# MVHuman: Tailoring 2D Diffusion with Multi-view Sampling For Realistic 3D Human Generation

Suyi Jiang   Haimin Luo   Haoran Jiang   Ziyu Wang   Jingyi Yu   Lan Xu

ShanghaiTech University

## Abstract

*Recent months have witnessed rapid progress in 3D generation based on diffusion models. Most advances require fine-tuning existing 2D Stable Diffsuions into multi-view settings or tedious distilling operations and hence fall short of 3D human generation due to the lack of diverse 3D human datasets. We present an alternative scheme named MVHuman to generate human radiance fields from text guidance, with consistent multi-view images directly sampled from pre-trained Stable Diffsuions without any fine-tuning or distilling. Our core is a multi-view sampling strategy to tailor the denoising processes of the pre-trained network for generating consistent multi-view images. It encompasses view-consistent conditioning, replacing the original noises with "consistency-guided noises", optimizing latent codes, as well as utilizing cross-view attention layers. With the multi-view images through the sampling process, we adopt geometry refinement and 3D radiance field generation followed by a subsequent neural blending scheme for free-view rendering. Extensive experiments demonstrate the efficacy of our method, as well as its superiority to state-of-the-art 3D human generation methods.*

## 1. Introduction

The 3D creation of us humans with photo-realism serves as the cornerstone for numerous applications like telepresence or immersive experiences in VR/AR. Early attempts [2] generally require expensive apparatus and immense artistic expertise and hence are limited to celebrities in feature films. Democratizing the accessible use of realistic human avatars to the mass crowd remains unsolved for the vision communities.

The Diffusion models [16], i.e., Stable Diffusion (SD) [66], have demonstrated high-fidelity and diverse 2D human generation, from simple text prompts. Using ControlNet [92], they can even generate hyper-realistic human images under various viewpoints and poses, yield-



Figure 1. Given text guidance, MVHuman generates free-view rendering results and fine-grained geometry with the aid of multi-view images sampled from pre-trained 2D diffusion models.

ing huge potential for 3D human generation. However, directly training an analogous Diffusion model under 3D representations [11, 57] for 3D human generation is infeasible, mainly due to the severe lack of diverse and high-quality human scans. Thus, various 2D-lifting approaches [12, 36, 62, 77, 81] explore to distill pre-trained 2D diffusion models to optimize certain 3D representations like mesh or NeRF [52]. However, they suffer from inefficient optimization with slow convergence and the multifaced Janus artifacts due to the lack of 3D awareness. Some recent methods [6, 21, 31, 34, 89] tailor such distillation scheme for 3D human generation. The utilization of human motion and shape priors [43] alleviates the Janus artifacts, but the inherent inefficiency and the cartoon-like saturated appearances caused by the distillation remain. Instead of directly using the pre-trained SD model [66], recent methods [38, 41, 71] turn to train multi-view diffusion models, which serves as multi-view priors to significantly improve the subsequent 3D content generation. However, training or fine-tuning such multi-view diffusion models still heavily relies on multi-view image datasets from real-world collection [87] or CG rendering [13]. As a result, it's still difficult to extend such multi-view strategies for 3D human generation due to the lack of human datasets.

In this paper, we present *MVHuman* – a novel scheme to generate human assets from text guidance, with the aid of consistent multi-view images directly sampled from pre-trained 2D diffusion models (see Fig. 1). In stark contrast to prior arts, we directly apply an existing 2D latent diffusion model [66] with ControlNet [92] to obtain such multi-view priors, without tedious fine-tuning or distilling.

The core of our MVHuman is a multi-view sampling process compatible with the 2D diffusion model [66], which bridges various views by carefully tailoring their input conditions, predicted noises, and the corresponding latent codes. Specifically, we first apply an off-the-shelf monocular reconstructor [84] on the generated front-view image to obtain a coarse geometry proxy and project it into various target views to obtain 2D skeletal poses, normal, and depth maps. These 2D attributes serve as the view-consistent conditions for utilizing the ControlNet [92]. Secondly, we introduce the concept of "consistency-guided noise". Note that in the deterministic sampler [44, 72], with a fixed latent code of a certain sampling step, we can obtain the predicted original signal from the corresponding predicted noise and vice verse. Thus, for sampling processes with different initial random noises, we can blend their predicted original signals at a sampling step, so as to obtain a consistent prediction and transform it back to the corresponding consistency-guided noises of various sampling processes. By replacing the original noises with such consistency-guided ones, various sampling processes will finally recover a consistent signal. We extend such consistency-guided noises into the multi-view setting by treating each view with an individual sampling process. Specifically, for a certain sampling step of a target view, we warp the decoded images from the predicted original signals of its adjacent views to the current view through depth-based warping and then encode these warped images back to the latent space of SD model [66]. We further apply an occlusion-aware blending strategy to obtain the consistent predictions and the subsequent consistency-guided noises, so as to replace the original predicted noises and generate consistent multi-view images. To further improve the view consistency, we explicitly optimize the latent codes of adjacent views under various sampling steps. We decode them into the image space and warp them to each other through depth-aware warping to calculate their consistency. Besides, inspired by recent video diffusion methods [7, 29], we modify the self-attention layer in the pre-trained SD network [66] to concatenate attention features from a reference view. Such a strategy further enhances the appearance similarity of the generated multi-view images.

Finally, with the above multi-view sampling, we generate the desired human images covering both full-body and upper-body views and subsequently generate radiance fields followed by a blending scheme for free-view rendering. We first add fine-grained details into the coarse proxy using the implicit geometric information within the images and then train a radiance field based on multi-plane features [5]. We further adapt the neural blending strategy [26, 73] into a two-stage coarse-to-fine setting. It blends the initial radiance rendering results with both the full-body and upper-body human images, so as to provide high-fidelity novel view synthesis. As an additional benefit, our MVHuman can seamlessly obtain many functions of the original SD model [66], such as text-based editing, or loading and upgrading a pre-trained LoRA [19] model into 3D results. To summarize, our main contributions include:

- We present a novel scheme to generate high-quality human assets, directly using pre-trained 2D diffusion models without fine-tuning or distilling.
- We introduce a dedicated multi-view sampling process with consistency-guided noise and latent code optimization, to generate view-consistent images.
- We utilize the generated multi-view images to refine geometry and adopt a tiered neural blending scheme on radiance fields to enable free-view rendering.

## 2. Related Work

**Text-guided 3D Content Generation.** Recent rapid progress in diffusion models within the text-to-image domain [16, 66, 92] has enhanced research interest in 3D content generation. Early works [50, 62, 76] propose Score Distillation Sampling (SDS) algorithm to lift pre-trained 2D diffusion models to optimize 3D NeRF [52]. The following works [12, 36] extend SDS to optimize other efficient 3D representations [32, 54, 69], or only generate textures [10, 45, 65] for existing meshes in pursuit of faster speed and quality. However, such SDS-based approaches suffer from oversaturation problems. ProlificDreamer [81] addresses such problems via Variational Score Distillation (VSD) but costs much longer optimization. Another line of works [22, 37, 38, 41, 48, 63, 64, 70, 74, 85, 88] turn to reconstruct 3D content from a single image by distillation. Their challenge evolves into altering the generation distribution by fine-tuning the diffusion network. Recent multi-view diffusion models [39, 71, 75] propose to generate view-consistent 2D images and subsequentially benefit 3D generation significantly, heavily relying on cleaned multi-view dataset.

**3D Human Generation.** Early 3D generative models based on textured mesh [17, 35, 59], point cloud [1, 86], and voxels [55, 56, 94] suffer from limited expressiveness while requiring 3D datasets. Recent GAN-based works [18, 25, 90, 91] utilize NeRF-like representations [8] to achieve directly 3D human generation using 2D images [14]. Recent diffusion-based works [6, 21, 24, 89] combine various 3D representation and human priors [3, 27, 42, 60] with SDS method [62] and achieve better quality. Avatar-
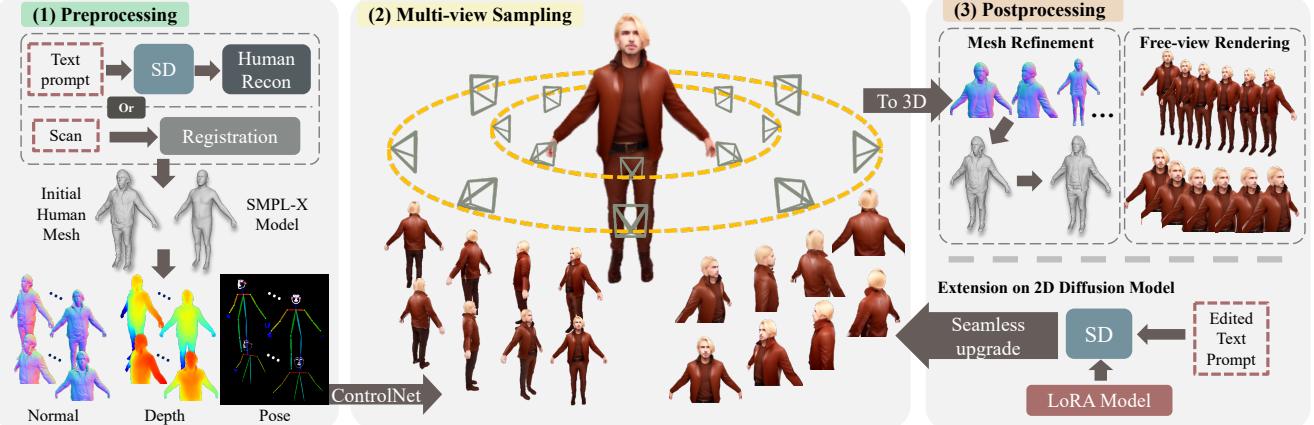
Figure 2. The overview of our MVHuman pipeline. Our method consists of three steps. The first preprocessing step obtains view-consistent conditions from the initial geometry and its aligned SMPL-X (Sec. 3.1). The second step performs the multi-view sampling process to generate view-consistent images (Sec. 3.2). The final postprocessing step includes geometry refinement and free-view rendering (Sec. 3.3).

