# Human-VDM: Learning Single-Image 3D Human Gaussian Splatting from Video Diffusion Models

**Zhibin Liu[1], Haoye Dong[2], Aviral Chharia[2], Hefeng Wu[1]**

[1]Sun Yat-Sen University  [2]Carnegie Mellon University

liuzhb26@mail2.sysu.edu.cn, {haoyed, achharia}@andrew.cmu.edu, wuhefeng@mail.sysu.edu.cn
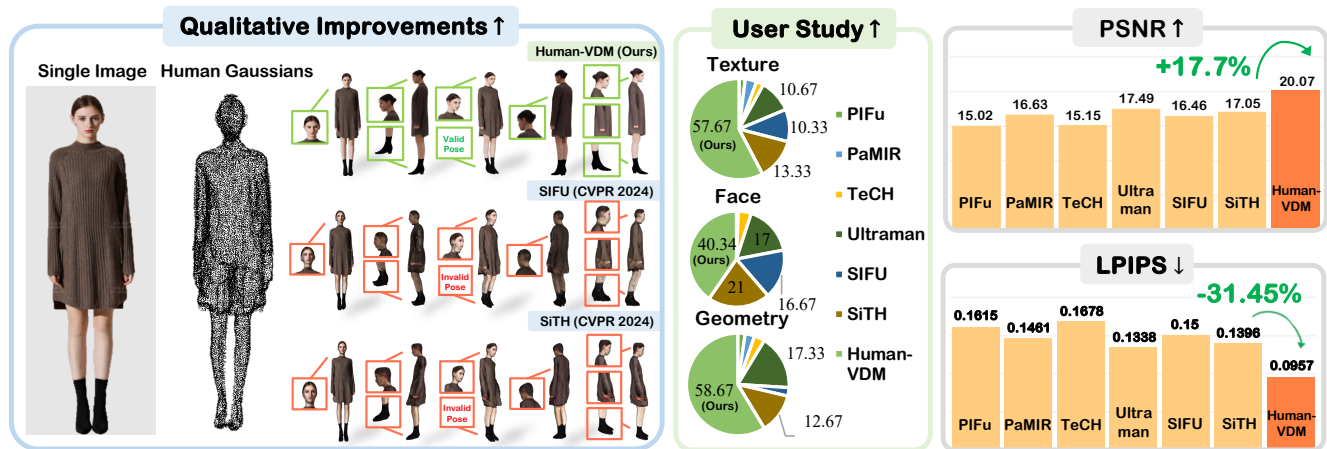
https://Human-VDM.github.io/Human-VDM

Figure 1: **Human-VDM for generating 3D humans from a single image.** Given a single RGB human image, Human-VDM aims to generate high-fidelity 3D human. Human-VDM preserves face identity, delivers realistic texture, ensures accurate geometry, and maintains a valid pose of the generated 3D human, surpassing the current state-of-the-art models.

## Abstract

Generating lifelike 3D humans from a single RGB image remains a challenging task in computer vision, as it requires accurate modeling of geometry, high-quality texture, and plausible unseen parts. Existing methods typically use multi-view diffusion models for 3D generation, but they often face inconsistent view issues, which hinder high-quality 3D human generation. To address this, we propose **Human-VDM**, a novel method for generating 3D **human** from a single RGB image using **V**ideo **D**iffusion **M**odels. Human-VDM provides temporally consistent views for 3D human generation using Gaussian Splatting. It consists of three modules: a view-consistent human video diffusion module, a video augmentation module, and a Gaussian Splatting module. First, a single image is fed into a human video diffusion module to generate a coherent human video. Next, the video augmentation module applies super-resolution and video interpolation to enhance the textures and geometric smoothness of the generated video. Finally, the 3D Human Gaussian Splatting module learns lifelike humans under the guidance of these high-resolution and view-consistent images. Experiments demonstrate that Human-VDM achieves high-quality 3D human from a single image, outperforming state-of-the-art methods in both generation quality and quantity.

## Introduction

Generating 3D humans from a single RGB image has gained significant attention in recent years due to its versatile applications in filmmaking, video games, human-robotic interaction, etc. However, existing approaches for 3D human generation largely rely on multi-view diffusion models, which often suffer from inconsistent views and lead to artifacts. To address this problem, we propose a 3D Human Gaussian Splatting framework that allows users to generate 3D humans from a single 2D image input while ensuring accurate geometry and realistic appearance. However, generating 3D humans using only a single RGB image presents a significant challenge due to its inherent ambiguity, which necessitates inferring unseen geometry and appearance that are not directly captured in a 2D image.

Current approaches address this challenge by incorporating parametric human shape models, such as SCAPE (Anguelov et al. 2005) and SMPL (Loper et al. 2023). However, these methods exclusively focus on reconstructing the human shape, neglecting the appearance details crucial for a fully realistic 3D representation. Earlier works, like PIFu (Saito et al. 2019), attempted to address this gap with a data-driven approach. They used Cy-

cleGAN (Zhu et al. 2017) and residual blocks (Johnson, Alahi, and Fei-Fei 2016) trained on image-3D pairs. However, such methods often struggle with novel appearances or poses mainly due to the lack of sufficient 3D training information. Subsequent methods, such as ECON (Xiu et al. 2023) and 2K2K (Han et al. 2023), enhanced performance by incorporating depth or normal estimation into the generation process. SIFU (Zhang, Yang, and Yang 2024) proposed a 3D human generation method using a side-view based Transformer with 3D aware Refinement. Despite the improvements, these methods often lack detail or result in inaccurate geometry, particularly with high-resolution input images.

Recently, SiTH (Ho et al. 2024) integrated a generative diffusion model into the 3D human generation pipeline to produce realistic textures and geometries, especially in unobserved regions. Ultraman (Chen et al. 2024) introduced a multi-view image generation model that helped in providing essential appearance priors aiding the generation process. Although diffusion models (Rombach et al. 2022), trained on extensive image datasets, have demonstrated potential for creating 3D humans, multi-view diffusion often struggles with generating view-consistent images and tends to introduce artifacts in the generated 3D humans.

This paper proposes Human-VDM, a novel Gaussian Splatting framework for generating 3D humans from a single image using video diffusion models. Human-VDM is comprised of three distinct modules: a view-consistent human video diffusion module, a video augmentation module, and a 3D human Gaussian Splatting module. Human-VDM first generates a 'view-consistent' human video, then enhances the quality of the frames through super-resolution and video frame interpolation, and finally employs 3D Gaussian Splatting (3DGS) (Kerbl et al. 2023) to effectively generate the 3D human model.

Initially, we fine-tune SV3D (Voleti et al. 2024), a latent video diffusion model specifically designed for generating object videos, to enable it to generate view-consistent human videos. However, a direct application of video diffusion models (Voleti et al. 2024) to the 3D human generation can result in geometric artifacts and blurry textures. Additionally, the generated video consists of only 21 frames at a low resolution of $576 \times 576$, which is insufficient for high-quality 3D human generation. To provide more view-consistent frames and realistic texture for 3D human generation, we carefully designed a video augmentation module that includes super-resolution and frame interpolation components. The generated human video is enhanced through this module by undergoing super-resolution and frame interpolation, which results in smooth, high-quality frames at a resolution of $1080 \times 1080$. Lastly, we introduce a 3D human Gaussian splatting module to generate realistic 3D human models. For this, we utilize SMPL (Loper et al. 2023) along with an optimizable feature tensor training strategy to optimize the parameters of the 3D Gaussians, thereby generating a high-quality 3D human from a single image. Figure 1 and 3 demonstrate that Human-VDM achieves state-of-the-art (SOTA) performance and generates realistic 3D humans from a single-view RGB image input. Our contributions can be summarized as follows:

- We propose a novel single-view 3D human generation framework that leverages the human video diffusion model to produce view-consistent human frames.
- We carefully designed a video augmentation model that consists of super-resolution and video frame interpolation to enhance the quality of the generated video.
- We introduce an effective Gaussian Splatting framework for 3D human reconstruction with offset prediction.
- Extensive experiments demonstrate that the proposed Human-VDM can generate realistic 3D humans from single-view images, outperforming state-of-the-art methods in both quality and effectiveness.

