mai, Yongliang Yang, and Niloy Mitra. Blockgan: Learning 3d object-aware scene representations from unlabelled images. *Advances in Neural Information Processing Systems*, 33:6767–6778, 2020. 2

[47] Michael Niemeyer, Jonathan T. Barron, Ben Mildenhall, Mehdi S. M. Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3

[48] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[49] Atsuhiro Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. Neural articulated radiance field. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5762–5772, 2021. 3

[50] Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo Strauss, and Andreas Geiger. Texture fields: Learning texture representations in function space. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4531–4540, 2019. 2

[51] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman.

Stylesdf: High-resolution 3d-consistent image and geometry generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13503–13513, 2022. 1, 2, 4, 6, 7

[52] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2

[53] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2085–2094, 2021. 1, 3

[54] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14314–14323, 2021. 3

[55] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *CVPR*, 2021. 3

[56] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *ACM Transactions on Graphics (TOG)*, 42(1):1–13, 2022. 1, 3

[57] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 3, 4, 6, 7

[58] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2, 3, 4, 7

[59] Kripasindhu Sarkar, Vladislav Golyanik, Lingjie Liu, and Christian Theobalt. Style and pose control for image synthesis of humans from a single monocular view. *arXiv preprint arXiv:2102.11263*, 2021. 3

[60] Kripasindhu Sarkar, Lingjie Liu, Vladislav Golyanik, and Christian Theobalt. Humangan: A generative model of human images. In *2021 International Conference on 3D Vision (3DV)*, pages 258–267. IEEE, 2021. 3

[61] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33:20154–20166, 2020. 1, 2

[62] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. Interfacegan: Interpreting the disentangled face representation learned by gans. *IEEE transactions on pattern analysis and machine intelligence*, 2020. 3

[63] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1532–1540, 2021. 1, 3

[64] Xin Suo, Yuheng Jiang, Pei Lin, Yingliang Zhang, Minye Wu, Kaiwen Guo, and Lan Xu. Neuralhumanfvv: Real-time neural volumetric human performance rendering using rgb cameras. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6226–6237, 2021. 3

[65] Attila Szabó, Givi Meishvili, and Paolo Favaro. Unsupervised generative 3d shape learning from natural images. *arXiv preprint arXiv:1910.00287*, 2019. 2

[66] Shuhei Tsuchida, Satoru Fukayama, Masahiro Hamasaki, and Masataka Goto. Aist dance video database: Multi-genre, multi-dancer, and multi-camera database for dance information processing. In *ISMIR*, volume 1, page 6, 2019. 1

[67] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. Ref-nerf: Structured view-dependent appearance for neural radiance fields. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5481–5490. IEEE, 2022. 3

[68] Liao Wang, Ziyu Wang, Pei Lin, Yuheng Jiang, Xin Suo, Minye Wu, Lan Xu, and Jingyi Yu. ibutter: Neural interactive bullet time generator for human free-viewpoint rendering. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 4641–4650, 2021. 3

[69] Liao Wang, Jiakai Zhang, Xinhang Liu, Fuqiang Zhao, Yanshun Zhang, Yingliang Zhang, Minye Wu, Jingyi Yu, and Lan Xu. Fourier plenoctrees for dynamic radiance field rendering in real-time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13524–13534, 2022. 3

[70] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 2021. 4

[71] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *CVPR*, 2021. 3

[72] Shaofei Wang, Katja Schwarz, Andreas Geiger, and Siyu Tang. Arah: Animatable volume rendering of articulated human sdfs. In *European Conference on Computer Vision*, 2022. 3

[73] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 2

[74] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 2

[75] Ziyu Wang, Yu Deng, Jiaolong Yang, Jingyi Yu, and Xin Tong. Generative deformable radiance fields for disentangled image synthesis of topology-varying objects. *Computer Graphics Forum*, 2022. 2

[76] Zongze Wu, Dani Lischinski, and Eli Shechtman. Stylespace analysis: Disentangled controls for stylegan image genera-

tion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12863–12872, 2021. 1, 3

[77] Jianfeng Xiang, Jiaolong Yang, Yu Deng, and Xin Tong. Gram-hd: 3d-consistent image generation at high resolution with generative radiance manifolds. *arXiv preprint arXiv:2206.07255*, 2022. 2

[78] Yang Xue, Yuheng Li, Krishna Kumar Singh, and Yong Jae Lee. Giraffe hd: A high-resolution 3d-aware generative model. In *CVPR*, 2022. 2

[79] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *CoRR*, abs/2106.12052, 2021. 4

[80] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. PlenOctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 3

[81] Polina Zablotskaia, Aliaksandr Siarohin, Bo Zhao, and Leonid Sigal. Dwnet: Dense warp-based network for pose-guided human video generation. *CoRR*, abs/1910.09139, 2019. 5

[82] Hongwen Zhang, Yating Tian, Xinchi Zhou, Wanli Ouyang, Yebin Liu, Limin Wang, and Zhenan Sun. 3d human pose and shape regression with pyramidal mesh alignment feedback loop. *CoRR*, abs/2103.16507, 2021. 5

[83] Jianfeng Zhang, Zihang Jiang, Dingdong Yang, Hongyi Xu, Yichun Shi, Guoxian Song, Zhongcong Xu, Xinchao Wang, and Jiashi Feng. Avatargen: a 3d generative model for animatable human avatars. *arXiv preprint arXiv:2208.00561*, 2022. 1, 2, 5