Craft [24] leverages NeuS[40, 78] and the SMPL model to facilitate the generation of avatars. DreamHuman [31] uses imGHUM [3] to constrain a deformable NeRF for human. Similarly, DreamAvatar [6], DreamWaltz [21], and Avatar-Verse [89] utilize SMPL as shape prior. TADA [34] creates displacement and texture layers based on the SMPL-X model to generate 3D avatars. These SDS-based methods also inherit limitations that tend to generate oversaturated appearance, which is a severe artifact for human generation. **3D Human Reconstruction.** Many previous works [15, 23, 67, 68] achieve monocular human geometry reconstruction using implicit function, but they are unstable for novel human poses. ICON [83] and ECON [84] take the SMPL to guide explicit normal estimation and achieve a better trade-off between geometry detail and pose stability. Other works [33, 58, 61, 80, 93] combine NeRF techniques [5, 9, 28, 46, 47, 49, 51, 53, 79] with human shape prior to model both human geometry and appearance, while a neural texture blending scheme [26, 47, 93] has been demonstrated to enhance the realisim of appearance. In this work, we adopt ECON to provide shape prior for human image generation and produce high-fidelity human assets.

# 3. Method

Here, we introduce our text-guided human radiance fields generation scheme, MVHuman, which seeks to directly utilize the 2D generative capability of the existing latent diffusion model [66] with ControlNet [92] without extra fine-tuning or distilling. As illustrated in Fig. 2, the core of our MVHuman is a novel multi-view sampling process to simultaneously sample multiple view-consistent images with the aid of a coarse human geometry proxy (Sec. 3.1). Specifically, we construct "consistency-guided noise" in sampling steps to gradually denoise the individual initial random noises of multiple views into consistent ones (Sec. 3.2). With the multi-view sampling above, we

carefully generate high-quality human images from multiple view points which enable reconstructing detailed 3D geometry and generating neural radiance fields followed by a neural blending scheme for free-view rendering (Sec.3.3).

## 3.1. Preprocessing

We propose to generate a rough human mesh as a geometry proxy from scratch guided by a text prompt. Specifically, we adopt the stable diffusion model to generate a realistic full-body human front-view image and then utilize the off-the-shelf ECON [84] model to reconstruct full-body mesh and aligned SPML-X [60] model providing 3D pose. Then we project them to desired target views to depth, normal maps, and 2D skeletal poses, providing view-consistent conditions for the ControNet [92] in the following sampling process. Note that we also support pre-provided geometry, e.g., human scans, by utilizing the registration technique [4] to obtain the corresponding SMPL-X model.

## 3.2. Multi-view Sampling Process

Here we describe our multi-view sampling process to generate view-consistent human images for various views, e.g., views evenly placed on circular tracks shown in Fig. 2, with the aid of geometry proxy and controlling conditions. The key observation is, for simultaneous sampling processes with different initial random noises, if we blend their predicted original signals and transform it back to a novel noise to replace the original predicted ones at each time step, such sampling processes will finally recover consistent signals (refer to the supplemental materials for more details). We apply such modified noise to a multi-view setting called "consistency-guided noise" so as to constrain the visible parts from multiple viewpoints to be as consistent as possible. Besides, we propose a latent codes optimization scheme and apply feature concatenation to the self-attention blocks in the SD model to further enhance the view consis-

tency of sampled human images.

**Preliminary.** Stable diffusion [66] is a diffusion model that learns a denoising backward process in the latent space of an encoder $\mathcal{E}$ and a decoder $\mathcal{D}$. It adopts a UNet-like network $\epsilon$ conditioned on a text-prompt embedding to predict the noise $\epsilon_t$ at each time step $t$ in the denoising process, which should follow a normal distribution of $\mathcal{N}(\mathbf{0}, \mathbf{I})$. A deterministic sampling process such as DDIM [72] is applied to denoise latent $x_t$ to $x_{t-1}$ as:

$$x_{t-1} = \sqrt{\alpha_{t-1}}x_{0\leftarrow t} + \sqrt{1-\alpha_{t-1}}\epsilon_t, t = T, \ldots, 1, \quad (1)$$

where $\{\alpha_t\}_{t=1}^T$ are constants from DDIM [72] and $x_{0\leftarrow t}$ is the predicted original signal at step $t$:

$$x_{0\leftarrow t} = \left( \frac{x_t - \sqrt{1-\alpha_t}\epsilon_t}{\sqrt{\alpha_t}} \right). \quad (2)$$

We denote Eqn. 1 as: $x_{t-1} = \mathcal{S}(x_t, \epsilon_t)$. Note that we omit the conditional inputs, e.g., text, depth, for brevity.

**Consistency-guided Noise.** As shown in Fig. 3, we propose to construct consistency-guided noise $\epsilon_t'^i$ at sampling steps for each view $v_i$ in total $N$ views thus gradually denoise the latent codes $x_t^i$ to view-consistent ones as: $x_{t-1}^i = \mathcal{S}(x_t^i, \epsilon_t'^i)$. To this end, we first conduct a warping-based blending scheme to fuse the predicted original signals $x_{0\leftarrow t}^{j_1}, x_{0\leftarrow t}^{j_2}, \ldots$ of adjacent source views $v_{j_1}, v_{j_2}, \ldots$ to the target view $v_i$ using the initial human mesh. Specifically, we warp the decoded images from source views to target view and then encode them back to latent space rather than warping in latent space to avoid misalignment:

$$x_{0\leftarrow t}'^j = \mathcal{E}(W_j^i(\mathcal{D}(x_{0\leftarrow t}^j))) \quad (3)$$

where $W_j^i$ denotes the warping operation from $v_j$ to $v_i$. Since the warping operation will not modify the pixel values, we can assume that the transformed original signal $x_{0\leftarrow t}'^j$ still follows a normal distribution of $\mathcal{N}(\mu_j, \sigma^2)$ (refer to the supplemental materials for more details). Then we generate weighted occlusion maps $M_j^i$ from $v_j$ to $v_i$ according to the visibility shown in Fig. 3:

$$M_j^i(x, y) = \begin{cases} 0, \text{if } \boldsymbol{p} \text{ is invisible to } v_j \\ (\boldsymbol{n}(\boldsymbol{p}) \cdot \frac{\boldsymbol{o}(v_j) - \boldsymbol{p}}{\|(\boldsymbol{o}(v_j) - \boldsymbol{p})\|}) * w_s + w_c, \text{else} \end{cases} \quad (4)$$

where $\boldsymbol{p}$ is the intersection point of the ray casting from the pixel $(x, y)$ of $v_i$ and the human mesh, $\boldsymbol{n}(\boldsymbol{p})$ is the corresponding normal and $\boldsymbol{o}(v_j)$ is the camera position of $v_j$. $w_s$ is a scaling factor and $w_c$ is a constant to avoid numerical instability. Now we can get the weighted average $\bar{x}_{0\leftarrow t}'^i$ of these processed original signals for target view $v_i$:

$$\bar{x}_{0\leftarrow t}'^i = \frac{M_i^i}{M_{sum}}x_{0\leftarrow t}'^i + \frac{M_{j_1}^i}{M_{sum}}x_{0\leftarrow t}'^{j_1} + \ldots, \quad (5)$$
$$M_{sum} = M_i^i + M_{j_1}^i + \ldots.$$

Here the $\bar{x}_{0\leftarrow t}'^i$ should follow the distribution of $\mathcal{N}(\frac{M_i^i}{M_{sum}}\mu_i + \frac{M_{j_1}^i}{M_{sum}}\mu_{j_1} + \ldots, \frac{M_i^{i\,2}+M_i^{j_1\,2}+\ldots}{M_{sum}^2}\sigma^2)$. Ideally we should scale the noise part of $x_{0\leftarrow t}'^i, x_{0\leftarrow t}'^{j_1}, \ldots$ with factor $\frac{M_{sum}}{\sqrt{M_i^{i\,2}+M_i^{j_1\,2}+\ldots}}$ to maintain the variance unchanged as $\sigma^2$. However, $\mu_i, \mu_{j_1}, \ldots$ and $\sigma^2$ cannot be easily measured because the decode-then-encode operation cannot guarantee the fidelity. Thus we empirically use the following formula to blend a new weighted average $\tilde{x}_{0\leftarrow t}'^i$ shown in Fig. 3:

$$\tilde{x}_{0\leftarrow t}'^i = \sum_{k=i,j_1,\ldots} \frac{M_k^i}{M_{sum}}(E x_{0\leftarrow t}'^k + (1-E)\bar{x}_{0\leftarrow t}'^k), \quad (6)$$

where $E = \frac{M_{sum}}{\sqrt{M_i^{i\,2}+M_i^{j_1\,2}+\ldots}}$.

