

# HVTR: Hybrid Volumetric-Textural Rendering for Human Avatars

Tao Hu<sup>1</sup>, Tao Yu<sup>2</sup>, Zerong Zheng<sup>2</sup>, He Zhang<sup>3</sup>, Yebin Liu<sup>2\*</sup>, Matthias Zwicker<sup>1</sup>  
<sup>1</sup>University of Maryland, College Park <sup>2</sup>Tsinghua University <sup>3</sup>Beihang University

## Abstract

We propose a novel neural rendering pipeline, Hybrid Volumetric-Textural Rendering (HVTR), which synthesizes virtual human avatars from arbitrary poses efficiently and at high quality. First, we learn to encode articulated human motions on a dense UV manifold of the human body surface. To handle complicated motions (e.g., self-occlusions), we then leverage the encoded information on the UV manifold to construct a 3D volumetric representation based on a dynamic pose-conditioned neural radiance field. While this allows us to represent 3D geometry with changing topology, volumetric rendering is computationally heavy. Hence we employ only a rough volumetric representation using a pose-conditioned downsampled neural radiance field (PD-NeRF), which we can render efficiently at low resolutions. In addition, we learn 2D textural features that are fused with rendered volumetric features in image space. The key advantage of our approach is that we can then convert the fused features into a high resolution, high quality avatar by a fast GAN-based textural renderer. We demonstrate that hybrid rendering enables HVTR to handle complicated motions, render high quality avatars under user-controlled poses/shapes and even loose clothing, and most importantly, be fast at inference time. Our experimental results also demonstrate state-of-the-art quantitative results. More results are available at our project page: <https://www.cs.umd.edu/~taohu/hvtr/>.

## 1. Introduction

Capturing and rendering realistic human appearance under varying poses and viewpoints is an important goal in computer vision and graphics. Recent neural rendering methods [65, 51, 26, 54] have made great progress in generating realistic images of humans, which are simple yet effective compared with traditional graphics pipelines [2, 3, 73].

Given a training dataset of multiple synchronized RGB videos of a human, the goal is to build an animatable vir-

tual avatar with pose-dependent geometry and appearance of the individual that can be driven by arbitrary poses from arbitrary viewpoints at inference time. We propose Hybrid Volumetric-Textural Rendering (HVTR). To represent the input to our system, including the pose and the rough body shape of an individual, we employ a skinned parameterized mesh (SMPL [31]) fitted to the training videos. Our system is expected to handle the articulated structure of human bodies, various clothing styles, non-rigid motions, and self-occlusions, and be fast at inference time. In the following, we will introduce how HVTR solves these challenges by proposing (1) effective pose encoding for better generalization, (2) rough yet effective volumetric representation to handle changing topology, and (3) hybrid rendering for fast and high quality rendering.

*Pose Encoding on a 2D Manifold.* The first challenge lies in encoding the input pose information so that it can be leveraged effectively by the image synthesis pipeline. Prior methods employ global pose parameter conditioning [75, 48, 35, 25], or Peng et al. [51] learn poses in a 3D sparse voxelized space by SparseConvNet [10]. For better pose generalization, [26, 50, 5] learn motions using skinning weights via a backward (or inverse) skinning step. Skinning weights are not powerful enough to represent complicated deformations due to arbitrary motions and various clothing styles, however, which may cause averaged and blurry results. In addition, changing topology is challenging for backward skinning used in [50], as it is not able to model one-to-many backward correspondences [6]. In contrast, we encode motions on a 2D UV manifold of the body mesh surface, and the dense representation enables us to utilize 2D convolutional networks to effectively encode pose-dependent features. We define a set of geometry and texture latents on the 2D manifold, which have higher resolution than the compressed latent vectors used in [30] to enable capturing local motion/appearance details for rendering. Since our method does not employ backward skinning, we also avoid the multi-correspondence problem.

*Rough Yet Effective Volumetric Representation.* Our input is a coarse SMPL mesh as used in [52, 54, 17], which cannot capture detailed pose- and clothing-dependent deformations. Inspired by the recent neural scene representations

\*Yebin Liu is the corresponding author.

[40, 26, 51], we model articulated humans with an implicit volumetric representation by constructing a dynamic pose-conditioned neural radiance field. This volumetric representation has the built-in flexibility to handle changing geometry and topology. Different from NeRF [40] for static scenes, we condition our proposed dynamic radiance field on our pose encoding defined on the UV manifold. This enables capturing pose- and view-dependent volumetric features. Constructing the radiance field is computationally heavy [40, 26, 51], however, hence we propose to learn only a rough volumetric representation by constructing a pose-conditioned downsampled neural radiance field (PD-NeRF). This allows us to balance the competing challenges of achieving computational complexity while still being able to effectively resolve self-occlusions. Yet learning PD-NeRF from low resolution images is challenging, and to address this, we propose an appropriate sampling scheme. We show that we can effectively train PD-NeRF from images with a size of only  $45 \times 45$  and as few as 7 sampled points along each query ray (see Fig. 3, 4 and Tab. 4).

**Hybrid Rendering.** The final challenge is to render full resolution images by combining the downsampled PD-NeRF and our learned latents on the 2D UV manifold. To solve this, we rasterize the radiance field into multi-channel (not just RGB) volumetric features in image space by volume rendering. The rasterized volumetric features preserve both geometric and appearance details [40]. In addition, we extract 2D textural features from our latents on the UV manifold for realistic image synthesis following the spirit of Deferred Neural Rendering (DNR) [66, 54, 17]. We fuse the 3D and 2D features by utilizing Attentional Feature Fusion (AFF[8]), and finally use a 2D GAN-based [9] textural rendering network (TexRenderer) to decode and supersample them into realistic images. Though TexRenderer works in image space, it is able to incorporate the rasterized volumetric features for geometry-aware rendering.

The hybrid rendering brings several advantages. 1) We are able to handle self-occlusions by volume rendering. 2) We can generate high quality details using GAN and adversarial training. This enables us to handle uncertainties involved in modeling dynamic details, and is well-suited for enforcing realistic rendered images [54, 18, 52]. In contrast, a direct deterministic regression of dynamic scenes often leads to blurry results as stated in [26]. 3) Both textural rendering, and volumetric rendering only from rough volumetric representations, are fast at training and inference time. In contrast, regressing a detailed geometry for view synthesis by volume rendering is time consuming. 4) Benefiting from (1) and (2) and leveraging the rough geometry and GAN-based rendering, we are able to handle loose clothing like skirts (Fig. 5).