[84] Jiakai Zhang, Xinhang Liu, Xinyi Ye, Fuqiang Zhao, Yanshun Zhang, Minye Wu, Yingliang Zhang, Lan Xu, and Jingyi Yu. Editable free-viewpoint video using a layered neural representation. *ACM Trans. Graph.*, 40(4), July 2021. 3

[85] Jichao Zhang, Enver Sangineto, Hao Tang, Aliaksandr Siarohin, Zhun Zhong, Nicu Sebe, and Wei Wang. 3d-aware semantic-guided generative model for human synthesis. In *European Conference on Computer Vision*, pages 339–356. Springer, 2022. 2

[86] Fuqiang Zhao, Wei Yang, Jiakai Zhang, Pei Lin, Yingliang Zhang, Jingyi Yu, and Lan Xu. Humannerf: Efficiently generated human radiance field from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7743–7753, 2022. 3

[87] Xiaoming Zhao, Fangchang Ma, David Güera, Zhile Ren, Alexander G Schwing, and Alex Colburn. Generative multi-plane images: Making a 2d gan 3d-aware. In *European Conference on Computer Vision*, pages 18–35. Springer, 2022. 2

[88] Peng Zhou, Lingxi Xie, Bingbing Ni, and Qi Tian. Cips-3d: A 3d-aware generator of gans based on conditionally-independent pixel synthesis. *arXiv preprint arXiv:2110.09788*, 2021. 1, 2

[89] Peihao Zhu, Rameen Abdal, Yipeng Qin, and Peter Wonka. Sean: Image synthesis with semantic region-adaptive normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5104–5113, 2020. 2

# 7. Supplemental Material

## 7.1. Implementation Details

We train our HumanGen using 4 NVIDIA A40 GPUs. We first train the geometry branch, which takes about 12 hours to converge. The trained implicit function $\mathbb{F}_{sdf}$ has hidden neurons of (273, 128, 32, 1). The loss weights are $\lambda_{mask} = 1$, $\lambda_{3D\_SDF} = 1$ and $\lambda_{eik} = 0.1$. For the texture branch, we first train the base model for around 24 hours and then continue to train each model for another roughly 18 hours. For generating the tri-plane feature, we generate planes of shape $256 \times 256$ with channel size of 32. We further add another two stylegan [30] synthesis blocks with upscale equals to 1 so as to extend the layer number of latent $w^+$ to 18, which is compatible with the layer number of $w^+$ from Stylegan-human [12] generator. The loss weights are $\lambda_{2D\_front} = 1$, $\lambda_{3D\_RGB} = 8$, $\lambda_b = 1\text{e}{-2}$, $\lambda_{2D\_back} = 8$ and the r1 regularization term of adversarial training (Eqn. 5) has weight $\lambda = 10$.

## 7.2. Additional Evaluation

We further conduct a texture fitting evaluation. Given the already-trained geometry branch, we only use anchor image and photometric loss to fit the frontal texture of given geometry. Let $M$ denotes the mapping network from Stylegan-human. As shown in Fig.13, without the mapping network from Stylegan-human (**w/o** $M$), we train a newly-initialized mapping network to map the same noise $z$ to some latent and synthesize the tri-plane feature. However, it generates results that are all blurred and have closing color with each other because the newly-initialized mapping network fails to recover the original latent space of Stylegan-human. Without the geometry embedding feature given to the decoder (**w/o geometry embedding**), though the fitting texture can follow the color consistency with anchor image, it tends to have no details. With both $M$ and geometry embedding (**full**), the fitting results can maintain better consistency with anchor image while recovering some geometry details as well. The corresponding quantitative results are provided in Tab.4. However, the texture fitting results are still less-detailed. Therefore, we further add the blending scheme and GAN training to improve the effects. We provide more generation results of our complete model in Fig.14. Note that our approach can generate photo-realistic results with detailed geometry and free-viewing ability.

Table 4. Quantitative evaluation of texture fitting.

| Method | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| w/o $M$ | 15.86 | 0.7835 | 0.2596 |
| w/o geometry embedding | 19.89 | 0.8236 | 0.2075 |
| full | **21.13** | **0.8434** | **0.1803** |



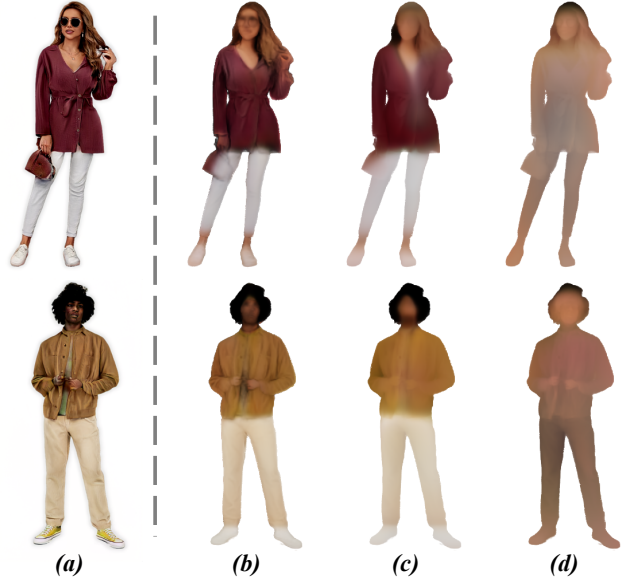*(a)*     *(b)*     *(c)*     *(d)*

Figure 13. Qualitative evaluation of texture fitting. (a) anchor image; (b) full; (c) w/o geometry embedding; (d) w/o $M$.

Figure 14. More generation results using our HumanGen. Note that our approach enables photo-realistic human generation with detailed geometry and free-view rendering.