## Related Works

**3D Human Generation.** PIFu (Saito et al. 2019) was among the first methods to introduce pixel-aligned features and neural fields (Xie et al. 2022) for reconstructing human figures from images by fitting parametric human shape models such as SMPL (Loper et al. 2023) and SCAPE (Anguelov et al. 2005). PIFuHD (Saito et al. 2020) further enhanced this framework with high-resolution normal guidance. Subsequent methods improved upon this initial approach by integrating additional human body priors. For instance, PaMIR (Zheng et al. 2021) and ICON (Xiu et al. 2022) utilized skinned body models to guide the reconstruction process, while ARCH (Huang et al. 2020), ARCH++ (He et al. 2021), and CAR (Liao et al. 2023) extended this approach by mapping global coordinates into canonical coordinates, enabling reposing. PHORUM (Alldieck, Zanfir, and Sminchisescu 2022) and S3F (Corona et al. 2023) introduced techniques to disentangle shading and albedo, facilitating relighting. Concurrently, another set of methods replaced neural representations with traditional Poisson surface reconstruction (Kazhdan and Hoppe 2013). Despite these advancements, such approaches have been primarily tailored to human bodies and often struggle with the complex topologies of loose clothing. To address this limitation, ECON (Xiu et al. 2023) and 2K2K (Han et al. 2023) integrated depth or normal estimation to enhance the reconstruction process. More recently, Ultraman (Chen et al. 2024) introduced a model to map texture thereby optimizing the texture details thus helping to maintain the color consistency during the final reconstruction. SIFU (Zhang, Yang, and Yang 2024) also proposed a novel approach that combined the 3D Consistent Texture Refinement pipeline with a side-view Decoupling Transformer.

**3D Human Generation with Diffusion models.** Diffusion models (Ramesh et al. 2022) trained on large image datasets have exhibited remarkable capabilities in generating 3D objects from text prompts. Earlier works, such as Fantasia3d (Chen et al. 2023) and Magic3d (Lin et al. 2023), predominantly followed an optimization-based workflow where 3D representations, such as NeRF (Mildenhall et al. 2021), were updated through neural rendering (Tewari et al. 2022). Although a few studies, such as TeCH (Huang et al. 2024), adapted this workflow for 3D human reconstruction, they
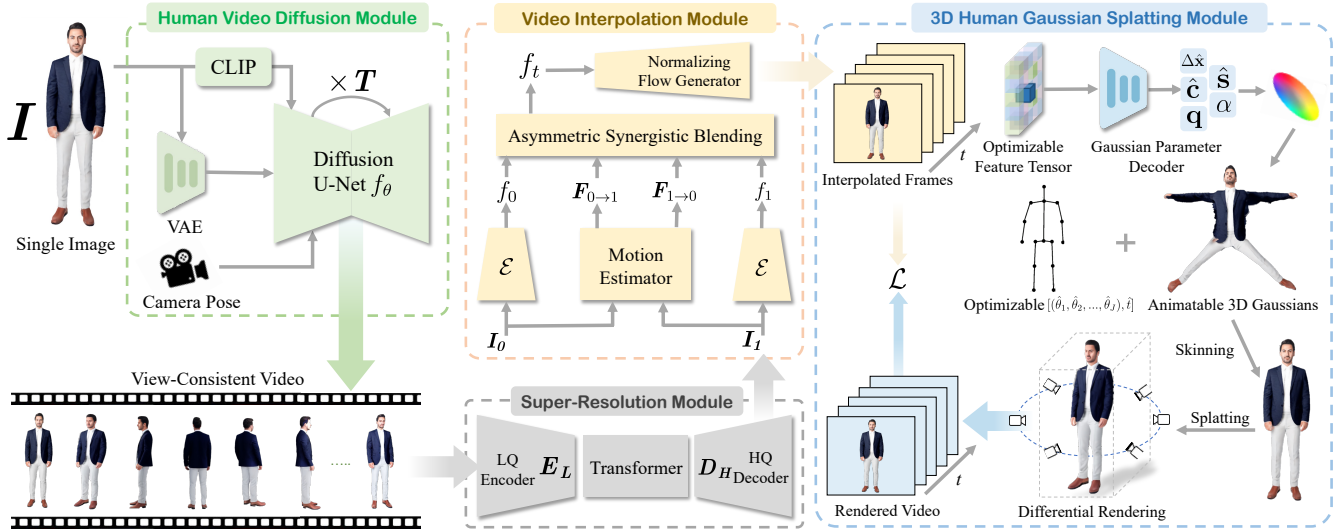
Figure 2: **Human-VDM model architecture.** An image $I$ is first input to a view-consistent human video diffusion module to generate a coherent human video. Next, the video augmentation module applies super-resolution and frame interpolation to enhance texture and generate high-quality interpolated frames. Finally, 3D Human Gaussian splatting learns lifelike 3D humans.

struggled to achieve accurate appearance and geometric representations of the human body due to the inherent ambiguities in text prompt condition. Recently, SiTH (Ho et al. 2024) integrated a generative diffusion model to produce full-body texture and geometry, including unobserved regions, within the reconstruction workflow. However, these methods still face challenges in capturing detailed clothing. In this paper, we leverage a video diffusion model (VDM) to generate an orbital video for 3D human reconstruction.

## Human-VDM

Given a single RGB image $I$ of a person, Human-VDM aims to generate its 3D human model (see Figure 2). Human-VDM comprises several key modules: (i) the Human Video Diffusion module, (ii) the Video augmentation module, which includes the super-resolution and frame interpolation sub-modules, and (iii) the Human Gaussian Splatting module. First, the Human Video Diffusion module generates view-consistent videos of the input image. This video is then processed by the Video Augmentation module, where super-resolution enhances the resolution to $1080 \times 1080$, while video frame interpolation (VFI) smoothens the video frames. Finally, the augmented video is fed into the Human Gaussian Splatting module to generate a high-fidelity 3D human model.

### Human Video Diffusion Module

To generate the video $\hat{V}$, we input the front image of a human, denoted as $I$, into a latent video diffusion model which we fine-tuned for high-quality human video generation. We specifically use SV3D (Voleti et al. 2024), a latent video diffusion model designed for generating videos from a single image, capable of producing consistent multi-view images. However, since SV3D was originally designed

for reconstructing general objects, its generated video quality for human body images is not satisfactory. Therefore, to enhance its capability for human video generation, we fine-tuned SV3D on Thuman 2.0 (Yu et al. 2021) dataset which includes a variety of high-quality human body scans. SV3D produces a raw orbital video, $\hat{V} = [\hat{f}_1, \hat{f}_2, \hat{f}_3, \ldots, \hat{f}_{21}]$, with a resolution of $576 \times 576$, illustrating the human from different viewpoints. The videos generated by the fine-tuned SV3D exhibit superior shape, appearance, and detailed rendering of areas not directly captured in a 2D image. We represent this generation process as follows:

$$\hat{V} = \text{SV3D}(I), \qquad (1)$$

where 'SV3D' denotes the generative process of the fine-tuned SV3D model.

### Video Augmentation Module

The 21-frame human video $\hat{V}$, with a resolution of $576 \times 576$, has limited expressive capacity for detailed 3D human reconstruction. To address this, we introduce the Video Augmentation Module, which includes super-resolution and frame interpolation. Super-resolution helps in improving the quality of textures while video frame interpolation improves the geometric smoothness of the 3D human and the quality of the previously invisible areas.