In summary, our contributions are: (1) We propose HVTR, a novel neural rendering pipeline, to generate hu-

Animatable Pipelines	Render- er	Geom- Recon	Fast Infer.
2D : EDN [4], vid2vid [69]	NIT	✗	✓
2D Plus : SMPLpix [52], DNR [66], ANR[54]	NIT	✗	✓
3D : NB[51], AniNeRF[50]	VolR	✓	✗
3D : Ours	Hybrid	✓	✓

Table 1: A set of recent human synthesis approaches classified by feature representations (2D or 3D) and rendering methods (NIT: neural image translation; VolR: volume rendering [22]). NB: Neural Body[51]. AniNeRF: Animatable NeRF[50].

man avatars from arbitrary skeleton motions using a hybrid strategy. HVTR achieves state-of-the-art performance, and is able to handle complicated motions, render high quality avatars even with loose clothing, generalize to novel poses, and support body shape control. Most importantly, it is fast at inference time. (2) HVTR uses an effective scheme to encode pose information on the UV manifold of body surfaces, and leverages this to learn a pose-conditioned downsampled NeRF (PD-NeRF) from low resolution images. Our experiments show how the rendering quality is influenced by the PD-NeRF resolution, and that even low resolution volumetric representations can produce high quality outputs at a small computational cost. (3) HVTR shows how to construct PD-NeRF and extract 2D textural features all based on pose encoding on the UV manifold, and most importantly, how the two can be fused and incorporated for fast, high quality, and geometry-aware neural rendering.

## 2. Related Work

**Neural Scene Representations.** Instead of explicitly modeling geometry, many neural rendering methods [62, 63, 41, 66] propose to learn implicit representations of scenes, such as DeepVoxels [62], Neural Volumes [30] SRNs [63], or NeRF [41]. In contrast to these static representations, we learn a dynamic radiance field on the UV manifold of human surfaces to model articulated human bodies.

**Shape Representations of Humans.** To capture detailed deformations of human bodies, most recent papers utilize implicit functions [37, 38, 7, 47, 56, 57, 19, 58, 39, 68, 45, 67, 78, 20, 77] or point clouds [34, 36] due to their topological flexibility. These methods aim at learning geometry from 3D datasets, whereas we synthesize human images of novel poses only from 2D RGB training images.

**Rendering Humans by Neural Image Translation (NIT).** Some existing approaches render human avatars by neural image translation, i.e. they map the body pose given in the form of renderings of a skeleton [4, 60, 53, 24, 79, 69], dense mesh [28, 70, 27, 59, 42, 12] or joint position heatmaps [32, 1, 33], to real images. As summarized in Tab. 1, EDN [4] and vid2vid [69] utilize GAN [9] net-

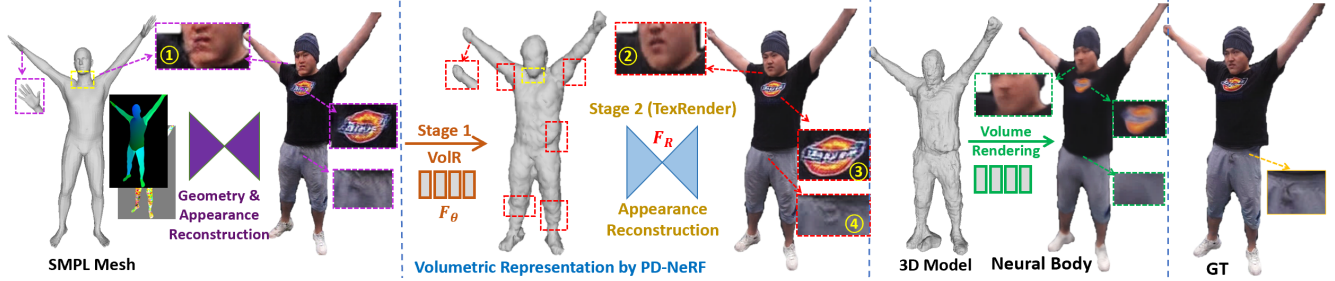


Figure 1: We illustrate the differences between (left) NIT methods (DNR), (middle) our hybrid approach, and (right) NeRF methods (Neural Body). DNR[66] and SMPLpix[52] are based on fixed mesh (SMPL[31] or SMPL-X[49]), and use a GAN for one-stage rendering without explicit geometry reconstruction. As a disadvantage, DNR needs to resolve geometric misalignments implicitly, which often leads to artifacts (see closeup ① in the figure). Yet our method (middle) works in two stages by first learning a downsampled volumetric representation (by PD-NeRF), and then utilizing a GAN for appearance synthesis. Though only learned from low resolution images ( $90 \times 90$  in this example), the rough volumetric representation still encodes more 3D pose-dependent features than SMPL, which enables us to handle self-occlusions (region ① vs ②), and preserve more details (③④) than DNR. In addition, our GAN-based renderer can generate high resolution wrinkles, whereas Neural Body cannot, although its geometry is more detailed. In addition, our approach is  $52\times$  faster than Neural Body at inference time (see Tab. 4).

works to learn a mapping from 2D poses to human images. To improve temporal stability and learn a better mapping, “2D Plus” methods [66, 52, 17] are conditioned on a coarse mesh (SMPL [31]), and take as input additional geometry features, such as DNR (UV mapped features) [66], SMPLpix (+ depth map) [52], and ANR (UV + normal map) [54]. A 2D convolutional network is often utilized for both shape completion and appearance synthesis in one stage [66, 52]. However, [66, 52, 54, 11] do not reconstruct geometry explicitly and cannot handle self-occlusions effectively. In contrast, our rendering is conditioned on a learned 3D volumetric representation using a two-stage approach (see Fig. 1). We show that our learned representation handles occlusions more effectively than other techniques [66, 52, 54, 11] that just take geometry priors (e.g., UV, depth or normal maps) from a coarse mesh as input (see Fig. 3, 7).

**Rendering Humans by Volume Rendering (VolR).** For stable view synthesis, recent papers [51, 26, 50, 44, 64, 5] propose to unify geometry reconstruction with view synthesis by volume rendering, which, however, is computationally heavy. In addition, the appearance synthesis (e.g., Neural Body [51]) largely relies on the quality of geometry reconstruction, which is very challenging for dynamic humans, and imperfect geometry reconstruction will lead to blurry images (Fig. 7). Furthermore, due to the difficulties of reconstructing geometry for loose clothing, most animatable NeRF methods (e.g., state-of-the-art Neural Actor [51]) cannot handle skirts. In contrast, our method utilizes a GAN to render high frequency details based on a downsampled radiance field, which makes our method more robust to geometry inaccuracies and fast at inference time. We can also render skirts (Fig. 5). A comparison of ours, NIT, and

VolR is shown in Fig. 1.

Ours is distinguished from [43] by rendering dynamic humans at high resolutions, and conditioning our rendering framework on the UV manifold of human body surfaces.