Finally we can get the consistency-guided noise using Eqn. 2: $\epsilon_t'^i = \frac{x_t^i - \sqrt{\alpha_t}\tilde{x}_{0\leftarrow t}'^i}{\sqrt{1-\alpha_t}}$ which can be used to sample the latent of step $t-1$ for view $v_i$: $x_{t-1}^i = \mathcal{S}(x_t^i, \epsilon_t'^i)$.

**Optimization of Latent Codes.** We further optimize the latent codes to enforce view consistency across different views. As shown in Fig. 3, for a target view $v_i$ and a source view $v_j$, we decode the latent code $x_t^j$ of $v_j$ to image space and then warp the image to $v_i$, and then apply L2 loss on the covered pixels in image space. We also conduct the same operation in reverse to formulate the following loss to optimize both $x_t^i$ and $x_t^j$:

$$\mathcal{L}_{latent} = \|\mathcal{D}(x_t^i) - W_j^i(\mathcal{D}(x_t^j))\|_2^2 + \|\mathcal{D}(x_t^j) - W_i^j(\mathcal{D}(x_t^i))\|_2^2. \quad (7)$$

**Self-attention Layer Modification.** To further enhance the view similarity between adjacent views, we also import consistency prior to views following recent video diffusion methods [7, 29] by modifying the self-attention blocks in the SD model. More specifically, for a view $v_i$, the input features to the self-attention block are query $Q^i$, key $K^i$, and value $V^i$. We extend the block to concatenate value and key from a predefined reference view, e.g., facing view, with the original ones:

$$\text{Ex-Attn}\left(Q^i, [K^{ref}, K^i], [V^{ref}, V^i]\right) =$$
$$\text{Softmax}\left( \frac{Q^i\left([K^{ref}, K^i]\right)^T}{\sqrt{c}} \right)[V^{ref}, V^i]. \quad (8)$$

**Human Images Generation.** As depicted in Fig.3, we configure two concentric circular tracks of viewpoints, each uniformly distributed with $N$ cameras (we set $N = 8$). The first track of views looks at the full body, denoted as $R_{fb} = \{v_i | i = 1, \ldots, N\}$ and the second focuses on the upper body only, denoted as $R_{ub} = \{v_i | i = N+1, \ldots, 2N\}$. For each view $v_i$ on $R_{ub}$ as the target view, its neighboring views on the same track within a range of $60°$ clockwise and counterclockwise are used as its source views. For
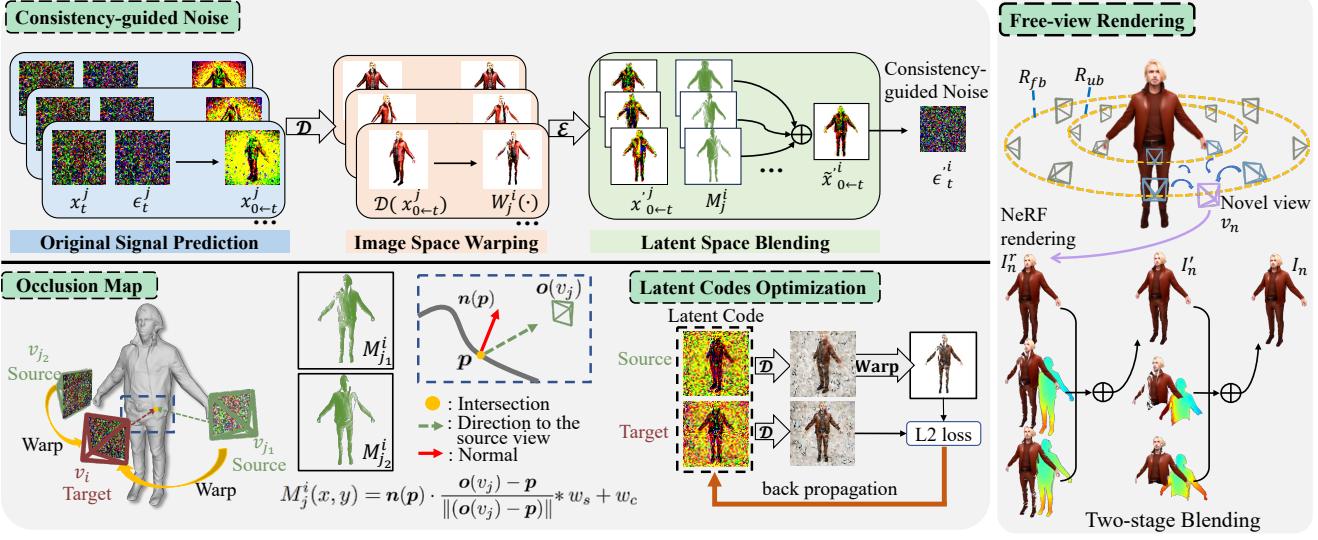
Figure 3. Illustration of our method. The left-top area illustrates the construction of the consistency-guided noise for each view with its adjacent source views at a sampling step. The left-down area illustrates the generation of weighted occlusion maps. The mid-down area illustrates the optimization of latent codes. And the right area illustrates the adopted two-stage blending scheme on human radiance fields for free-view rendering.

each view $v_i$ on $R_{fb}$ as the target view, in addition to its neighboring views on the same track, we set the view $v_j$ $(j = i+N)$ from $R_{ub}$ with the same azimuth angle as an additional source view. We replace the $\tilde{x}'^i_{0 \leftarrow t}$ with $W^i_j(x'^j_{0 \leftarrow t})$ in the region where $M^i_j > 0$ to maintain the distinguishable details, e.g., facial details, from the closeup views. Throughout the multi-view sampling process, we alternate between leveraging predicted noise and consistency-guided noise. To be precise, we sample with $\epsilon^i_t$ when $t$ is odd and $\epsilon'^i_t$ when $t$ is even. We find it helps enhance the fidelity of the generated details.

For the latent optimization, we initially compute the loss following Eqn. 7 for each pair of adjacent views on the same track $R_{fb}$ or $R_{ub}$. To improve generated details, we lock the latent codes from the front-view and back-view and further lock those from the side-views upon completion of 20% of the sampling process. Subsequently, we optimize between the two tracks, for each view $v_i$ on $R_{fb}$, we compute the loss with the locked latent code from $v_{i+N}$ on $R_{ub}$. More details are available in the supplementary material.

### 3.3. Postprocessing for Human Assets

**Geometry Optimization.** In order to facilitate the generation of 3D human assets, we leverage the implicit geometric information within the human images to refine the initial mesh. For the obtained desired high-quality human images as $\{I_i | i = 1, \ldots, 2N\}$ from the multi-view sampling process, we further utilize the human normal predictor from ECON to predict corresponding normals, denoted as $\{\mathcal{N}^p_i | i = 1, \ldots, 2N\}$. We then utilize a differentiable rasterizer [32, 82] to render the normal map of the initial mesh to each view, denoted as $\{\mathcal{N}^m_i | i = 1, \ldots, 2N\}$. Sub-

sequently, we optimize the position of vertices to minimize the difference between the predicted normal maps and rendered ones [30], however, the predicted normals might not conform strictly to the value range typical of an authentic normal map. Thus, we approximate gradients of normal maps using the Sobel operator $G$ and conduct L2 loss in gradient space as:

$$\mathcal{L}_{normal} = \sum_{i=1}^{2N} \|G(\mathcal{N}^m_i) - G(\mathcal{N}^p_i)\|^2_2. \quad (9)$$