**Video Super-resolution sub-module.** For image super-resolution on each frame of $\hat{V}$, we employ Code-Former (Zhou et al. 2022), a transformer-based model designed primarily for enhancing facial image resolution. CodeFormer performs Low Quality (LQ) to High Quality (HQ) mapping by first learning a discrete codebook and an HQ decoder $D_H$ through self-reconstruction learning. During Codebook Lookup, a transformer and an LQ encoder $E_L$
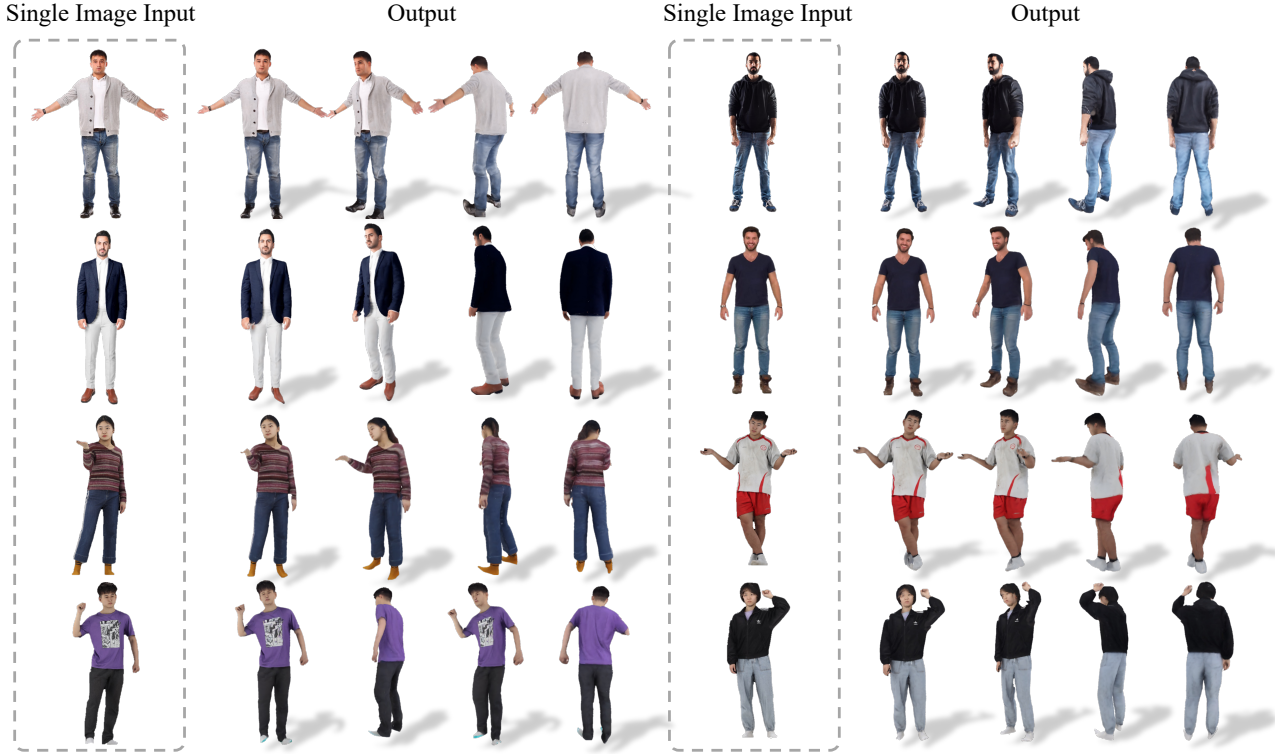
Figure 3: **Qualitative Results.** Novel view results from Human-VDM with various poses, genders, diverse clothing, and different hairstyles demonstrate the robustness of the proposed Human-VDM model. It consistently achieves high photo-realistic quality and precise geometric accuracy. 🔍 **zoom in** for details.

are additionally introduced to accurately model the cookbook code combination. For facial images, increasing the resolution of each frame of $\hat{V}$ by $4\times$ and then resizing it to $1080 \times 1080$ yields clear and realistic images that significantly benefit 3D reconstruction. Similarly, we increase the resolution of each frame in the raw orbital video $\hat{V}$ by $4\times$ and resize it to $1080 \times 1080$, resulting in a high-resolution video $V^{'} = [f_1^{'}, f_2^{'}, ..., f_{21}^{'}]$ with improved texture quality. This process is formulated as follows:

$$f_i^{'} = \text{Resize}(\text{SuperResolution}(\hat{f}_i)), \ 1 \leq i \leq 21, \quad (2)$$

where 'SuperResolution' denotes the operation of CodeFormer, while 'Resize' denotes the operation of resizing the image to $1080 \times 1080$.

**Video Frame Interpolation (VFI) sub-module.** To enhance video consistency and interpolate frames, we employ PerVFI (Wu et al. 2024). VFI provides additional visual information from diverse angles, improving the geometric smoothness of the 3D human and the quality of the invisible areas. PerVFI performs perception-oriented VFI and inputs two reference frame images $I_0$ and $I_1$ to reconstruct intermediate frames. First, bidirectional optical flows, i.e., $F_{0\rightarrow1}$ and $F_{1\rightarrow0}$ are estimated using a motion estimator. Additionally, two encoders capture multi-scale features. These features are then blended using asymmetric synergistic blending to obtain intermediate features $f_t$. These features are finally de-

coded to obtain the intermediate frame using a conditional flow generator, which samples from a normal distribution. We input the 21-frame high-resolution video frames $V^{'}$ into PerVFI, resulting in an 81-frame high-resolution augmented video $V = [f_1, f_2, ..., f_{81}]$. This is formulated as follows:

$$f = \text{VFI}(f_j^{'}), \ 1 \leq j \leq 81, \quad (3)$$

where 'VFI' denotes the frame interpolation operation.

### 3D Human Gaussian Splatting Module

We leverage 3D Gaussian Splatting (Kerbl et al. 2023) to model the 3D human from the augmented human video $V$. 3D Gaussian Splatting employs point-based representation, which facilitates high-quality real-time rendering by modeling the 3D object as a collection of parameterized static 3D Gaussians. Each Gaussian is characterized by a color $c \in \mathbb{R}^3$, a 3D center position $x \in \mathbb{R}^3$, opacity $\alpha \in \mathbb{R}$, a 3D scaling factor $s \in \mathbb{R}^3$, and a 3D rotation $q \in \mathbb{R}^4$.