### 3. Method

*Problem setup.* Our goal is to render pose- and view-dependent avatars of an individual from an arbitrary pose  $\mathbf{p}$  (represented by a coarse human mesh) and an arbitrary viewpoint (position  $\mathbf{o}$ , view direction  $\mathbf{d}$ ):

$$I_{\mathbf{p} \rightarrow \mathbf{t}} = \text{HTVR}(\mathbf{p}, \mathbf{o}, \mathbf{d}, K)$$

where  $K$  denotes the camera intrinsic parameters, and  $I_{\mathbf{p} \rightarrow \mathbf{t}} \in \mathbb{R}^{W \times H \times 3}$  is the output image.

We first introduce our pose encoding method (Sec. 3.1), and based on this we present how to extract 2D textural features (Sec. 3.2) and 3D pose-dependent volumetric features (Sec. 3.3). Finally, we describe how we fuse these features and synthesize the final RGB avatars (Sec. 3.4). Fig. 2 shows an outline of the proposed framework.

#### 3.1. Pose Encoding

Given a skeleton of pose  $\mathbf{p}$  and posed SMPL mesh  $I_{\mathbf{p}}$  as input, we first project each surface point  $s_i \in I_{\mathbf{p}}$  from 3D space to its UV manifold, represented by a UV positional map  $M \in \mathbb{R}^{U \times U \times 3}$ ,  $M_{u_i} = s_i$ , where  $u_i$  describes the relative location of the point on the body surface manifold. With this, we define a set of geometry latents  $Z_G \in \mathbb{R}^{U \times U \times C_g}$  to represent the intrinsic local geometry features, and texture latents  $Z_T \in \mathbb{R}^{U \times U \times C_t}$  to represent high-dimensional neural textures as used in [66, 54, 17]. Both latents are defined

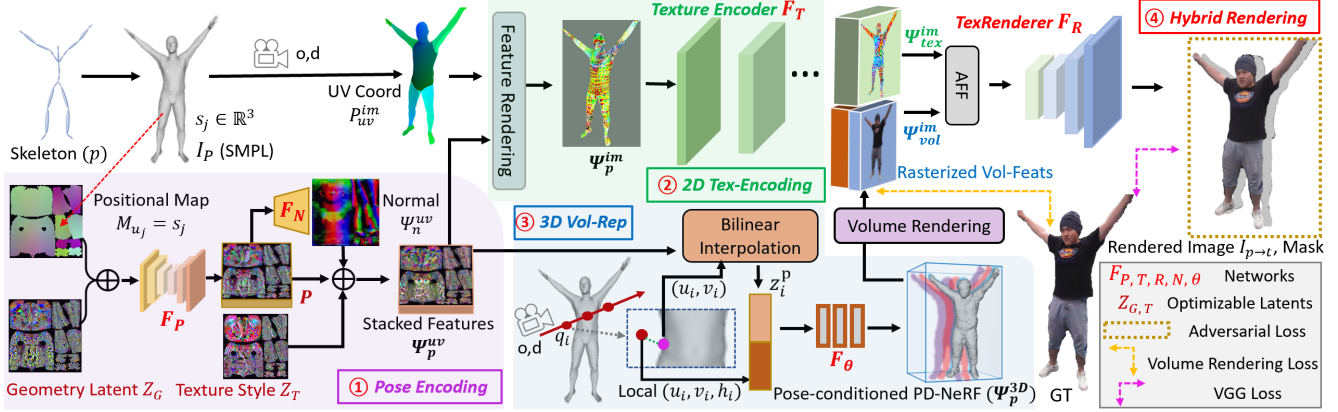


Figure 2: Pipeline overview. Given a coarse SMPL mesh  $I_p$  with pose  $(p)$  and a target viewpoint  $(o, d)$ , our system renders a detailed avatar  $I_{p \rightarrow t}$  using four main components: ① pose encoding, ② 2D textural feature encoding, ③ 3D volumetric representation, and ④ hybrid rendering. ① **Pose Encoding** in UV space: We record the 3D positions of the mesh  $p_j \in I_p$  on a UV positional map  $M$ . We stack it with a geometry latent ( $Z_G$ ) and encode them into a pose feature tensor  $P$ . We then construct the pose-dependent features  $\Psi_p^{uv}$  by stacking  $P$ , an optimizable texture style latent  $Z_T$ , and the estimated normals  $\Psi_n^{uv}$  in UV space. ② **2D Tex-Encoding**: A *Feature Rendering* module renders the coarse mesh with  $\Psi_p^{uv}$  into image features  $\Psi_p^{im}$  by utilizing a rasterized UV coordinate map ( $P_{uv}^{im}$ ). The image features are then encoded as 2D textural features  $\Psi_{tex}^{im}$  by the *Textural Encoder*  $F_T$ . ③ **3D Vol-Rep**: To capture the rough geometry and address self-occlusion problems, we further learn a volumetric representation by constructing a pose-conditioned downsampled neural radiance field (PD-NeRF)  $\Psi_p^{3D}$  to encode 3D pose-dependent features. ④ **Hybrid Rendering**: PD-NeRF is rasterized into image space  $\Psi_{vol}^{im}$  by volume rendering, where 3D volumetric features are also preserved. Both the 2D  $\Psi_{tex}^{im}$  and 3D features  $\Psi_{vol}^{im}$  are pixel-aligned in image space, fused by Attentional Feature Fusion (AFF), and then converted into a realistic image  $I_{p \rightarrow t}$  and a mask by *TexRenderer*  $F_R$ .

in UV space, and shared across different poses. Our geometry and texture latents have higher resolution than the compressed representation used in other works (e.g., latent vectors used in [30]), which enables us to capture local details, and the rendering pipeline can leverage them to infer local geometry and appearance changes.

To extract geometric features, the projected poses  $M$  and the geometry latents  $Z_G$  are convolved by a ConvNet  $F_P$  to obtain pixel-aligned pose features  $P$ . In addition, to enforce learning geometric features by  $F_P$ , we predict the normal  $\Psi_n^{uv} \in \mathbb{R}^{U \times V \times 3}$  of the posed mesh in UV space using a shallow ConvNet  $F_N$ :  $\Psi_n^{uv} = F_N(P)$ . We then concatenate the geometric features  $\Psi_n^{uv}$  and  $P$ , and the texture features  $Z_T$  to obtain our pose-dependent features  $\Psi_p^{uv} = \text{Cat}(\Psi_n^{uv}, P, Z_T)$ .

Note that though the UV positional map used is similar to [34, 36], ours is distinguished by learning pose-dependent features from 2D images instead of 3D point clouds, and we have a normal estimation network to enforce geometric learning, whereas [34, 36] did not.