**Free-view Rendering.** With multi-view images, we are already able to train an efficient neural radiance field utilizing a multi-plane representation [5]. However, the rendering quality has yet to reach the level of fidelity exhibited by the images themselves. Thus we further implement a neural blending strategy characterized by a two-stage, coarse-to-fine approach that initially integrates the full-body views followed by the inclusion of the upper-body views. Following [73], we train a UNet-like network $\Theta$ that takes warped RGB and depth disparity maps as input and outputs two blending weights: $[W_1, W_2] = \Theta(\hat{I}_1, \hat{I}_2, O_1, O_2)$. As shown in Fig. 3, given a novel view $v_n$, we first obtain the NeRF rendering result $I^r_n$ at $v_n$ and the rendered depth map $D_n$ from mesh. In the first stage, we find two nearest views to $v_n$ from full body views $R_{fb}$, denoted as $v_{j_1}, v_{j_2}$. We then warp their images to the novel view and obtain the depth disparity maps from the warped depth and $D_n$ as inputs to $\Theta$: $[W_{j_1}, W_{j_2}] = \Theta(\hat{I}_{j_1}, \hat{I}_{j_2}, O_{j_1}, O_{j_2})$. The blended result of the first stage is denoted as $I'_n$. In the second stage, we further find two nearest views to $v_n$ from upper body views $R_{ub}$, denoted as $v_{j_3}, v_{j_4}$. Employ-

5

Figure 4. The texture and geometry results of MVHuman, including characters from games/movies, celebrities, and customized humans.

ing a similar procedure, we calculate their blending weights $W_{j_3}, W_{j_4}$. To this end, we get the final blended result $I_n$ as:

$$I_n^{'} = W_{j_1} \odot \hat{I}_{j_1} + W_{j_2} \odot \hat{I}_{j_2} + (1 - W_{j_1} - W_{j_2}) \odot I_n^r$$
$$I_n = W_{j_3} \odot \hat{I}_{j_3} + W_{j_4} \odot \hat{I}_{j_4} + (1 - W_{j_3} - W_{j_4}) \odot I_n^{'}$$
$$(10)$$

## 4. Experimental Results

In this section, we first demonstrate the capability of our MVHuman, and then compare our method against the state-of-the-art 3D human generation methods by conducting a user study. We further evaluate each component of our method on view consistency and quality with quantitative and qualitative results.

### 4.1. Comparison

We compare our method with state-of-the-art text-guided 3D generation methods, TEXTure [65], Fantasia3D [12], DreamWaltz [21], MVDream [71], AvatarVerse [89] and TADA [34]. As illustrated in Fig. 5, TEXTure generates texture with ghosting and suffers from Janus Failure. Fantasia3D fails to generate complete human body geometry. MVDream generates reasonable geometries but its results lack detail, especially in the head region. DreamWaltz generates blurry texture with unclear edges. TADA takes more than 3 hours for each prompt, and it fails to generate geometries with correct proportions. The authors of AvatarVerse provide us with the cases for qualitative comparison. Their method generates detailed and realistic geometry, but the color of its texture is oversaturated.

We also conduct a comprehensive user study to evaluate the performance of our generated human assets in terms of whether they match the given prompts and their overall quality. 30 prompts describing well-known people and characters are designed and used to generate and render 360-degree videos with full-body and upper-body views. We randomly select 10 samples for each user and ask them to comprehensively consider and select the best results for each sample from four aspects: conforming to the prompt text, geometric and texture quality, face details, and body proportions. As illustrated in Fig. 6, 26 volunteers majoring in computer vision completed the survey, and our method shows significant advantages in the result. Please refer to the supplemental materials for more comparisons.

### 4.2. Ablation Study

**Multi-view Sampling.** We evaluate the multi-view sampling process. In Fig. 7, first we condition the SD with depth, normal map and Openpose image of each view, and conduct a vanilla sampling process (**+ conditions**). The generated results align with the conditions with no view consistency. Then we incorporate the consistency-guided noise (**+ C-G noise**), and the results show better consistency, yet there remain some blending artifacts in the textures. When we incorporate the optimization of latent codes (**+ optim**), the blending artifacts are alleviated. Finally, after the self-attention block is modified (**+ attn**), our full method achieves the best visual effect. We further quantitatively evaluate the view consistency. For each view, we warp the generated image to its neighboring views to compute the PSNR value on the overlapped region. As shown in Tab. 1, our full method achieves the best score.

**Number of Views.** We evaluate the number of views $N$ on each track $R_{fb}, R_{ub}$ qualitatively. As shown in Fig. 8, when $N = 6$, insufficient views can lead to artifacts on NeRF

Figure 5. Qualitative comparison between TEXTure, Fantasia3D, DreamWaltz, MVDream, AvatarVerse, TADA and our MVHuman. Our method balances geometric quality and appearance details while avoiding oversaturation and Janus Failure. The prompts from (a) to (c) are *"Ian McKellen, Magneto"*, *"Gal Gadot, Wonder Woman"*, *"Lara Croft, Tomb Raider"*. (* Our second is a case using scan mesh.)



Figure 6. Result of user study over Fantasia3D, TEXTure, MV-Dream, DreamWaltz, TADA, and our method.

training and final blending results. When $N = 12$, the running time can cost more than one hour, and some blurry artifacts may occur due to the imbalanced optimization of latent codes between views. Our setting $N = 8$ ensures stable generation while keeping high quality and efficiency. **Blending Scheme.** We evaluate our blending scheme qualitatively. As shown in Fig. 9, the directly generated full-body image (**w/o blending**) is a little inconsistent with results from the closeup view, especially on the face region. Our blending scheme (**w blending**) helps improve the generation quality at full-body views. More evaluations are detailed in the supplementation.

### 4.3. Application

**Text-guided Editing.** As one of the inherent capabilities of the SD model, text-guided editing can be seamlessly inte-



(a) + conditions  (b) + C-G noise  (c) + optim  (d) + attn (full)

Figure 7. Qualitative evaluation of the multi-view sampling. Note the change in the chest area. (*"Loki, green and gold armor"*)

Table 1. Quantitative evaluation of view consistency.

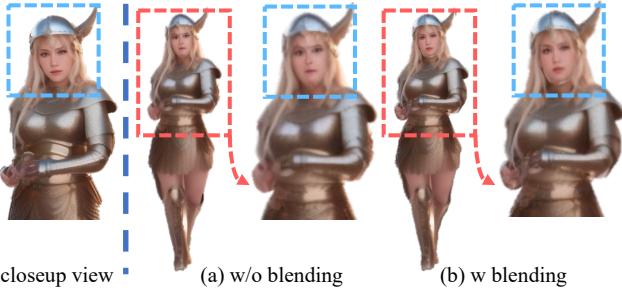| Method | PSNR ↑ | SSIM ↑ |
|---|---|---|
| + conditions | 19.297 | 0.9092 |
| + C-G noise | 28.010 | 0.9678 |
| + optimization | 33.444 | 0.9851 |
| full | 34.074 | 0.9860 |

grated into our MVHuman. As shown in Fig. 10, by modifying the textual description, we can manipulate the color of a character's clothing or hair.

**Style Transfer with LoRA.** Our MVHuman is based on a pre-trained 2D SD model, which leverages an extensive repository of LoRA models. This allows us to seamlessly integrate LoRA models into our method, facilitating tasks

(a) 6 views            (b) 12 views            (c) 8 views

Figure 8. Qualitative evaluation of the number of views on each track. (*"The Flash"* & *"Harrison Ford"*)



closeup view        (a) w/o blending            (b) w blending

Figure 9. Qualitative evaluation of the neural blending scheme. (*"Valkyrie, in chain mail, skirt armor, and helmet"*)



Figure 10. Application: enabling text-guided human asset editing.

such as style transfer. As shown in Fig. 11, with two LoRA models loaded to the SD model, the style of the generated results has changed respectively.

## 4.4. Discussion

**Limitations.** Although MVDream shows promising human generation results, it still has some limitations. First, our method relies on the accuracy of the alignment between the initial coarse mesh and SMPL-X model. It is meaningful to find a way to efficiently obtain these prerequisites. Second, textual descriptions often encounter challenges in articulating specific details. This problem may benefit from a combination with image-conditioned diffusion models. Moreover, our method does not support the relighting tasks because it does not take albedo, material into consideration. Thus it is meaningful to have a diffusion model that can generate albedo and various materials given conditions such



Figure 11. Application: style transfer with LoRA models. (*"Leo Tolstoy, a writer"* & *"A woman wearing ski clothes"*)

as semantic segmentation. Also, it could be interesting to combine our method with other 3D representations to better deal with specific generations like hair.

**Social Impact.** When researching generative technologies, we concern about their potential infringements on intellectual property. There should be heightened legal restrictions on the applications of such technologies. Furthermore, gender and cultural diversity are crucial. It is necessary for any generative technology to ensure inclusivity and avoid stereotypes. In this paper, all our results are carefully selected based on these principles.

## 5. Conclusion

We have presented a novel approach to generate human radiance fields from text prompts, with consistent multi-view images directly sampled from pre-trained 2D diffusion models without any fine-tuning or distilling. The core of our approach is a multi-view sampling process compatible with the 2D diffusion model, which carefully tailors the input conditions, predicted noises and the corresponding latent codes. It enables consistent multi-view generation by producing the "consistency-guided noise" to replace the original predicted noise, explicitly optimizing the latent codes of adjacent views, and modifying the self-attention layer of the pre-trained SD network. Such generated images further benefit producing exquisite human assets with refined geometry and free-view rendering. Our experimental results demonstrate the effectiveness of our approach and we also showcase various 3D applications by seamlessly lifting many editing functions of the original 2D diffusion models to 3D. With the above characteristics, we believe that our approach is a critical step towards high-quality 3D human generation, providing a new understanding of the relationship between 2D and 3D generation.