In this module, we incorporate an appearance network in conjunction with an optimizable feature tensor to enhance the representation of 3D Gaussian models refined from video data (Hu et al. 2024). For each $i^{th}$ frame $f_i$ in the augmented video $V$, we first extract the SMPL model of the human body. We then sample points on the surface of this model and map their positions onto a UV position map, denoted by $m$. We introduce an optimizable feature tensor to
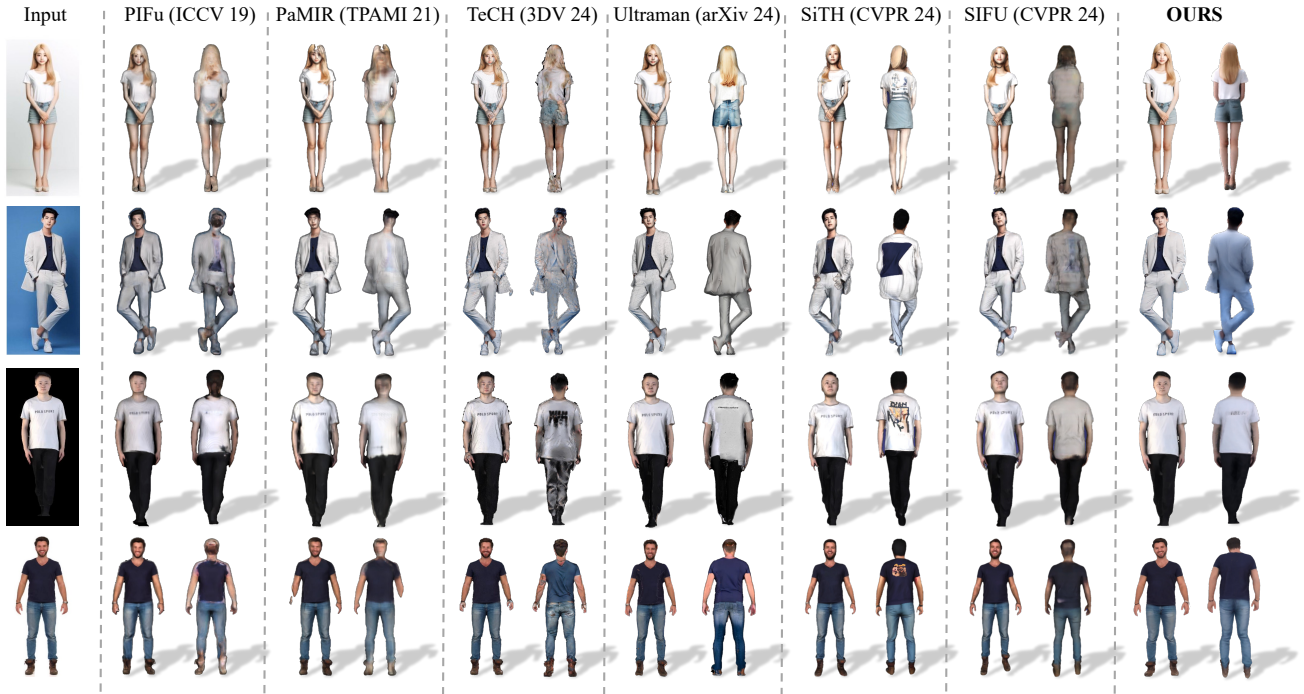
Figure 4: **Qualitative Comparison.** Human-VDM compared to other SOTA models including PIFu (Saito et al. 2019), PaMIR (Zheng et al. 2021), TeCH (Huang et al. 2024), Ultraman (Chen et al. 2024), SiTH (Ho et al. 2024), and SIFU (Zhang, Yang, and Yang 2024). The results demonstrate that Human-VDM achieves superior 3D human generation quality. Note that recent SOTAs fail to predict the unseen back view as shown above. 🔍 **zoom in** for details.

capture the appearance of the reconstructed human. The parameters for each Gaussian are predicted by a Gaussian parameter decoder using the optimizable feature concatenated with $m$ as input. These predictions form the 3D Gaussians in the canonical space. Using Linear Blend Skinning (LBS), these canonical 3D Gaussians can be reposed into motion space for rendering. This is formulated as follows:

$$m = M(\tilde{\theta}, \beta)$$
$$P = Decode(cat(t, m)), \quad (4)$$
$$f_i^r = \text{Splatting}(\text{LBS}(D, J(\beta), \hat{\theta}_i), P),$$

where $\tilde{\theta}$ is the pose parameters of the SMPL model in canonical space and $\beta$ is the average shape parameters calculated from $V$, respectively. $M$ is the operation of mapping the positions of the sampled points on the surface of the SMPL model onto a UV map; $t$ denotes the optimizable feature tensor, $Decode$ means the process of decoding the aligned feature tensors to predict the parameters of Gaussians $P$. $D = T(\beta) + dT$ denotes the locations of 3D Gaussians in canonical space, formed by adding corrective point displacements dT on the template mesh surface $T(\beta)$, $J(\beta)$ produces 3D joint locations, $\hat{\theta}_i$ represents the refined pose parameter optimized from $\theta_i$, which denotes the pose parameters obtained from $f_i$, 'LBS' is the operation of Linear Blend Skinning; 'Splatting' denotes the render process, resulting in a rendered image $f_i^r$.

**Training Objectives.** For formulating the loss function, we take the current frame image $f_i$ as the ground truth and calculate the loss with the rendered image $f_i^r$ for optimization. This is formulated as follows:

$$\mathcal{L} = \lambda_{\text{RGB}}\mathcal{L}_{\text{RGB}} + \lambda_{\text{SSIM}}\mathcal{L}_{\text{SSIM}} + \lambda_{\text{LPIPS}}\mathcal{L}_{\text{LPIPS}}$$
$$+ \lambda_{\text{Offset}}\mathcal{L}_{\text{Offset}} + \lambda_{\text{Scale}}\mathcal{L}_{\text{Scale}} + \lambda_f\mathcal{L}_f, \quad (5)$$

where $\mathcal{L}_{\text{RGB}}$ is the L1-loss between the ground truth and the rendered frame. $\mathcal{L}_{\text{SSIM}}$ and $\mathcal{L}_{\text{LPIPS}}$ denotes the SSIM and LPIPS losses, respectively. $\mathcal{L}_{\text{Offset}}$, $\mathcal{L}_{\text{Scale}}$ and $\mathcal{L}_f$ calculate the L2-norm of predicted offsets and scales, and the feature map, respectively. The weight coefficients $\lambda_{\text{RGB}}$, $\lambda_{\text{SSIM}}$, $\lambda_{\text{LPIPS}}$, $\lambda_{\text{Offset}}$, $\lambda_{\text{Scale}}$ and $\lambda_f$, are set to 0.8, 0.2, 0.2, 10, 1.0 and 1.0 respectively.

## Experiments and Results

**Dataset.** Most works use the popular Thuman 2.0 dataset (Yu et al. 2021), which comprises 2,500 high-quality human body scans, each accompanied by a detailed 3D model and texture mapping. The dataset includes a wide range of action poses and provides the SMPL-X (Pavlakos et al. 2019) parameters along with corresponding grids.

**Evaluation Metrics.** Following previous works on 3D human generation, we use the four major metrics to evaluate the performance of Human-VDM. These include CLIP-Similarity (Radford et al. 2021), LPIPS (Learned Perceptual Image Patch Similarity) (Zhang et al. 2018), SSIM (Wang

Table 1: **User study and Quantitative Comparisons.** Human-VDM compared to recent single-image based 3D human generation SOTAs. Top two results are colored as `first` `second`.

| Method | Venue | User Study | | | | Quantitative Evaluation | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Geometry (%) | Texture (%) | Face (%) | Which is best (%) | CLIP Sim. ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ |
| PIFu (Saito et al. 2019) | ICCV 2019 | 2.33 | 2.00 | 0.33 | 1.67 | 0.8501 | 0.8884 | 0.1615 | 15.0248 |
| PaMIR (Zheng et al. 2021) | TPAMI 2021 | 3.00 | 3.67 | 0.33 | 2.33 | 0.8861 | 0.8924 | 0.1461 | 16.6267 |
| TeCH (Huang et al. 2024) | 3DV 2024 | 3.33 | 2.33 | 4.33 | 4.00 | 0.8875 | 0.8709 | 0.1678 | 15.1464 |
| Ultraman (Chen et al. 2024) | arXiv 2024 | 17.33 | 10.67 | 17.00 | 11.00 | 0.9131 | 0.8958 | 0.1338 | 17.4877 |
| SIFU (Zhang, Yang, and Yang 2024) | CVPR 2024 | 2.67 | 10.33 | 16.67 | 15.67 | 0.8663 | 0.7931 | 0.1500 | 16.4600 |
| SiTH (Ho et al. 2024) | CVPR 2024 | 12.67 | 13.33 | 21.00 | 11.67 | 0.8978 | 0.8963 | 0.1396 | 17.0533 |
| **Human-VDM** | **Ours** | **58.67** | **57.67** | **40.34** | **53.66** | **0.9235** | **0.9228** | **0.0957** | **20.068** |

Table 2: **Ablation studies.** Human-VDM's ablation experiments to verify the effect of proposed components. Without is abbreviated as 'w/o'.