### 3.2. 2D Textural Feature Encoding

Given a viewpoint  $(o, d)$ , we also render a UV coordinate map  $P_{uv}^{im}$  [13] to encode the shape and pose features of  $I_p$  in image space. This allows us to transform  $\Psi_p^{uv}$  from UV space to pose-dependent features  $\Psi_p^{im}$  in image space

using a *Feature Rendering* module [66, 54, 17]. We further encode  $\Psi_p^{im}$  as a high-dimensional textural feature using a *Texture Encoder*  $F_T$  implemented by a 2D ConvNet.

### 3.3. 3D Volumetric Representation

Though existing methods achieve compelling view synthesis by just rendering with 2D textural features [66, 17, 54], they cannot handle self-occlusions effectively since they do not reconstruct the geometry. We address this by learning a 3D volumetric representation using a pose-conditioned neural radiance field  $\Psi_p^{3D}$  (PD-NeRF in Fig. 2).

We include pose information to learn the volumetric representation by looking up the encoded pose feature in  $\Psi_p^{uv}$  corresponding to each 3D query point. To achieve this, we project each query point  $q_i$  in the posed space of  $I_p$  to a local point  $\hat{q}_i = (u_i, v_i, h_i)$  in an UV-plus-height space,

$$(u_i, v_i, f_i) = \arg \min_{u, v, f} \|q_i - B_{u, v}(V_{[Tri(f)]})\|_2, \quad (1)$$

where  $f \in \{1, \dots, N_F\}$  is the triangle index,  $Tri(f)$  is the triangle (face),  $V_{[Tri(f)]}$  are the three vertices of  $Tri(f)$ ,  $(u, v)$  are the barycentric coordinates of the face, and  $B_{u, v}(\cdot)$  is the barycentric interpolation function. The height  $h_i$  is given by the signed distance of  $q_i$  to the nearest face  $Tri(f_i)$ .

With this, we sample the local feature  $z_i^p$  of  $\hat{q}_i$  from the encoded pose features  $\Psi_p^{uv}$ :  $z_i^p = B_{u_i, v_i}(\Psi_p^{uv})$ . Given a cam-

era position  $\mathbf{o}$  and view direction  $\mathbf{d}$ , we predict the density  $\sigma$  and appearance features  $\xi$  of  $q_i$  as

$$F_\theta : (\gamma(\hat{q}_i), \gamma(\mathbf{d}), z_i^p) \rightarrow (\sigma, \xi), \quad (2)$$

where  $\gamma$  is a positional encoder. Note that  $\xi \in \mathbb{R}^h$  is a high-dimensional feature vector, where the first three channels are RGB colors. A key property of our approach is that  $F_\theta$  is conditioned on high resolution encoded pose features  $z_i^p$  instead of pose parameters  $\mathbf{p}$ .

### 3.4. Hybrid Volumetric-Textural Rendering

Though the radiance field PD-NeRF can be directly rendered into target images by volume rendering [22], this is computationally heavy. In addition, a direct deterministic regression using RGB images often leads to blurry results in dynamic scenes as stated in [26].

**Volumetric Rendering.** To address this, we use PD-NeRF to render downsampled images by a factor  $s$  for fast inference. We rasterize PD-NeRF into multi-channel volumetric features  $\Psi_{vol}^{im} \in \mathbb{R}^{W_d \times H_d \times h}$ , and each pixel  $\mathfrak{R}(r, p)$  is predicted by  $N$  consecutive samples  $\{x_1, \dots, x_N\}$  along the corresponding ray  $r$  through volume rendering [22],

$$\mathfrak{R}(r, p) = \sum_{n=1}^N \left( \prod_{i=1}^{n-1} e^{-\sigma_i \delta_i} \right) \cdot (1 - e^{-\sigma_n \delta_n}) \cdot \xi_n, \quad (3)$$

where  $\delta_n = \|x_n - x_{n-1}\|_2$ , and density and appearance features  $\sigma_n, \xi_n$  of  $x_n$  are predicted by Eq. 2. Note the first three channels of  $\Psi_{vol}^{im}$  are RGB, which are supervised by downsampled ground truth images (see Fig. 2).

**Attentional Volumetric Textural Feature Fusion.** With both the 2D textural features  $\Psi_{tex}^{im}$  and the rasterized 3D volumetric features  $\Psi_{vol}^{im}$ , the next step is to fuse them and leverage them for 2D image synthesis. This poses several challenges. First,  $\Psi_{tex}^{im}$  is trained in 2D, which converges faster than  $\Psi_{vol}^{im}$ , since  $\Psi_{vol}^{im}$  needs to regress a geometry by optimizing downsampled images, and NeRF training generally converges more slowly for dynamic scenes [51]. Second,  $\Psi_{tex}^{im}$  has higher dimensions (both resolution and channels) than  $\Psi_{vol}^{im}$ , because  $\Psi_{vol}^{im}$  is learned from downsampled images with relatively weak supervision. Due to this, the system may tend to ignore volumetric features of  $\Psi_{vol}^{im}$  at this stage. To solve this problem, we first use a ConvNet to downsample  $\Psi_{tex}^{im}$  to the same size as  $\Psi_{vol}^{im}$ . We also extend the channels of  $\Psi_{vol}^{im}$  to the same dimensionality as  $\Psi_{tex}^{im}$  using a ConvNet. Another approach would be to upsample the resolution of  $\Psi_{vol}^{im}$  instead of extending channels, but we found this destabilizes the training of PD-NeRF.

Finally, we fuse the resized features by Attentional Feature Fusion (AFF [8]):  $\Psi_{vt}^{im} = AFF(\Psi_{vol}^{im}, \Psi_{tex}^{im})$ .  $\Psi_{vt}^{im}$  has the same size as  $\Psi_{vol}^{im}$ . AFF is also learned, and we include it in  $F_R$  in Fig. 2. See [8] for more details about AFF.

**Textural Rendering.** The TexRenderer net  $F_R$  converts the fused features  $\Psi_{vt}^{im}$  into the target avatar  $I_{p \rightarrow t}$  and a mask.

$F_R$  has a similar architecture as Pix2PixHD [71]. See the Appendices for more details.

### 3.5. Optimization

HVTR is trained end-to-end by optimizing networks  $F_{P,T,N,R,\theta}$  and latent codes  $Z_{G,T}$ . Given a ground truth image  $I_t$  and mask  $M_t$ , downsampled ground truth image  $I_t^D$ , and predicted image  $I_{p \rightarrow t}$  and mask  $M_{p \rightarrow t}$ , we use the following loss functions:

**Volume Rendering Loss.** We utilize  $\mathcal{L}_{vol}$  to supervise the training of volume rendering, which is applied on the first three channels of  $\Psi_{vol}^{im}$ ,  $\mathcal{L}_{vol} = \|\Psi_{vol}^{im}[:3] - I_t^D\|_2^2$ .