# References

[1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and generative models for 3d point clouds. In *International conference on machine learning*, pages 40–49. PMLR, 2018. 2

[2] Oleg Alexander, Mike Rogers, William Lambeth, Matt Chiang, and Paul Debevec. The digital emily project: photoreal facial modeling and animation. In *Acm siggraph 2009 courses*, pages 1–15. 2009. 1

[3] Thiemo Alldieck, Hongyi Xu, and Cristian Sminchisescu. imghum: Implicit generative models of 3d human shape and articulated pose. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5441–5450, 2021. 2, 3

[4] Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Loopreg: Self-supervised learning of implicit surface correspondences, pose and shape for 3d human mesh registration. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 3

[5] Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. *CVPR*, 2023. 2, 3, 5

[6] Yukang Cao, Yan-Pei Cao, Kai Han, Ying Shan, and Kwan-Yee K Wong. Dreamavatar: Text-and-shape guided 3d human avatar generation via diffusion models. *arXiv preprint arXiv:2304.00916*, 2023. 1, 2, 3

[7] Duygu Ceylan, Chun-Hao P Huang, and Niloy J Mitra. Pix2video: Video editing using image diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23206–23217, 2023. 2, 4

[8] 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. 2

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

[10] Dave Zhenyu Chen, Yawar Siddiqui, Hsin-Ying Lee, Sergey Tulyakov, and Matthias Nießner. Text2tex: Text-driven texture synthesis via diffusion models. *arXiv preprint arXiv:2303.11396*, 2023. 2

[11] Hansheng Chen, Jiatao Gu, Anpei Chen, Wei Tian, Zhuowen Tu, Lingjie Liu, and Hao Su. Single-stage diffusion nerf: A unified approach to 3d generation and reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2416–2425, 2023. 1

[12] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 1, 2, 6

[13] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of*

[14] 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. 2

[15] 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

[16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1, 2

[17] 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

[18] Fangzhou Hong, Zhaoxi Chen, Yushi LAN, Liang Pan, and Ziwei Liu. EVA3d: Compositional 3d human generation from 2d image collections. In *International Conference on Learning Representations*, 2023. 2

[19] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. 2

[20] Xin Huang, Ruizhi Shao, Qi Zhang, Hongwen Zhang, Ying Feng, Yebin Liu, and Qing Wang. Humannorm: Learning normal diffusion model for high-quality and realistic 3d human generation. *arXiv preprint arXiv:2310.01406*, 2023. 3

[21] Yukun Huang, Jianan Wang, Ailing Zeng, He Cao, Xianbiao Qi, Yukai Shi, Zheng-Jun Zha, and Lei Zhang. Dreamwaltz: Make a scene with complex 3d animatable avatars. *arXiv preprint arXiv:2305.12529*, 2023. 1, 2, 3, 6

[22] Yangyi Huang, Hongwei Yi, Yuliang Xiu, Tingting Liao, Jiaxiang Tang, Deng Cai, and Justus Thies. Tech: Text-guided reconstruction of lifelike clothed humans. *arXiv preprint arXiv:2308.08545*, 2023. 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] Ruixiang Jiang, Can Wang, Jingbo Zhang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Avatarcraft: Transforming text into neural human avatars with parameterized shape and pose control, 2023. 2, 3

[25] Suyi Jiang, Haoran Jiang, Ziyu Wang, Haimin Luo, Wenzheng Chen, and Lan Xu. Humangen: Generating human radiance fields with explicit priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12543–12554, 2023. 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*

the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 13142–13153, 2023. 1

*on Computer Vision and Pattern Recognition*, pages 6155–6165, 2022. 2, 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

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

[29] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. *arXiv preprint arXiv:2303.13439*, 2023. 2, 4

[30] Byungjun Kim, Patrick Kwon, Kwangho Lee, Myunggi Lee, Sookwan Han, Daesik Kim, and Hanbyul Joo. Chupa: Carving 3d clothed humans from skinned shape priors using 2d diffusion probabilistic models, 2023. 5

[31] Nikos Kolotouros, Thiemo Alldieck, Andrei Zanfir, Eduard Gabriel Bazavan, Mihai Fieraru, and Cristian Sminchisescu. Dreamhuman: Animatable 3d avatars from text. *arXiv preprint arXiv:2306.09329*, 2023. 1, 3

[32] Samuli Laine, Janne Hellsten, Tero Karras, Yeongho Seol, Jaakko Lehtinen, and Timo Aila. Modular primitives for high-performance differentiable rendering. *ACM Transactions on Graphics*, 39(6), 2020. 2, 5

[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] Tingting Liao, Hongwei Yi, Yuliang Xiu, Jiaxiang Tang, Yangyi Huang, Justus Thies, and Michael J Black. Tada! text to animatable digital avatars. *arXiv preprint arXiv:2308.10899*, 2023. 1, 3, 6

[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] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1, 2

[37] Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization, 2023. 2

[38] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9298–9309, 2023. 1, 2

[39] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453*, 2023. 2, 3

[40] Xiaoxiao Long, Cheng Lin, Peng Wang, Taku Komura, and Wenping Wang. Sparseus: Fast generalizable neural surface reconstruction from sparse views. *ECCV*, 2022. 3

[41] Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin Ma, Song-Hai Zhang, Marc Habermann, Christian Theobalt, and Wenping Wang. Wonder3d: Single image to 3d using cross-domain diffusion, 2023. 1, 2, 3

[42] 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, 2015. 2

[43] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 851–866. 2023. 1

[44] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022. 2

[45] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022. 2

[46] 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 Photography (ICCP)*, pages 1–12, Los Alamitos, CA, USA, 2021. IEEE Computer Society. 3

[47] 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), 2022. 3

[48] Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. Realfusion: 360deg reconstruction of any object from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8446–8455, 2023. 2

[49] 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

[50] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for shape-guided generation of 3d shapes and textures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12663–12673, 2023. 2

[51] 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

[52] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view syn-

thesis. In *Computer Vision – ECCV 2020*, pages 405–421, Cham, 2020. Springer International Publishing. 1, 2

[53] 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, 2022. 3

[54] Jacob Munkberg, Jon Hasselgren, Tianchang Shen, Jun Gao, Wenzheng Chen, Alex Evans, Thomas Müller, and Sanja Fidler. Extracting Triangular 3D Models, Materials, and Lighting From Images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8280–8290, 2022. 2

[55] 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

[56] 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

[57] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-e: A system for generating 3d point clouds from complex prompts. *arXiv preprint arXiv:2212.08751*, 2022. 1

[58] 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

[59] 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

[60] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3d hands, face, and body from a single image. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 10975–10985, 2019. 2, 3

[61] 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

[62] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *ICLR*, 2023. 1, 2

[63] Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skorokhodov, Peter Wonka, Sergey Tulyakov, and Bernard Ghanem. Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. *arXiv preprint arXiv:2306.17843*, 2023. 2

[64] Amit Raj, Srinivas Kaza, Ben Poole, Michael Niemeyer, Nataniel Ruiz, Ben Mildenhall, Shiran Zada, Kfir Aberman, Michael Rubinstein, Jonathan Barron, et al. Dreambooth3d: Subject-driven text-to-3d generation. *arXiv preprint arXiv:2303.13508*, 2023. 2

[65] Elad Richardson, Gal Metzer, Yuval Alaluf, Raja Giryes, and Daniel Cohen-Or. Texture: Text-guided texturing of 3d shapes. In *ACM SIGGRAPH 2023 Conference Proceedings*, New York, NY, USA, 2023. Association for Computing Machinery. 2, 6

[66] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 2, 3, 4

[67] 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)*, 2019. 3

[68] 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)*, 2020. 3

[69] Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2

[70] Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu, Chao Xu, Xinyue Wei, Linghao Chen, Chong Zeng, and Hao Su. Zero123++: a single image to consistent multi-view diffusion base model, 2023. 2

[71] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023. 1, 2, 6

[72] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv:2010.02502*, 2020. 2, 4

[73] 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. 2, 5

[74] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen. Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 22819–22829, 2023. 2

[75] Shitao Tang, Fuyang Zhang, Jiacheng Chen, Peng Wang, and Yasutaka Furukawa. Mvdiffusion: Enabling holistic multi-view image generation with correspondence-aware diffusion. *arXiv preprint arXiv:2307.01097*, 2023. 2

[76] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A. Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12619–12629, 2023. 2

11

[77] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12619–12629, 2023. 1

[78] 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. 3

[79] 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

[80] 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

[81] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *arXiv preprint arXiv:2305.16213*, 2023. 1, 2