| Ablation | CLIP Sim.↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | CLIP Sim. (Front View)↑ | SSIM (Front View)↑ | LPIPS (Front View)↓ | PSNR (Front View)↑ |
|---|---|---|---|---|---|---|---|---|
| w/o frame interpolation | 0.9234 | 0.9216 | 0.0973 | 20.030 | 0.9286 | 0.9122 | 0.0930 | 19.75 |
| w/o super-resolution | 0.9231 | 0.8981 | 0.0865 | 20.076 | 0.9448 | 0.8857 | 0.0767 | 19.615 |
| w/o fine-tuned SV3D | 0.9146 | 0.9145 | 0.1062 | 18.726 | 0.9449 | 0.9095 | 0.0933 | 19.615 |
| Full | 0.9235 | 0.9228 | 0.0957 | 20.068 | 0.9607 | 0.9257 | 0.0846 | 21.184 |

et al. 2004) and PSNR. CLIP (Radford et al. 2021) measures the similarity between two images, providing a more representative evaluation of image feature similarity. LPIPS (Zhang et al. 2018), measures differences based on learned perceptual image patch similarity, aligning more closely with human perception. Likewise, SSIM (Structural Similarity Index) (Wang et al. 2004) is used to compare the luminance, contrast, and structure between two images. Lastly, PSNR (Peak Signal-to-Noise Ratio) assesses image quality based on pixel-level error, making it an error-sensitive evaluation metric.

**Training details.** To produce high-quality human videos, we fine-tuned SV3D using the Thuman 2.0 dataset (Yu et al. 2021) to enhance its 3D human video generation capabilities. We selected 475 samples from Thuman 2.0, excluding those used in subsequent quantitative comparisons. For each sample, 21 images were rendered from various angles following (Xiu et al. 2022). All images corresponding to a sample are rendered at the same horizontal position with a constant angular interval of $360/21$ degree to ensure the consistency of rendered multi-view images. The first rendered image of each body was employed as the input, while the remaining images served as ground truth for fine-tuning SV3D. We freeze the image encoder and decoder of the original SV3D (Voleti et al. 2024) model and optimize the U-Net weights (Ronneberger, Fischer, and Brox 2015). The learning rate was set to `5e-6` and fine-tuned on one NVIDIA A800 GPU with a batch size of 13.

## Qualitative Comparison

Figure 3 presents the qualitative 3D human generation results from Human-VDM on a variety of input images that differ in gender, body posture, lighting, color, and clothing styles. The results demonstrate Human-VDM's signif-

icant performance with high appearance consistency, texture, and geometry qualities. Next, we compare Human-VDM with recent SOTA works on single-image based 3D human generation (see Figure 4), including PIFu (Saito et al. 2019), PaMIR (Zheng et al. 2021), TeCH (Huang et al. 2024), Ultraman (Chen et al. 2024), SiTH (Ho et al. 2024) and SIFU (Zhang, Yang, and Yang 2024). Compared to Human-VDM, PaMIR (Zheng et al. 2021) exhibits significant shortcomings in the geometry of the generated 3D human, e.g., the body of the generated human is incomplete for the first image. On the other hand, TeCH (Huang et al. 2024), PIFu (Saito et al. 2019), and SiTH (Ho et al. 2024) reconstruct remarkable geometries but contain apparent artifacts. Likewise, SIFU (Zhang, Yang, and Yang 2024) displays misalignment in character motion and suboptimal texture quality on the back of the generated human. While Ultraman (Chen et al. 2024) obtains good geometry but fails to predict the realistic appearance of unseen view. Therefore, the proposed Human-VDM outperforms SOTA models in terms of texture quality and appearance consistency.

## Quantitative Comparison

Following previous methods (Chen et al. 2024), we randomly selected 50 samples from Thuman 2.0 (Yu et al. 2021). Four views of the ground truth (GT), i.e., front, back, left, and right, were used to compute scores between the reconstructed results and the GT across these views. As reported in Table 1, Human-VDM achieves the lowest LPIPS and the highest CLIP score, indicating that the rendered images produced by our method are highly consistent with the input images. Additionally, Human-VDM achieves the highest SSIM and PSNR scores, further demonstrating that the rendered images of the generated 3D human are most closely aligned with the ground truth. All reported scores demon-

strate the superiority of the proposed Human-VDM over existing SOTA methods.

## User Study

The discussed metrics may not always fully capture the quality of generated 3D humans in terms of realism and other details. Thus following previous works, a user preference study was conducted to evaluate the performance of Human-VDM against existing SOTA methods. We compare Human-VDM with six recent SOTA models using 10 different samples, each with four views of generated 3D humans in different samples. For each sample, 30 volunteers were asked to vote on their impressions regarding four key aspects: geometry quality, texture quality, face quality, and overall quality. For a fair comparison, the results for the other six SOTA models were generated using their official code, with all settings left at their default values. As shown in Table 1, the proposed Human-VDM surpasses SOTA models in the aforementioned aspects.

Most volunteers considered Human-VDM to generate the best results, especially in terms of geometry and texture. Though Human-VDM does not particularly dominate in face quality relatively, it performs the best face consistency with the input image as shown in Figure 4. More than $53\%$ of the volunteers confirm that Human-VDM outperforms other SOTA models, which confirms Human-VDM's superiority.

## Ablation Study

We performed ablation studies by systematically excluding various components to assess the effectiveness of the proposed modules through both quantitative and qualitative comparisons. For this analysis, we randomly selected 30 samples from the Thuman 2.0 dataset (Yu et al. 2021). We compared the full model with the variants excluding the proposed modules using the CLIP Similarity (Radford et al. 2021), SSIM (Wang et al. 2004), LPIPS (Zhang et al. 2018), and PSNR metrics. The evaluation covered rendered results from four viewpoints: front, back, left, and right. We additionally report results solely for the front view as well. Table 2 presents the quantitative comparisons, while the qualitative visual comparisons are illustrated in Figure 5.

Quantitative results demonstrate that the proposed full model achieves superior CLIP Similarity and SSIM across both the single view and four views. The visual ablation results further establish that the 3D human generated by the full model exhibits more photorealistic textures and precise geometry. Results produced without finetuned SV3D are less lifelike and realistic since the videos generated by the original SV3D are not satisfactory. Without Super-Resolution, the video frames are not distinct enough for the Human Gaussian Splatting module, which results in blurs and artifacts of the reconstructed humans. Due to the lack of features presented by only 21 frames, results generated without frame interpolation are not good enough yet, which has apparent artifacts in novel views. This confirms the significance and contribution of the video augmentation module. In general, the finetuned SV3D provides high-quality human orbital video for realistic reconstruction; the super-resolution module enhances the quality of video frames to

generate more distinct results, and the VFI module enables the model to generate remarkable results in novel views. Although the full model shows a slight decrease in LPIPS and PSNR, the visual results indicate that the 3D human reconstructed by the complete model is of higher quality. Overall, the full model achieves better performance i.e., when including the proposed components. This confirms the effectiveness of the proposed modules.
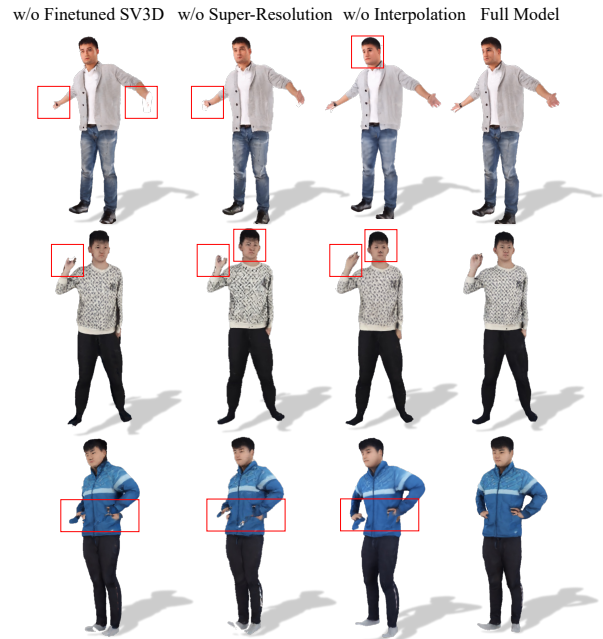


Figure 5: **Qualitative Visual Ablation Comparisons.** Compared to other variants, the proposed full model achieves highly realistic textures and accurate geometry.