**Normal Loss.** To enforce learning of geometric features by  $F_P$ , we employ a normal loss  $\mathcal{L}_n$ :  $\mathcal{L}_n = \|\Psi_n^{uv} - N_t^{uv}\|_1$ , where  $N_t^{uv}$  is the ground truth normal of mesh  $I_P$  projected into UV space.

**Feature Loss.** We use a feature loss [21] to measure the differences between the activations on different layers of the pretrained VGG network [61] of the generated image  $I_{p \rightarrow t}$  and ground truth image  $I_t$ ,

$$\mathcal{L}_{feat} = \sum \frac{1}{N^j} \|p^j(I_{p \rightarrow t}) - p^j(I_t)\|_2, \quad (4)$$

where  $p^j$  is the activation and  $N^j$  the number of elements of the  $j$ -th layer in the pretrained VGG network.

**Mask Loss.** The mask loss is  $\mathcal{L}_{mask} = \|M_{p \rightarrow t} - M_t\|_1$ .

**Pixel Loss.** We also enforce an  $\ell_1$  loss between the generated image and ground truth as  $\mathcal{L}_{pix} = \|I_{p \rightarrow t} - I_t\|_1$ .

**Adversarial Loss.** We leverage a multi-scale discriminator  $D$  [71] as an adversarial loss  $\mathcal{L}_{adv}$ .  $D$  is conditioned on both the generated image and feature image  $\Psi_p^{im}$ .

**Face Identity Loss.** We use a pre-trained network to ensure that TexRenderer preserves the face identity on the cropped face of the generated and ground truth image,

$$\mathcal{L}_{face} = \|N_{face}(I_{p \rightarrow t}) - N_{face}(I_t)\|_2, \quad (5)$$

where  $N_{face}$  is the pretrained SphereFaceNet [29].

**Total Loss.**  $\mathcal{L}_{total} = \sum_{i \in \{vol, n, feat, mask, pix, adv, face\}} \lambda_i \mathcal{L}_i$ .

The networks were trained using the Adam optimizer [23]. See the Appendices for more details.

## 4. Experiments

**Dataset.** We evaluate our method on 10 datasets, denoted R1-6, Z1-3, and M1. We captured R1-6, and each dataset has 5 cameras at a resolution of  $1280 \times 720$  (yet big human bounding box) with 800-2800 frames. Z1-3 [51] have 24 cameras ( $1024 \times 1024$ , 620-1400 frames each), and we use splits of 10/7, 12/8, 5/5 separately for training/test cameras. M1 [14] has 101 cameras ( $1285 \times 940$ , 20K frames each), and we utilize 19/8 training/test cameras. For these datasets, we select key sequences to include various motions and use a split of 80%/20% for training and testing. All the tested poses are novel, and rendered viewpoints for R4, Z1-Z3, M1 are new. Yet since these methods render humans in local

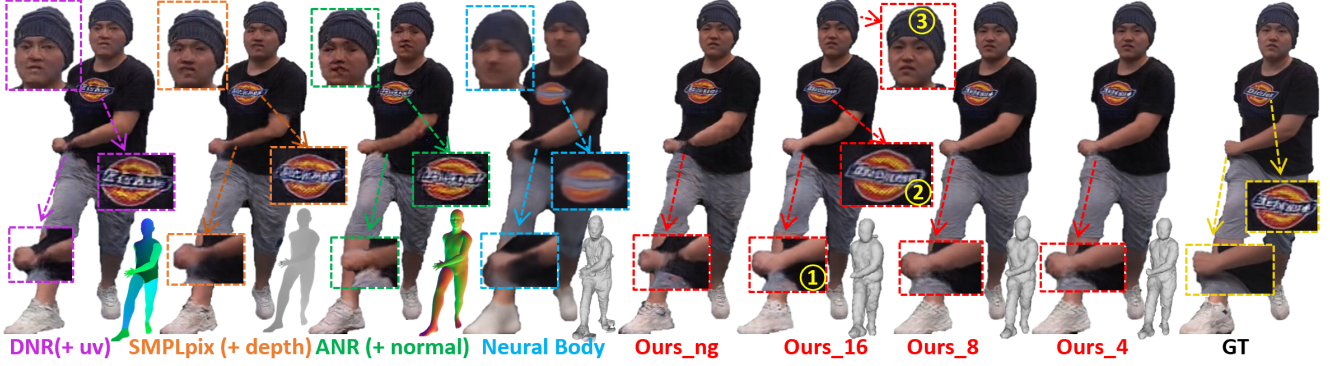


Figure 3: Qualitative results of our variants by changing downsampling factor  $S$  of PD-NeRF. Though existing SMPL based 2D-Plus NIT methods take as input extra geometry priors, such as DNR(+ UV)[66], SMPLpix(+ depth)[52], ANR(UV + normal)[54], they fail to fully utilize the priors for geometry-aware rendering. Instead, ours can handle self-occlusions better ① and also improve the rendering quality (②③) by learning a 1/16 downsampled PD-NeRF (Ours\_16,  $45 \times 45$ ). Note for ours and Neural Body, the learned geometries are shown.

	R1				R2				R3			
Models	LPIPS↓	FID↓	SSIM↑	PSNR↑	LPIPS	FID	SSIM	PSNR	LPIPS	FID	SSIM	PSNR
DNR	.102	75.02	.831	25.73	.125	98.86	.820	27.92	.108	80.33	.809	24.05
SMPLpix	.100	69.81	.835	25.93	.124	94.81	.826	27.97	.104	74.57	.810	24.16
ANR	.117	78.50	.830	26.02	.129	101.72	.825	28.30	.098	69.14	.813	24.29
Neural Body	.212	155.84	.833	26.17	.218	161.99	<b>.833</b>	28.61	.240	165.03	.811	24.16
Ours	<b>.090</b>	<b>62.33</b>	<b>.842</b>	<b>26.22</b>	<b>.108</b>	<b>84.43</b>	<b>.833</b>	<b>28.62</b>	<b>.093</b>	<b>66.01</b>	<b>.823</b>	<b>24.55</b>
	R4				R5				R6			
DNR	.108	93.16	.833	23.34	.136	121.50	.817	24.06	.088	74.77	.864	25.81
SMPLpix	.107	88.14	.837	23.37	.131	118.64	.818	24.10	.077	64.33	.875	26.14
ANR	.138	91.92	.812	23.26	.140	123.55	.823	24.67	.083	63.16	.875	26.61
Neural Body	.198	126.26	<b>.856</b>	<b>24.26</b>	.220	161.93	.816	24.25	.142	94.96	.880	27.19
Ours	<b>.096</b>	<b>78.79</b>	.849	23.98	<b>.117</b>	<b>93.56</b>	<b>.827</b>	<b>24.84</b>	<b>.070</b>	<b>57.00</b>	<b>.891</b>	<b>27.42</b>
	Z1				Z2				Z3			
DNR	.145	92.78	.797	22.06	.145	87.27	.774	25.04	.109	82.79	.826	23.16
SMPLpix	.150	90.90	.797	22.14	.144	81.78	.774	25.18	.113	83.96	.827	22.92
ANR	.205	171.69	.775	<b>22.35</b>	.159	110.85	.778	25.41	.173	123.84	.790	22.14
Neural Body	.215	163.83	.789	22.16	.238	155.27	<b>.792</b>	<b>25.88</b>	.204	167.66	.825	<b>23.89</b>
Ours	<b>.143</b>	<b>90.43</b>	<b>.805</b>	22.31	<b>.132</b>	<b>79.14</b>	.785	25.69	<b>.105</b>	<b>78.03</b>	<b>.829</b>	23.23