[82] Markus Worchel, Rodrigo Diaz, Weiwen Hu, Oliver Schreer, Ingo Feldmann, and Peter Eisert. Multi-view mesh reconstruction with neural deferred shading. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6187–6197, 2022. 5

[83] Yuliang Xiu, Jinlong Yang, Dimitrios Tzionas, and Michael J. Black. ICON: Implicit Clothed humans Obtained from Normals. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13296–13306, 2022. 3

[84] Yuliang Xiu, Jinlong Yang, Xu Cao, Dimitrios Tzionas, and Michael J. Black. ECON: Explicit Clothed humans Optimized via Normal integration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 3

[85] Dejia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Yi Wang, and Zhangyang Wang. Neurallift-360: Lifting an in-the-wild 2d photo to a 3d object with 360deg views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4479–4489, 2023. 2

[86] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4541–4550, 2019. 2

[87] Xianggang Yu, Mutian Xu, Yidan Zhang, Haolin Liu, Chongjie Ye, Yushuang Wu, Zizheng Yan, Chenming Zhu, Zhangyang Xiong, Tianyou Liang, et al. Mvimgnet: A large-scale dataset of multi-view images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9150–9161, 2023. 1

[88] Bohan Zeng, Shanglin Li, Yutang Feng, Hong Li, Sicheng Gao, Jiaming Liu, Huaxia Li, Xu Tang, Jianzhuang Liu, and Baochang Zhang. Ipdreamer: Appearance-controllable 3d object generation with image prompts. *arXiv preprint arXiv:2310.05375*, 2023. 2

[89] Huichao Zhang, Bowen Chen, Hao Yang, Liao Qu, Xu Wang, Li Chen, Chao Long, Feida Zhu, Kang Du, and Min Zheng. Avatarverse: High-quality & stable 3d avatar creation from text and pose. *arXiv preprint arXiv:2308.03610*, 2023. 1, 2, 3, 6

[90] 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. In *European Conference on Computer Vision*, pages 668–685. Springer, 2022. 2

[91] 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

[92] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 1, 2, 3

[93] 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

[94] 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. 2

# MVHuman: Tailoring 2D Diffusion with Multi-view Sampling For Realistic 3D Human Generation

## Supplementary Material

In this supplementary material, we first provide additional analysis on the multi-view sampling process in Sec. 6. Then we provide details of our implementation in Sec. 7. More evaluations are discussed in Sec. 8 and Sec. 9. We further provide more results in Sec. 10.

## 6. Additional Analysis on Multi-view Sampling

In this section, we discuss more about the consistency-guided noise of the multi-view sampling process. In sec. 6.1, we showcase a special scenario of multi-view sampling in 2D in which all views are identical. We demonstrate that, with consistency-guided noise, for sampling processes with different initial random noises, they can finally recover consistent results. In sec. 6.2, we further analyze how such multi-view sampling process can be transferred into a very different 3D scenario to support generating multi-view consistent images in 3D.

### 6.1. Multi-view Sampling in 2D Scenario

In stable diffusion models, for multiple sampling processes with different initial random noises, if we denoise them with their own vanilla predicted noise, the final results will be various because these sampling processes are individual. Yet, we find that we can blend their predicted original signals at sampling timesteps, so as to obtain a consistent prediction and transform it back to the corresponding consistency-guided noises. By replacing the original predicted noises with such consistency-guided ones, various sampling processes can finally recover a consistent signal.

Specifically, assuming we have $N$ identical views, which means there are no warping operations, the scenario would degrade to 2D. At timestep $t$, we have their latent codes: $\{x_t^k | k = 1, \ldots, N\}$ and corresponding predicted noises: $\{\epsilon_t^k | k = 1, \ldots, N\}$. These noises follow a normal distribution of $\mathcal{N}(\mathbf{0}, \mathbf{I})$. With following equation:

$$x_{0\leftarrow t}^k = \frac{x_t^k - \sqrt{1-\alpha_t}\epsilon_t^k}{\sqrt{\alpha_t}}, \tag{11}$$

we can obtain corresponding predicted original signals at step t: $\{x_{0\leftarrow t}^k | k = 1, \ldots, N\}$. For each $x_{0\leftarrow t}^k$, it should follow the distribution of $\mathcal{N}(\frac{x_t^k}{\sqrt{\alpha_t}}, \frac{1-\alpha_t}{\alpha_t}\mathbf{I})$ respectively. Note that these predicted original signals have various means but share the same variance. Since there is no warping operation, we can directly calculate the average result $\bar{x}_{0\leftarrow t}$ following:

$$\bar{x}_{0\leftarrow t} = \frac{1}{N}(x_{0\leftarrow t}^1 + \ldots + x_{0\leftarrow t}^N) \tag{12}$$

which should follow the following distribution:

$$\bar{x}_{0\leftarrow t} \sim \mathcal{N}(\frac{1}{N\sqrt{\alpha_t}}(x_t^1 + \ldots + x_t^N), \frac{1 \times N}{N^2}\frac{1-\alpha_t}{\alpha_t}\mathbf{I})$$
$$\sim \mathcal{N}(\frac{x_t^1 + \ldots + x_t^N}{N\sqrt{\alpha_t}}, \frac{1}{N}\frac{1-\alpha_t}{\alpha_t}\mathbf{I}). \tag{13}$$

In order to maintain sampling quality, we need to maintain the variance of the distribution unchanged as $\frac{1-\alpha_t}{\alpha_t}\mathbf{I}$. Note that for each $x_{0\leftarrow t}^k$, its mean is known and related to $x_t^k$. Thus we can easily scale its noise part by $\sqrt{N}$ so that each new $x_{0\leftarrow t}^k$ follows a distribution of $\mathcal{N}(\frac{x_t^k}{\sqrt{\alpha_t}}, N\frac{1-\alpha_t}{\alpha_t}\mathbf{I})$. Then with Eqn. 12, we can get a new average result $\tilde{x}_{0\leftarrow t}$:

$$\tilde{x}_{0\leftarrow t} \sim \mathcal{N}(\frac{1}{N\sqrt{\alpha_t}}(x_t^1 + \ldots + x_t^N), \frac{1 \times N}{N^2}N\frac{1-\alpha_t}{\alpha_t}\mathbf{I})$$
$$\sim \mathcal{N}(\frac{x_t^1 + \ldots + x_t^N}{N\sqrt{\alpha_t}}, \frac{1-\alpha_t}{\alpha_t}\mathbf{I}). \tag{14}$$

With $\tilde{x}_{0\leftarrow t}$ and $x_t^k$, we can then obtain $N$ new consistency-guided noises with Eqn. 11, denoted as $\epsilon_t^{'k}$, $k = 1, \ldots, N$:

$$\epsilon_t^{'k} = \frac{x_t^k - \sqrt{\alpha_t}\tilde{x}_{0\leftarrow t}}{\sqrt{1-\alpha_t}}$$
$$\sim \mathcal{N}(\frac{1}{\sqrt{1-\alpha_t}}(x_t^k - \frac{x_t^1 + \ldots + x_t^N}{N}), \mathbf{I}). \tag{15}$$

We can observe that the mean of each consistency-guided noise $\epsilon_t^{'k}$ is not $\mathbf{0}$ and thus provides a direction towards a common target for each sampling process.

As shown in Fig. 12(a), with vanilla predicted noise for sampling, several simultaneous sampling processes with different initial noise will generate various individual results. In Fig. 12(b), with the consistency-guided noise for sampling, these sampling processes can finally recover high-quality consistent results.

### 6.2. Multi-view Sampling in 3D Scenario

With the multi-view sampling in 2D, we then aim to lift it into the 3D scenario so as to generate multi-view consistent images, which can further help reconstruct 3D results. In 3D scenario, the calculation of consistency-guided noise differs in three aspects:

1) The overlapped region differs between each pair of views. Thus, for each view as a target view, it only blends its predicted original signal with that from neighboring views
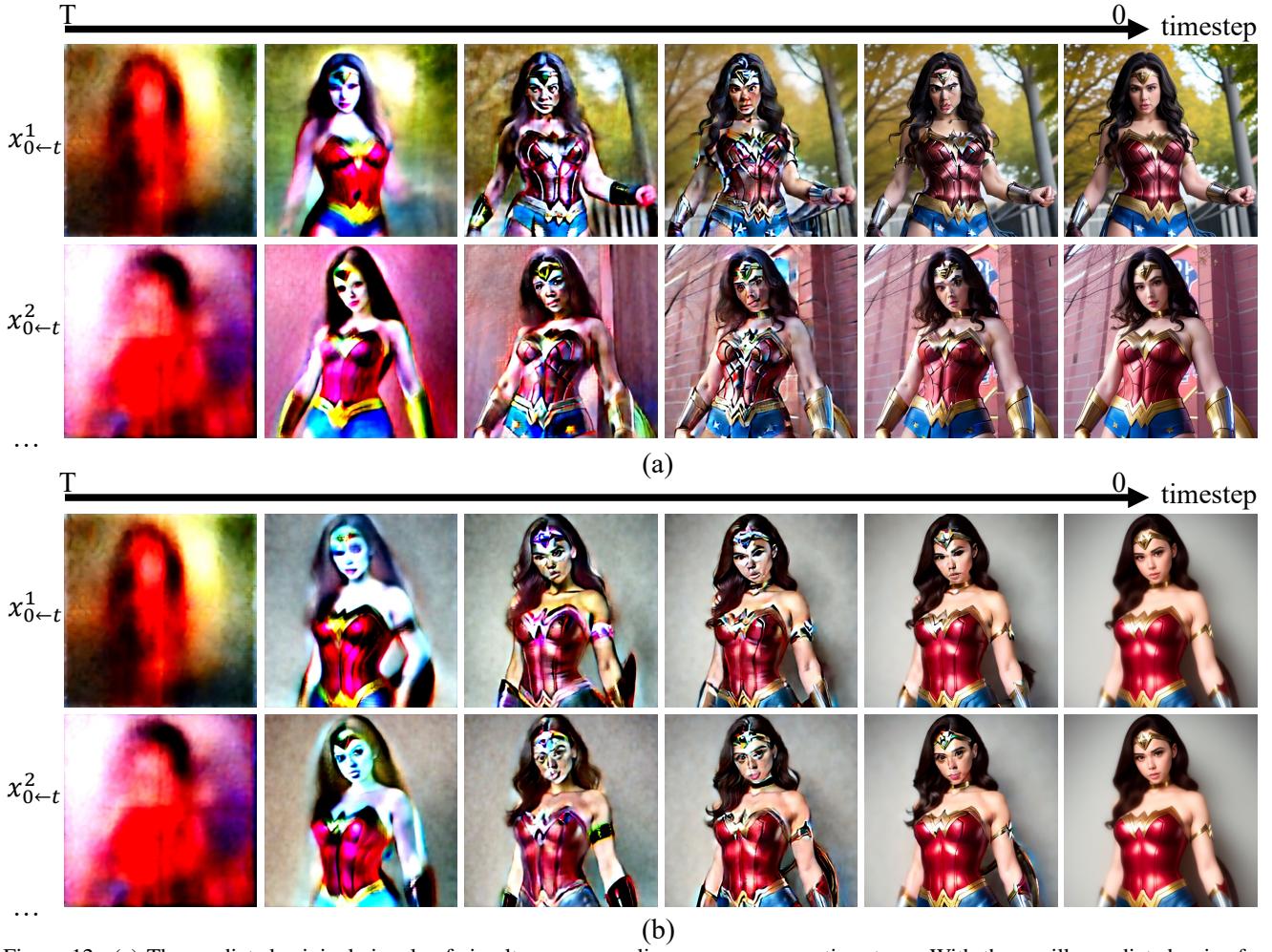
Figure 12. (a) The predicted original signals of simultaneous sampling processes over timesteps. With the vanilla predicted noise for sampling, several simultaneous sampling processes generate various outputs. (b) With the consistency-guided noise for sampling, several simultaneous sampling processes can recover consistent results.

that have common overlapped regions to obtain a consistent prediction.

2) In 3D, we cannot directly warp latent codes and predicted original signals between views because view consistency in image space is not equivalent to that in latent space. Also, it is $64 \times 64$ low resolution in latent space, and warping operations in such low resolution can lead to misalignment. In order to obtain view consistency in image space, we transform the predicted original signals from source views $v_j, j = j_1, j_2, \ldots$ to target view $v_i$ with the following decode-warp-encode operation:

$$x_{0\leftarrow t}^{'j} = \mathcal{E}(W_j^i(\mathcal{D}(x_{0\leftarrow t}^j))). \qquad (16)$$

3) Notably, the decode-warp-encode operation will somehow change the distributions of the predicted original signals $x_{0\leftarrow t}^j$, making it not easy to measure the means of $x_{0\leftarrow t}^{'j}$. Since $x_{0\leftarrow t}^j$ have the same variance $\frac{1-\alpha_t}{\alpha_t}\mathbf{I}$, we can still assume $x_{0\leftarrow t}^{'j}$ follow normal distributions of the same

variance $\sigma^2$ because the warping operation will not modify pixel value. Based on such assumption, we design an empirical formula to enhance the variance of the distribution of the weighted average which is then used to calculate consistency-guided noise.

## 7. Implementation Details

We test our MVHuman using NVIDIA RTX 3090 GPUs. The ControlNet models we use are the official version of v1.1. The stable diffusion models used are the official version of v1.5 and the open-source ChilloutMix model which is a fine-tuned SD v1.5 model with a better ability to generate human images. The sampler tested are DDIM [72] and Dpm-solver++ [44], and the number of sampling timesteps is set as 150. For all the results, we generate images of $512 \times 512$ resolution. The parameters in Eqn. 4 are $w_s = 0.2$, $w_c = 1.0$. Upon completion of $0\% - 20\%$ of the multi-view sampling process, we initiate the adoption

2

Figure 13. Qualitative comparison with DreamAvatar, DreamHuman, HumanNorm. Our results demonstrate photo-realistic appearance while avoiding oversaturation and Janus failure. The prompts of the first column from top to down are *"Doctor Strange"*(* case using scan mesh), *"Superman"*, *"Ronald Weasley"*. The prompts of the second column from top to down are *"a black female surgeon"*, *"a man wearing a bomber jacket"*, *"a cowboy"*. The prompts of the third column from top to down are *"Donald Trump"*, *"Joe Biden"*, *"Leonardo DiCaprio in a maroon long sleeve top"*.

of consistency-guided noise for sampling at timesteps (abbreviated as a C-G step), but the original predicted noise is still adopted one step every one C-G step to improve quality. For the latent codes optimization, we perform it every 4 timesteps in our method. Our full method takes about 46 minutes to generate a case.

## 8. More Comparisons

In this section, we provide more comparison results with SOTA prompt-guided human generation methods (8.1), single-image human reconstruction methods (8.2), and recent multi-view diffusion model (8.3).

### 8.1. Comparison with Human Generation Methods

We further compare our method with more state-of-the-art text-guided 3D human generation methods, DreamAvatar [6], DreamHuman [31], HumanNorm [20]. These methods are not open-source currently, so the compared cases are selected from their galleries. As illustrated in Fig. 13, DreamAvatar suffers from oversaturation and Janus problem. DreamHuman demonstrates good color saturation, however, its low resolution imposes constraints on the final quality. HumanNorm produces mesh with texture, but its results lack photo-realistic fidelity.

### 8.2. Comparison with Image-to-3D Methods

We further qualitatively compare our method with state-of-the-art image-to-3D methods. As illustrated in Fig. 14, SyncDreamer [39] and Wonder3D [41] both take an image as input and output multi-view images for reconstruction. They can generate human body images from various views, but the appearance has blurry artifacts, particularly in the facial part. TeCH [22] takes more than 3 hours on its network fine-tuning and SDS-based optimization. It faithfully recovers the input image on frontal view, but exhibits some inconsistencies on other views. Our method achieves high-quality generation while maintaining high view consistency.

### 8.3. Comparison with Multi-view Diffusion Model

We further compare our method with the direct outcomes of MVDream. Given a text prompt, MVDream can generate four images of $256 \times 256$ resolution in less than 1 minute. MVHuman can also be easily customized to generate specific four views. As illustrated in Fig. 15, our method can generate high-quality results while the results of MVDream tend to be cartoon-like. Notably, our method takes less than 7 minutes for a case under the situation without any code optimization.
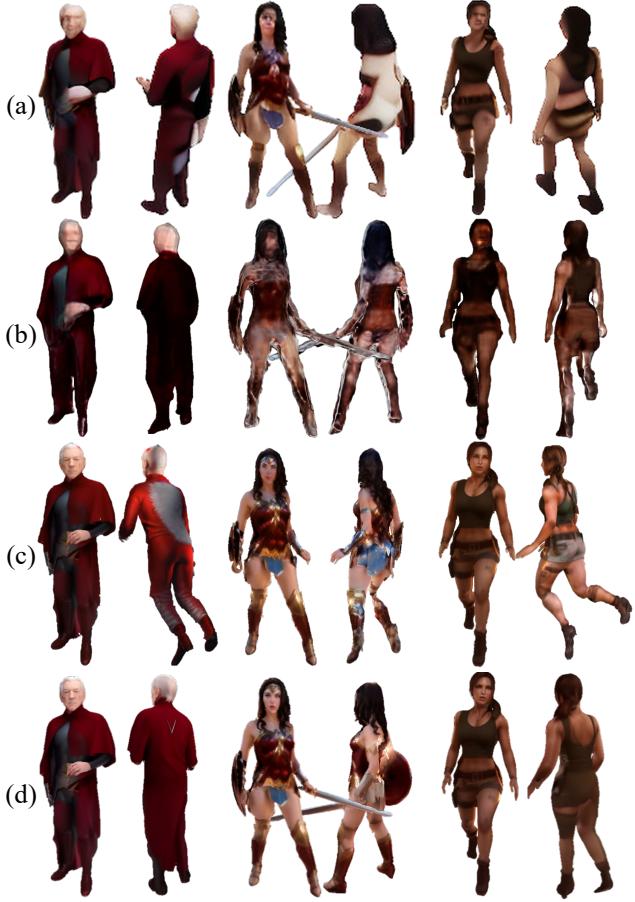
3