## Conclusion and Future Work

We propose a novel 3DGS-based framework for generating 3D humans from a single RGB image leveraging human video diffusion models. We first generate a view-consistent orbital video around the human and then augment the video through super-resolution and video frame interpolation. Finally, we reconstruct a remarkable 3D human using 3D Gaussian with the enhanced video. Both quantitative and qualitative experiments demonstrate that Human-VDM excels in generating 3D humans from a single image, outperforming state-of-the-art methods.

**Limitations and Future works.** Human-VDM has two limitations. First, it is challenging to accurately generate precise finger geometry due to the intricate and small size of finger poses. Second, applying large video diffusion models limits the model's overall ability to achieve a real-time 3D human generation. Future works can focus on addressing these limitations by enhancing geometry generation for complex and small finger poses, as well as developing more efficient models that can achieve real-time 3D human generation.

# References

Alldieck, T.; Zanfir, M.; and Sminchisescu, C. 2022. Photorealistic monocular 3d reconstruction of humans wearing clothing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1506–1515.

Anguelov, D.; Srinivasan, P.; Koller, D.; Thrun, S.; Rodgers, J.; and Davis, J. 2005. Scape: shape completion and animation of people. In *ACM SIGGRAPH 2005 Papers*, 408–416. ACM.

Blattmann, A.; Dockhorn, T.; Kulal, S.; Mendelevitch, D.; Kilian, M.; Lorenz, D.; Levi, Y.; English, Z.; Voleti, V.; Letts, A.; et al. 2023. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*.

Chen, M.; Chen, J.; Ye, X.; Gao, H.-a.; Chen, X.; Fan, Z.; and Zhao, H. 2024. Ultraman: Single Image 3D Human Reconstruction with Ultra Speed and Detail. *arXiv preprint arXiv:2403.12028*.

Chen, R.; Chen, Y.; Jiao, N.; and Jia, K. 2023. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *Proceedings of the IEEE/CVF international conference on computer vision*, 22246–22256.

Corona, E.; Zanfir, M.; Alldieck, T.; Bazavan, E. G.; Zanfir, A.; and Sminchisescu, C. 2023. Structured 3d features for reconstructing controllable avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16954–16964.

Deitke, M.; Schwenk, D.; Salvador, J.; Weihs, L.; Michel, O.; VanderBilt, E.; Schmidt, L.; Ehsani, K.; Kembhavi, A.; and Farhadi, A. 2023. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 13142–13153.

Han, S.-H.; Park, M.-G.; Yoon, J. H.; Kang, J.-M.; Park, Y.-J.; and Jeon, H.-G. 2023. High-fidelity 3d human digitization from single 2k resolution images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12869–12879.

He, T.; Xu, Y.; Saito, S.; Soatto, S.; and Tung, T. 2021. Arch++: Animation-ready clothed human reconstruction revisited. In *Proceedings of the IEEE/CVF international conference on computer vision*, 11046–11056.

Ho, I.; Song, J.; Hilliges, O.; et al. 2024. Sith: Single-view textured human reconstruction with image-conditioned diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 538–549.

Hu, L.; Zhang, H.; Zhang, Y.; Zhou, B.; Liu, B.; Zhang, S.; and Nie, L. 2024. Gaussianavatar: Towards realistic human avatar modeling from a single video via animatable 3d gaussians. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 634–644.

Huang, Y.; Yi, H.; Xiu, Y.; Liao, T.; Tang, J.; Cai, D.; and Thies, J. 2024. Tech: Text-guided reconstruction of lifelike clothed humans. In *2024 International Conference on 3D Vision (3DV)*, 1531–1542. IEEE.

Huang, Z.; Xu, Y.; Lassner, C.; Li, H.; and Tung, T. 2020. Arch: Animatable reconstruction of clothed humans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3093–3102.

Johnson, J.; Alahi, A.; and Fei-Fei, L. 2016. Perceptual losses for real-time style transfer and super-resolution. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*, 694–711. Springer.

Kazhdan, M.; and Hoppe, H. 2013. Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(3): 1–13.

Kerbl, B.; Kopanas, G.; Leimkühler, T.; and Drettakis, G. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Trans. Graph.*, 42(4): 139–1.

Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25.

Liao, T.; Zhang, X.; Xiu, Y.; Yi, H.; Liu, X.; Qi, G.-J.; Zhang, Y.; Wang, X.; Zhu, X.; and Lei, Z. 2023. High-fidelity clothed avatar reconstruction from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8662–8672.

Lin, C.-H.; Gao, J.; Tang, L.; Takikawa, T.; Zeng, X.; Huang, X.; Kreis, K.; Fidler, S.; Liu, M.-Y.; and Lin, T.-Y. 2023. Magic3d: High-resolution text-to-3d content creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 300–309.

Loper, M.; Mahmood, N.; Romero, J.; Pons-Moll, G.; and Black, M. J. 2023. SMPL: A skinned multi-person linear model. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, 851–866. ACM.

Mildenhall, B.; Srinivasan, P. P.; Tancik, M.; Barron, J. T.; Ramamoorthi, R.; and Ng, R. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1): 99–106.

Pavlakos, G.; Choutas, V.; Ghorbani, N.; Bolkart, T.; Osman, A. A.; Tzionas, D.; and Black, M. J. 2019. Expressive body capture: 3d hands, face, and body from a single image. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10975–10985.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.

Ramesh, A.; Dhariwal, P.; Nichol, A.; Chu, C.; and Chen, M. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2): 3.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10684–10695.

Ronneberger, O.; Fischer, P.; and Brox, T. 2015. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, 234–241. Springer.

Saito, S.; Huang, Z.; Natsume, R.; Morishima, S.; Kanazawa, A.; and Li, H. 2019. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2304–2314.

Saito, S.; Simon, T.; Saragih, J.; and Joo, H. 2020. 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*, 84–93.

Tewari, A.; Thies, J.; Mildenhall, B.; Srinivasan, P.; Tretschk, E.; Wang, Y.; Lassner, C.; Sitzmann, V.; Martin-Brualla, R.; Lombardi, S.; Simon, T.; Theobalt, C.; Niessner, M.; Barron, J. T.; Wetzstein, G.; Zollhoefer, M.; and Golyanik, V. 2022. Advances in Neural Rendering. arXiv:2111.05849.

Vaswani, A. 2017. Attention is all you need. *arXiv preprint arXiv:1706.03762*.

Voleti, V.; Yao, C.-H.; Boss, M.; Letts, A.; Pankratz, D.; Tochilkin, D.; Laforte, C.; Rombach, R.; and Jampani, V. 2024. Sv3d: Novel multi-view synthesis and 3d generation from a single image using latent video diffusion. *arXiv preprint arXiv:2403.12008*.

Wang, Z.; Bovik, A. C.; Sheikh, H. R.; and Simoncelli, E. P. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4): 600–612.

Wu, G.; Tao, X.; Li, C.; Wang, W.; Liu, X.; and Zheng, Q. 2024. Perception-Oriented Video Frame Interpolation via Asymmetric Blending. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2753–2762.

Xie, Y.; Takikawa, T.; Saito, S.; Litany, O.; Yan, S.; Khan, N.; Tombari, F.; Tompkin, J.; Sitzmann, V.; and Sridhar, S. 2022. Neural Fields in Visual Computing and Beyond. arXiv:2111.11426.

Xiu, Y.; Yang, J.; Cao, X.; Tzionas, D.; and Black, M. J. 2023. Econ: Explicit clothed humans optimized via normal integration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 512–523.