Table 2: Quantitative comparisons on nine datasets (averaged on all test views and poses). To reduce the influence of the background, all scores are calculated from images cropped to 2D bounding boxes. LPIPS[76] and FID[16] capture human judgement better than per-pixel metrics such as SSIM[72] or PSNR. All poses are novel, and R4, Z1-Z3 are tested on new views. The pose variations are relatively small in Z1-Z3 datasets, for which we mainly evaluate the capability of capturing/rendering high frequency details instead of pose generalization.

space and the captured human characters move, we found that novel poses mattered more than novel viewpoints for quantitative results. See R1 in Fig. 3, and R2-R4, Z1, Z3 in Fig. 7, M1 in Fig. 5, and the Appendices for more details.

**Baselines.** We compare our method with NIT-based methods (DNR[66], ANR[54], SMPLpix[52]), and NeRF-based Neural Body [51] (as used in [26, 74] for animation synthesis), and Animatable NeRF [50] (see the Appendices). For fair comparisons, DNR, ANR, SMPLpix all have the same network architectures as ours, the same SMPL model as input, and were trained with the losses mentioned in their papers. ANR: Since the code of ANR was not released when

this work was developed, we cannot guarantee our reproduced ANR achieves the performance as expected, though it converges and generates reasonable results. SMPLpix: We follow the author’s recent update<sup>1</sup> to strengthen SMPLpix by rasterizing the SMPL mesh instead of the sparse SMPL vertices [52]. Neural Body[51] and Animatable NeRF[50] are trained with their provided code and setup separately.

**Annotation:** Ours\_ $S(N)$  indicates the variant of our method, where  $S$  is the downsampling factor of PD-NeRF, and  $N$  is the number of sample points along each ray. By default, we use the setting of Ours\_8(12) as our method for

<sup>1</sup><https://github.com/sergeyprokudin/smplpix>

comparisons. Ours\_ng is the variant without PD-NeRF.

#### 4.1. Evaluations

**Differences to the Baseline Methods.** As shown in Fig. 3, compared with 2D-Plus methods (DNR, SMPLpix, ANR), we can handle self-occlusions better and generate more details than Neural Body. We also compare the architecture of ours, DNR, and Neural Body in Fig. 1.

**Comparisons.** We evaluate our methods on the 10 datasets, shown in Fig. 3, 7, 5 (see R5, R6, Z2 in the Appendices). We summarize the quantitative results in Tab. 2, 3, where we achieve the best performances on 34/40 evaluations metrics, and on all the 20 LPIPS/FID scores.

**Rendering Skirts.** Our method is capable of rendering loose clothing like skirts as shown in Fig. 5 with rough geometry reconstruction, whereas solo volume rendering methods (e.g., Neural Actor stated in [26]) generally cannot because they rely more on the quality of geometry reconstruction, which is also challenging for dynamic skirts.

**Accuracy, Inference Time, GPU Memory.** See the accuracy and inference time in Tab. 4. Ours can improve the performance over DNR and SMPLpix by about 10% (even 14% by Ours.4) at a small computational cost, and is almost  $52\times$  faster than Neural Body. For fair comparisons, we evaluate Tab. 4 on R1 (Fig. 1) dataset (about 8k frames for training, 2k for testing), where each frame was cropped to  $720 \times 720$  close to the human bounding box (bbox), to reduce the influence of white background. Yet one **limitation** is that we require more GPU memory in training, Ours.4(20)-most GPU-consuming version: 21GB; Neural Body: 5GB; ANR:11GB. However, in inference, Ours.4(20): 4GB, Neural Body: 15GB. Note that this was evaluated on the cropped bbox with downsampled  $S=4$ , and we can process high resolution like  $1024 \times 1024$  (Z1-3) or  $1285 \times 940$  (M1). See the Appendices for more details.

**Applications.** We can render avatars under user-controlled novel views, poses, and shapes for **Novel View Synthesis, Animation** (Fig. 3-7), and **Shape Editing** (Fig. 6).

M1	LPIPS ↓	FID ↓	SSIM ↑	PSNR ↑
DNR	.195	144.78	.687	19.96
Ours	<b>.179</b>	<b>132.83</b>	<b>.696</b>	<b>20.18</b>

Table 3: Comparisons on M1 dataset under novel poses and views.

#### 4.2. Ablation Study

We analyze how PD-NeRF affects the final rendering quality and inference time by evaluating two parameters: the resolution represented by a downsampling factor  $1/S$ , and the number of sampled points  $N$  along each ray, as shown in Fig. 3, 4. Fig. 3 shows that we can improve the capability of solving self-occlusions by just incorporating a  $1/16$  ( $45 \times 45$ ) downsampled PD-NeRF (Ours\_16 vs



Figure 4: The effectiveness of PD-NeRF vs. downsampling factor  $S$  and sample points  $N$ . The 2nd-7th are variants of our method (see Annotation at Sec. 4). The first three channels of the volumetric features ( $\Phi_{vol}^{im}$  in Fig. 2) are shown at the bottom right. PD-NeRF improves the capability of handling self-occlusions (e.g., cheeks {3④⑤⑥⑦} vs {1②}), and the quality (⑥ vs ⑤) can significantly be improved by increasing  $S$ , which is also shown in Tab. 4.

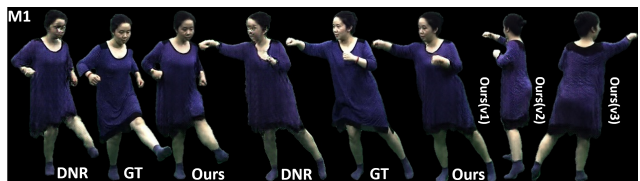


Figure 5: We can render skirts on novel poses and view-points.



Figure 6: Rendering results of HVTR for different body shapes of the same individual. Top-left: SMPL shapes (visualized as UV coordinate maps); Bottom-left: renderings of PD-NeRF; Middle: normal shape. Both PD-NeRF and HVTR generate reasonable results. Not just a straightforward texture to shape mapping, HVTR can generate some shape-dependent wrinkles (marked in red for big models), though these shapes were not seen in training.