Figure 14. Qualitative comparison with image-to-3D methods. (a) SyncDreamer; (b) Wonder3D; (c) TeCH; (d) Ours.
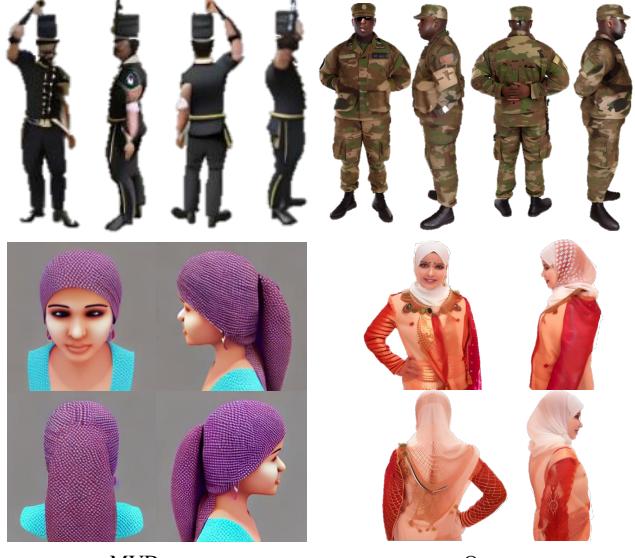


Figure 15. Qualitative comparison with MVDream. Our appearance is more photo-realistic. The prompts from top to down are *"a black army soldier"*, *"an Indian woman wearing a hood, upper body"*. (* Both of ours are cases using scan meshes)



(a)        (b)        (c)

Figure 16. Qualitative evaluation of the multi-view sampling. (a) w/o C-G noise; (b) w/o optimization; (c) full

Table 2. Quantitative evaluation of view consistency.

| Method | PSNR ↑ | SSIM ↑ |
|---|---|---|
| w/o C-G noise | 29.497 | 0.9630 |
| w/o optimization | 28.121 | 0.9699 |
| full | 34.074 | 0.9860 |



(a)        (b)        (c)

Figure 17. Qualitative evaluation of mixed sampling steps. (a) one original sampling step every one C-G step; (b) one original sampling step every two C-G steps; (c) all with C-G steps.



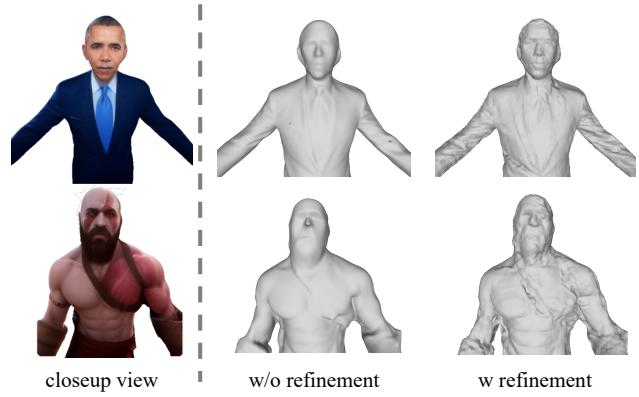closeup view     w/o refinement     w refinement

Figure 18. Qualitative evaluation of geometry refinement. With geometry refinement, we can obtain fine-grained geometry having details aligned with the appearance in the images.
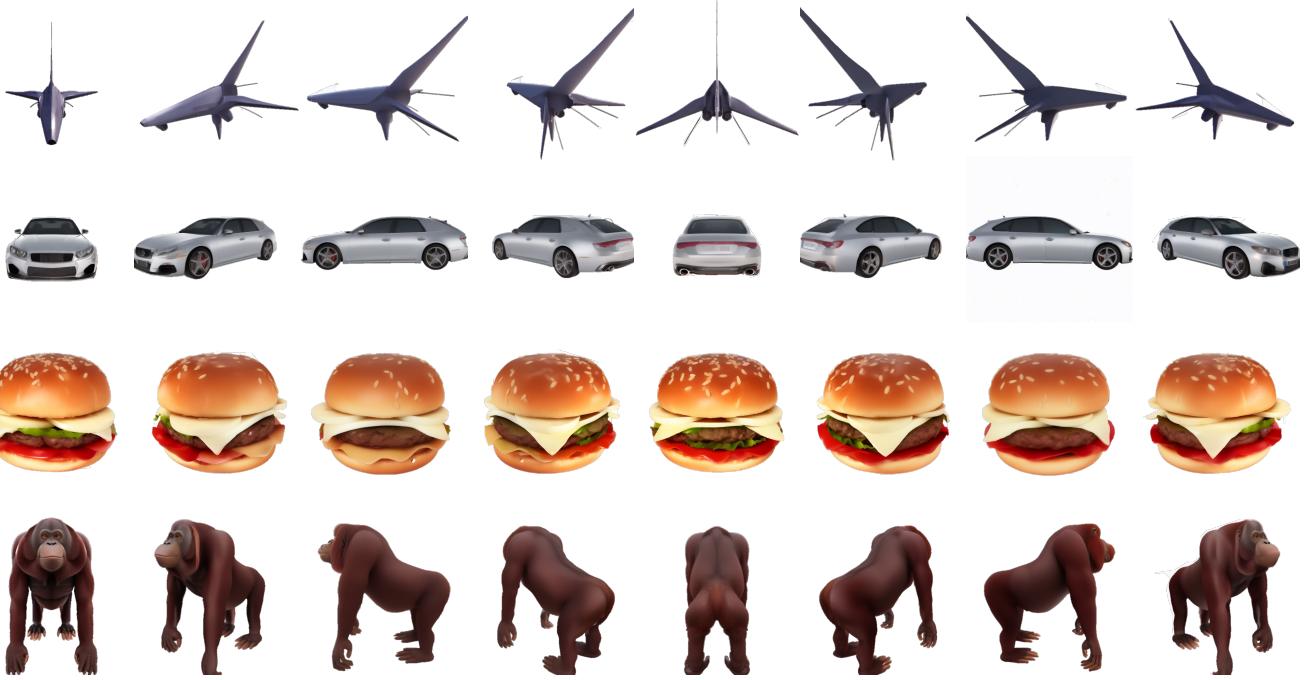
Figure 19. More results on general objects. The prompts from top to down are *"a starship"*, *"a silver car"*, *"a delicious hamburger"*, *"an orangutan"*.

## 9. More Ablation Study

### 9.1. Multi-view Sampling

We further evaluate the multi-view sampling process. In Fig. 16, without the consistency-guided noise (**w/o C-G noise**), though the outcomes appear satisfactory from each view, the overall consistency is lacking because the optimization alone can not provide enough constraints on view consistency. Without the latent codes optimization (**w/o optimization**), since we mix the use of the original predicted noise and the C-G noise, the outcomes are also not strictly view-consistent. Our full method achieves high view consistency and maintains high quality at the same time. As shown in Tab. 2, the full method achieves the best score.

We then evaluate the ratio between sampling steps using the original predicted noise (original sampling step) and steps using the C-G noise (C-G step). Our full sampling process adopts one original sampling step every one C-G step (a step using C-G noise for sampling). As shown in Fig. 17(a), with this setting, we can obtain results that are both high-quality and view-consistent. In Fig. 17(b), the sampling process adopts one original sampling step every two C-G steps, and in Fig. 17(c), all sampling steps are C-G steps. We can observe that as the proportion of C-G steps increases, it sacrifices some visual effects in exchange for higher view consistency.

### 9.2. Geometry Refinement

We additionally evaluate the geometry refinement of the post-process. As illustrated in Fig. 18. The initial geometry appears to be smooth. The geometry refinement adds fine-grained details to the initial geometry. Meanwhile, these details also align well with the appearance depicted in the images.

## 10. More Results

We also test our method on generating general objects by customizing the multi-view sampling process to generate images of 8 views in a single circular track. The conditions for ControlNet are depth and normal maps. As illustrated in Fig. 19, our method can also be applied to generate general objects. However, it may fail when the given conditions are ambiguous (e.g. the provided geometry is too smooth).

## 11. Ethics Statement

Though human eyes can distinguish real-world photos and results of MVHuman, we worry about the potential ethical risk. Thus we want to state again that MVHuman should not be utilized to deceive those who are unfamiliar with this domain. Moreover, like all other generative methods, MVHuman should also not be utilized to generate results that violate the diversity of gender, race, and culture.