Xiu, Y.; Yang, J.; Tzionas, D.; and Black, M. J. 2022. Icon: Implicit clothed humans obtained from normals. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 13286–13296. IEEE.

Yu, T.; Zheng, Z.; Guo, K.; Liu, P.; Dai, Q.; and Liu, Y. 2021. Function4D: Real-time Human Volumetric Capture from Very Sparse Consumer RGBD Sensors. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR2021)*.

Zhang, R.; Isola, P.; Efros, A. A.; Shechtman, E.; and Wang, O. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 586–595.

Zhang, Z.; Yang, Z.; and Yang, Y. 2024. Sifu: Side-view conditioned implicit function for real-world usable clothed human reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9936–9947.

Zheng, Z.; Yu, T.; Liu, Y.; and Dai, Q. 2021. Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. *IEEE transactions on pattern analysis and machine intelligence*, 44(6): 3170–3184.

Zhou, S.; Chan, K.; Li, C.; and Loy, C. C. 2022. Towards robust blind face restoration with codebook lookup transformer. *Advances in Neural Information Processing Systems*, 35: 30599–30611.

Zhu, J.-Y.; Park, T.; Isola, P.; and Efros, A. A. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, 2223–2232.

In the supplementary material, we provide a more detailed explanation of the model architecture, as well as training specifics, such as loss function weights, dataset descriptions, and definitions of the evaluation metrics. Additionally, we include further visual results and an analysis of failure cases.

## Model Architecture Details

### Human Video Diffusion Module

**Module Architecture.** The Video Diffusion Module of Human-VDM is based on SV3D (Voleti et al. 2024). SV3D's architecture builds upon SVD (Blattmann et al. 2023) and consists of a UNet (Ronneberger, Fischer, and Brox 2015) model with multiple layers. Each layer comprises a sequence of 1 residual block with Conv3D layers, followed by spatial and temporal transformer blocks integrated with attention layers. After being embedded into the latent space via the visual autoencoder (VAE) of SVD, the conditioning image is concatenated with the noisy latent state input $z_t$ at noise timestep $t$ before being fed into the UNet. The CLIP-embedding (Radford et al. 2021) matrix of the input image is provided to the cross-attention layers of each transformer block (Vaswani 2017), serving as the key and value, with the layer's feature acting as the query. Along with the diffusion noise timestep, the camera trajectory is also incorporated into the residual blocks. The camera pose angles $e_i$ and $a_i$ are first embedded into the position embeddings. These camera pose embeddings are then concatenated, linearly transformed, and combined with the noise timestep embedding. The composite embedding is fed into every residual block, where it is added to the block's output after another linear transformation to match the feature size.

**Static Orbits.** The original SV3D model (Voleti et al. 2024) consists of two main orbits: (1) the static orbit and (2) the dynamic orbit. Our study utilizes the static orbit, where the camera moves around the object at evenly spaced azimuth angles while maintaining the same elevation angle as in the conditioning image.

**Fine-tuning SV3D for Human Video Diffusion.** The original SV3D is fine-tuned upon SVD-xt (Blattmann et al. 2023) on the Objaverse dataset (Deitke et al. 2023), which contains synthetic 3D objects covering a wide diversity. For each object, (Voleti et al. 2024) renders 21 frames around it on a random color background at $576 \times 576$ resolution, field-of-view of 33.8 degrees. We adopt the same rendering strategy for the Thuman 2.0 dataset (Yu et al. 2021) to fine-tune SV3D for high-quality human video generation.

### Video Augmentation Module

**Video Super-Resolution sub-module.** CodeFormer (Zhou et al. 2022) is a transformer-based model (Vaswani 2017) to enhance the resolution of human images. Upon learning a discrete codebook, an encoder $E_H$ embed the high-quality human image $I_h \in \mathbb{R}^{H \times W \times 3}$ as a compressed feature $Z_h \in \mathbb{R}^{m \times n \times d}$ by an encoder $E_H$. Each "pixel" in $Z_h$ is then replaced by the nearest entry in the learnable codebook $\mathcal{C} = c_k \in \mathbb{R}^{d N}_{k=0}$. Afterward, the quantized feature $Z_c \in \mathbb{R}^{m \times n \times d}$ along with the code token sequence $s \in 0, \cdots, N - 1^{m \cdot n}$ are produced as the following:

$$Z_c^{(i,j)} = \arg \min_{c_k \in \mathcal{C}} \| Z_h^{(i,j)} - c_k \|_2,$$
$$s^{(i,j)} = \arg \min_k \| Z_h^{(i,j)} - c_k \|_2. \tag{6}$$

Given $Z_c$, the high-quality human image $I_{rec}$ is reconstructed by the decoder $D_H$. The $m \times n$ code token sequence, denoted as $s$, constitutes a novel latent discrete representation, which encodes the specific indices corresponding to entries in the learned codebook, i.e., $Z_c^{(i,j)} = c_k$ when $s^{(i,j)} = k$.

Subsequently, with the codebook $\mathcal{R}$ and decoder $D_H$ held constant, a Transformer module (Vaswani 2017) is introduced for predicting the code sequence, capturing the global human composition from low-quality inputs. To extract the low-quality features $Z_l \in \mathbb{R}^{m \times n \times d}$ using $E_L$, the features are first unfolded to $m \cdot n$ vectors $Z_l^v \in \mathbb{R}^{(m \cdot n) \times d}$, which are subsequently fed into the Transformer. In the transformer, the $s^{th}$ self-attention block performs the below operation:

$$X_{s+1} = \text{Softmax}(Q_s K_s)V_s + X_s, \tag{7}$$

where $X_0 = Z_l^v$. $X_s$ is used to get the queries $Q$, key $K$, and value $V$ through linear layers.

**Video Frame Interpolation (VFI) sub-module.** PerVFI is a novel model of frame interpolation. Given two reference frame images, $I_0$ and $I_1 \in \mathbb{R}^{H \times W \times 3}$, with height $H$ and width $W$, PerVFI is designed for reconstructing the intermediate frame $I_t$ within the target time $t \in (0, 1)$. It incorporates an asymmetric synergistic blending (ASB) module and a conditional normalizing flow-based generator.

After estimating bidirectional optical flows, PerVFI presents a pyramidal architecture, which can better capture multiscale information to extract features at different scales. Specifically, a feature encoder $E_\theta$ is used to encode the two images into pyramid features with $L$ levels, which can be denoted as $f_i = E_{\theta(I_i)}$, $i = 0, 1$. Subsequently, a feature blending module, denoted as $B_\theta$, blends the pyramidal features to produce intermediate pyramid features. Afterward, a conditional normalizing flow-based generator $G_\phi$, which is invertible, decodes $f_t$ into the output frame $I_t$. The output is formulated as $I_t = G_\phi^{-1}(r; f_t)$, where $r \sim \mathcal{N}(0, \tau) \in \mathbb{R}^{H \times W \times 3}$ represents a variable drawn from a normal distribution with a temperature parameter $\tau$; $f_t$ is the feature pyramid with $L$ levels.

### 3D Human Gaussian Splatting Module

In 3D Gaussian, human appearances are determined by point displacements $dT$ and properties $\mathbf{P}$. Modeling dynamic human appearances involves estimating these evolving properties. We propose a dynamic appearance network coupled

Figure 6: Additional results comparing Human-VDM with SOTA models. The results demonstrate that Human-VDM achieves superior 3D human generation quality. 🔍 **zoom in** for details.