Ours\_ng). Tab. 4 shows that the quantitative results can be improved by increasing  $S$  (e.g., Ours.8(12) vs Ours\_16(12), or ⑥ vs ④ in Fig. 4) or sampling more points (e.g., Ours.16(12) vs Ours\_16(7)), which illustrate the effectiveness of PD-NeRF. Yet  $S$  contributes more than  $N$  as shown in Tab. 4 and Fig. 4. It seems that  $N$  easily reaches the peak where the performance cannot be improved obviously, such as Ours.16(12) vs Ours\_16(20), as listed in Tab. 4. Yet Ours.8(12) significantly outperforms Ours\_16(12) by doubling the resolution as seen in Tab. 4 and Fig. 4, which illustrate the effectiveness of higher resolution PD-NeRF.

Models	LPIPS	FID	Time (s)	VR_T(%)
DNR	0.102	75.015	.184	-
SMPLpix	0.100	69.812	.198	-
ANR	0.117	78.501	.224	-
NeuralBody	0.212	155.838	18.20	-
Ours_ng	0.099	70.528	.257	-
Ours_16(7)	0.097	64.871	.292	11.99
Ours_16(12)	0.096	63.792	.295	12.88
Ours_16(20)	0.096	63.834	.305	15.41
Ours_8(12)	0.090	62.333	.349	26.36
Ours_4(20)	<b>0.086</b>	<b>60.788</b>	.464	44.61

Table 4: Accuracy and inference time on novel poses. VR\_T(%) indicates the percentages of the volume rendering time. We test the end-to-end inference time on a GeForce RTX 3090, and the time for rendering the required maps are also counted, such as DNR (UV coord maps), SMPLpix (depth maps), ANR (UV coord + normal maps), ours: UV coord + depth maps (used in PD-NeRF to sample query points). PyTorch3D[55] is used for rendering.

The ablation study of face identity loss and feature fusion can be found in the Appendices.

## 5. Discussion and Conclusion

**Potential Societal Impact.** Our method enables a digital portrait copy which can be reenacted by another portrait video. Therefore, given a portrait video of a specific person, it can be used to generate portrait videos, which need to be addressed carefully before deploying the technique.

**Conclusion.** We introduce Hybrid Volumetric-Textural Rendering (HVTR), a novel neural rendering pipeline, to generate human avatars under user-controlled poses, shapes and viewpoints. HVTR can handle complicated motions, render loose clothing, and provide fast inference. The key is to learn a pose-conditioned downsampled neural radiance field to handle changing geometry, and to incorporate both neural image translation and volume rendering techniques for fast geometry-aware rendering. We see our framework as a promising component for real-time telepresence.

## References

- [1] Kfir Aberman, M. Shi, Jing Liao, Dani Lischinski, B. Chen, and D. Cohen-Or. Deep video-based performance cloning. *Computer Graphics Forum*, 38, 2019. 2
- [2] George Borshukov, Dan Piponi, Oystein Larsen, J. P. Lewis, and Christina Tempelaar-Lietz. Universal capture: image-based facial animation for "the matrix reloaded". In *SIGGRAPH '03*, 2003. 1
- [3] Joel Carranza, Christian Theobalt, Marcus A. Magnor, and Hans-Peter Seidel. Free-viewpoint video of human actors. *ACM SIGGRAPH 2003 Papers*, 2003. 1
- [4] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A. Efros. Everybody dance now. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5932–5941, 2019. 2
- [5] Jianchuan Chen, Ying Zhang, Di Kang, Xuefei Zhe, Linchao Bao, and Huchuan Lu. Animatable neural radiance fields from monocular rgb video. *ArXiv*, abs/2106.13629, 2021. 1, 3
- [6] Xu Chen, Yufeng Zheng, Michael J. Black, Otmar Hilliges, and Andreas Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. *ArXiv*, abs/2104.03953, 2021. 1
- [7] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5932–5941, 2019. 2
- [8] Yimian Dai, Fabian Gieseke, Stefan Oehmcke, Yiquan Wu, and Kobus Barnard. Attentional feature fusion. In *IEEE Winter Conference on Applications of Computer Vision, WACV 2021*. 2, 5, 15, 16
- [9] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014. 2
- [10] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9224–9232, 2018. 1
- [11] Artur Grigorev, Karim Isakov, Anastasia Ianina, Renat Bashirov, Ilya Zakharkin, Alexander Vakhitov, and Victor S. Lempitsky. Stylepeople: A generative model of fullbody human avatars. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5147–5156, 2021. 3
- [12] A. K. Grigor'ev, Artem Sevastopolsky, Alexander Vakhitov, and Victor S. Lempitsky. Coordinate-based texture inpainting for pose-guided human image generation. *Computer Vision and Pattern Recognition (CVPR)*, pages 12127–12136, 2019. 2
- [13] R. A. Güler, N. Neverova, and I. Kokkinos. Densepose: Dense human pose estimation in the wild. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7297–7306, 2018. 4
- [14] Marc Habermann, Lingjie Liu, Weipeng Xu, Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. Real-time deep dynamic characters. *ACM Transactions on Graphics (TOG)*, 40:1 – 16, 2021. 5
- [15] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 13
- [16] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NIPS*, 2017. 6, 15
- [17] Tao Hu, Kripasindhu Sarkar, Lingjie Liu, Matthias Zwicker, and Christian Theobalt. Egorenderer: Rendering human avatars from egocentric camera images. In *Proceedings of*