Figure 7: In-the-wild testing results comparing Human-VDM with SOTA models. The results demonstrate that Human-VDM achieves superior 3D human generation quality. 🔍 **zoom in** for details.
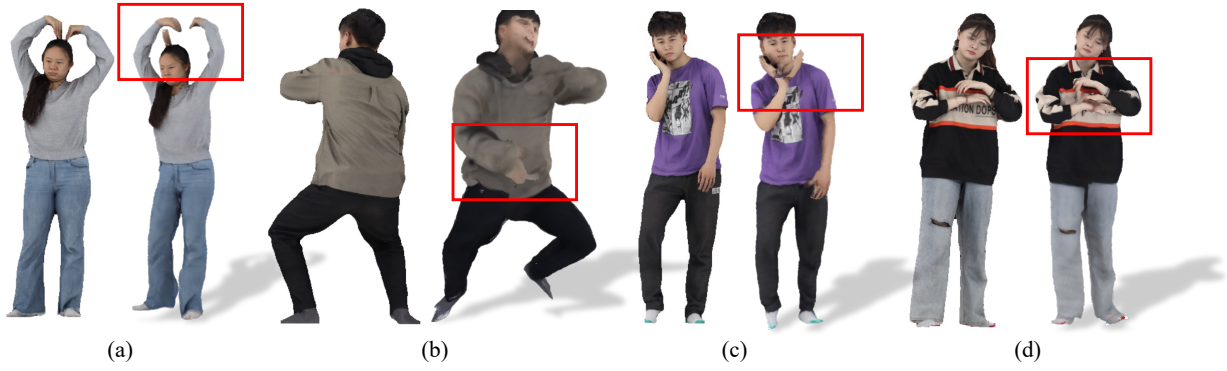
Figure 8: Failure cases of Human-VDM. The intricate and small size of fingers makes it challenging to accurately generate precise finger geometry, as shown in (a) and (b). Moreover, we can see the case of (c) hand-face and (d) hand-hand interactions, which remain challenging in 3D human generation. The left image shows the input, while the right is the generated 3D human.

with an optimizable feature tensor to effectively capture dynamic human appearances across various poses. The dynamic appearance network is designed to learn a mapping from a 2D manifold representing the underlying human shape to the dynamic properties of 3D Gaussians as follows:

$$f_\phi : \mathcal{S}^2 \in \mathbb{R}^3 \to \mathbb{R}^7, \quad (8)$$

the 2D human manifold $\mathcal{S}^2$ is depicted by a UV positional map $I \in \mathbb{R}^{H \times W \times 3}$, where each valid pixel stores the position $(x, y, z)$ of one point on the posed body surface. The final predictions consist of per point offset $\Delta\hat{\mathbf{x}} \in \mathbb{R}^3$, color $\hat{\mathbf{c}} \in \mathbb{R}^3$, and scale $\hat{s} \in \mathbb{R}$.

Human poses $\boldsymbol{\theta}$ and translations $t$ estimated from monocular videos are usually inaccurate. Hence, the 3D Gaussians reposed in motion space may be inaccurately represented, potentially resulting in unsatisfactory rendering outcomes. To address this issue, we jointly optimize human motions and appearances. We update the estimated body poses and translations by calculating $(\Delta\boldsymbol{\theta}, \Delta\mathbf{t})$ to refine human motions, which can be formulated as follows:

$$\hat{\boldsymbol{\Theta}} = (\boldsymbol{\theta} + \Delta\boldsymbol{\theta}, \mathbf{t} + \Delta\mathbf{t}). \quad (9)$$

We modify $\theta$ in the equation of animatable Gaussians in the main article using $\hat{\boldsymbol{\Theta}}$ to render the proposed animatable 3D Gaussians differentiable with respect to the motion conditions. Finally, the current frame image is taken as the ground truth to calculate the loss with the rendered image.

**Training Objectives**

We use the current frame image, i.e., $f_i$, and the rendered image, i.e., $f_i^r$, for supervising the Human-VDM model. The total loss consists of six different loss functions which include $\mathcal{L}_{\text{RGB}}$, $\mathcal{L}_{\text{SSIM}}$, $\mathcal{L}_{\text{LPIPS}}$, $\mathcal{L}_{\text{Offset}}$, $\mathcal{L}_{\text{Scale}}$ and $\mathcal{L}_f$. In this section, we describe the loss functions in greater detail.

$\mathcal{L}_{\text{RGB}}$ is the L1-loss between the ground truth and the rendered frame and is formulated as:

$$\mathcal{L}_{\text{RGB}}(x, y) = \frac{1}{HW} \sum_{h,w}^{HW} |y_{hw} - x_{hw}|, \quad (10)$$

$\mathcal{L}_{\text{SSIM}}$ (Wang et al. 2004), or the Structural Similarity Index Metric Loss is a perceptual metric to measure the similarity between two images, taking luminance, contrast, and structure into account. We define the SSIM loss as follows:

$$\mathcal{L}_{\text{SSIM}}(x, y) = 1 - \text{SSIM}(x, y)$$
$$= 1 - \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)}, \quad (11)$$

where $\mu_x$ and $\mu_y$ stands for the mean of $x$ and $y$; $\sigma_x$ and $\sigma_y$ represent the variance of $x$ and $y$, while $\sigma_{xy}$ denote the covariance of $x$ and $y$.

$\mathcal{L}_{\text{LPIPS}}$ (Zhang et al. 2018) measures image similarity, which evaluates the perceptual difference between two images through deep learning models. In this paper, we utilize AlexNet (Krizhevsky, Sutskever, and Hinton 2012) for extracting features of images. We calculate $\mathcal{L}_{\text{LPIPS}}$ as:

$$\mathcal{L}_{\text{LPIPS}}(x, y) = \sum_l \frac{1}{H_l W_l} \sum_{h,w} ||w_l \odot (\hat{f}_{xhw}^l - \hat{f}_{yhw}^l)||_2^2, \quad (12)$$

where $\hat{f}_{xhw}^l$ represents the feature output of image $x$ in layer $l$ at the pixel $hw$, and $\hat{f}_{yhw}^l$ means the same of image $y$. $w_l$ is a trainable parameter in layer $l$.

$\mathcal{L}_{\text{Offset}}$, $\mathcal{L}_{\text{Scale}}$ and $\mathcal{L}_f$ calculate the L2-norm of the feature map, predicted offsets and scales on the canonical surface, respectively. We formulate them as follows:

$$\mathcal{L}_{\text{Offset}} = \frac{1}{N} \sum_{i=1}^N (\Delta\hat{x}_i)^2, \quad (13)$$

where $\Delta x_i$ denote the predicted offset of $i^{th}$ gaussian.

$$\mathcal{L}_{\text{Scale}} = \frac{1}{N} \sum_{i=1}^N (\hat{s}_i)^2, \quad (14)$$

where $s_i$ denotes the predicted scale of $i^{th}$ gaussian.

$$\mathcal{L}_{\text{f}} = \frac{1}{F} \sum_{i=1}^F (t_i)^2, \quad (15)$$

where $t_i$ denotes the optimized feature.

# Implementation Details

In this section, we present additional details on the model implementation. The Gaussian decoder is implemented as an MLP. A total of 202,738 Gaussians were initially sampled on the surface of the canonical SMPL model. The adjustable coefficient $w$, which presents the reliance on input low-quality image, is set to $0.7$ in the Super-Resolution module. For each sample, we train the dynamic appearance network on a single NVIDIA RTX 3090 GPU for 1000 epochs with a batch size of 2. The learning rate of the network is set to `3e-3`.

# Additional Results

In this section, we present additional results, including in-the-wild testing and failure cases.

### In-the-wild visual results

To demonstrate the superiority of Human-VDM, we provide more visual comparison results. This includes additional results as shown in Figure 6, including results on challenging in-the-wild cases illustrated in Figure 7.

### Failure Cases

In this subsection, we present several cases of failure in Human-VDM. Although Human-VDM performs exceptionally well in generating 3D humans from a single RGB image, it still has a few limitations and failure cases, as discussed in the main text. Figure 8 shows the failure cases of Human-VDM. For example, when the human in the input image interacts with their hands against their body, some artifacts may appear at the contact region.