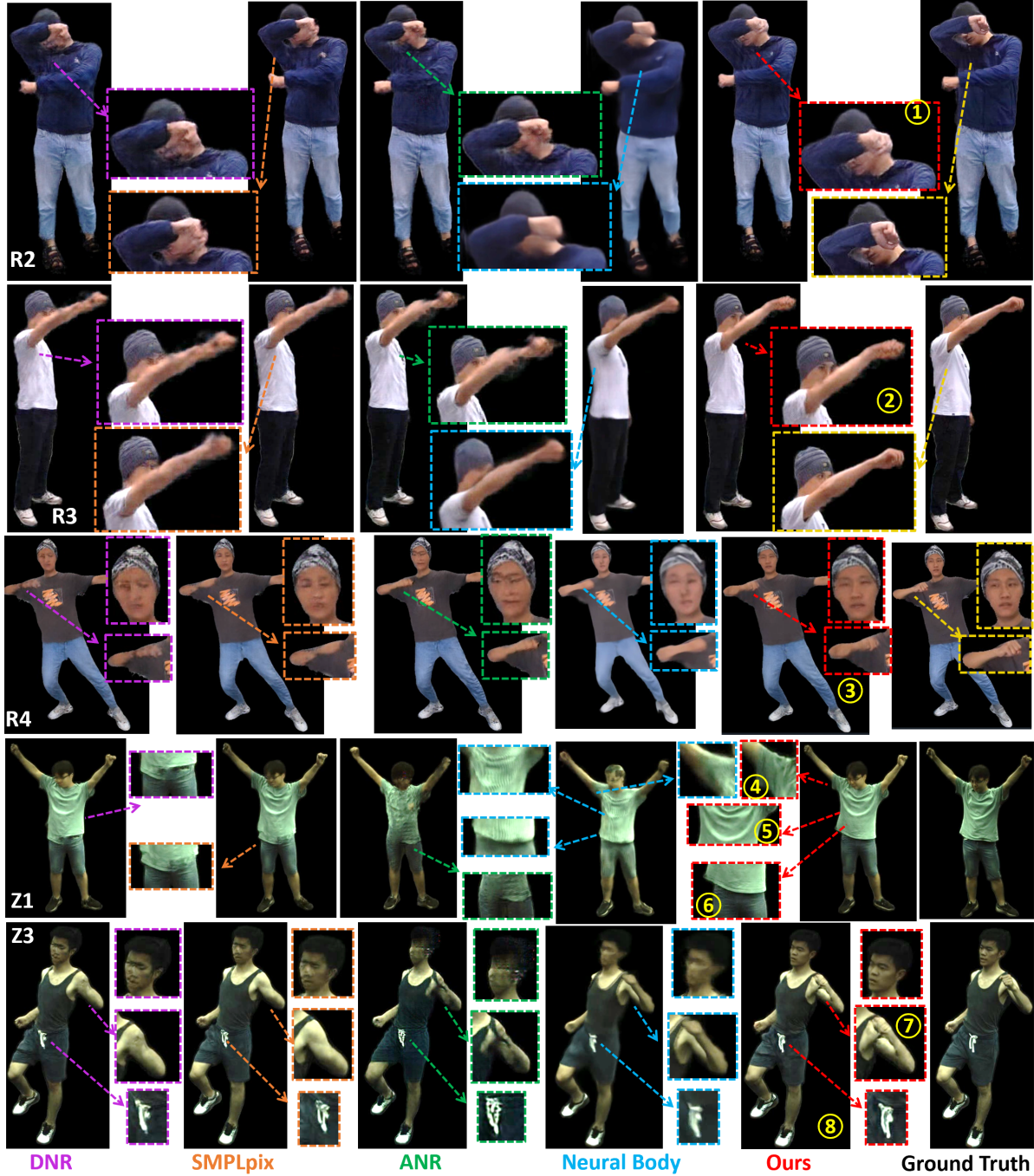


Figure 7: Comparisons with NIT methods (DNR[66], SMPLpix [52], ANR [54]), and a NeRF method (Neural Body [51]) on R2-4, Z1, and Z3. Our method can generate different levels of pose-dependent details: ⑥ offsets, ⑤ big wrinkles, ④ tiny wrinkles. We handle self-occlusions better (①②③⑦) compared to NIT methods, generates high-quality details (④⑤⑧), and preserves thin parts (③⑦) and facial details better. All the poses are novel, and R4, Z1, Z3 are novel views. Note that we cannot guarantee our reproduced ANR achieves the expected performance as stated in Sec. 4.

the IEEE/CVF International Conference on Computer Vision (ICCV), pages 14528–14538, October 2021. 1, 2, 3, 4, 14

- [18] Jingwei Huang, Justus Thies, Angela Dai, Abhijit Kundu, Chiyu Max Jiang, Leonidas J. Guibas, Matthias Nießner, and Thomas A. Funkhouser. Adversarial texture optimization

from rgb-d scans. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 1556–1565, 2020. 2

- [19] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. Arch: Animatable reconstruction of clothed hu-

- mans. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3090–3099, 2020. 2
- [20] Timothy Jeruzalski, Boyang Deng, Mohammad Norouzi, J. P. Lewis, Geoffrey E. Hinton, and Andrea Tagliasacchi. Nasa: Neural articulated shape approximation. *ArXiv*, abs/1912.03207, 2020. 2
- [21] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. volume 9906, pages 694–711, 10 2016. 5
- [22] James T. Kajiya and Brian Von Herzen. Ray tracing volume densities. *Proceedings of the 11th annual conference on Computer graphics and interactive techniques*, 1984. 2, 5
- [23] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2015. 5, 13
- [24] Bernhard Kratzwald, Zhiwu Huang, Danda Pani Paudel, and Luc Van Gool. Towards an understanding of our world by GANing videos in the wild. *arXiv:1711.11453*, 2017. 2
- [25] Zorah Löhner, Daniel Cremers, and Tony Tung. Deepwrinkles: Accurate and realistic clothing modeling. In *ECCV*, 2018. 1
- [26] Lingjie Liu, Marc Habermann, V. Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. *ArXiv*, abs/2106.02019, 2021. 1, 2, 3, 5, 6, 7
- [27] Lingjie Liu, Weipeng Xu, Marc Habermann, Michael Zollhöfer, Florian Bernard, Hyeonwoo Kim, Wenping Wang, and Christian Theobalt. Neural human video rendering by learning dynamic textures and rendering-to-video translation. *IEEE Transactions on Visualization and Computer Graphics*, PP:1–1, 05 2020. 2
- [28] Lingjie Liu, Weipeng Xu, Michael Zollhoefer, Hyeonwoo Kim, Florian Bernard, Marc Habermann, Wenping Wang, and Christian Theobalt. Neural rendering and reenactment of human actor videos. *ACM Transactions on Graphics (TOG)*, 2019. 2
- [29] Weiyang Liu, Y. Wen, Zhiding Yu, Ming Li, B. Raj, and Le Song. Sphereface: Deep hypersphere embedding for face recognition. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6738–6746, 2017. 5
- [30] Stephen Lombardi, Tomas Simon, Jason M. Saragih, Gabriel Schwartz, Andreas M. Lehrmann, and Yaser Sheikh. Neural volumes. *ACM Transactions on Graphics (TOG)*, 38:1 – 14, 2019. 1, 2, 4
- [31] M. Loper, Naureen Mahmood, J. Romero, Gerard Pons-Moll, and Michael J. Black. Smpl: a skinned multi-person linear model. *ACM Trans. Graph.*, 34:248:1–248:16, 2015. 1, 3, 13
- [32] Liqian Ma, Xu Jia, Qianru Sun, Bernt Schiele, Tinne Tuytelaars, and Luc Van Gool. Pose guided person image generation. In *Advances in Neural Information Processing Systems*, pages 405–415, 2017. 2
- [33] Liqian Ma, Qianru Sun, Stamatios Georgoulis, Luc van Gool, Bernt Schiele, and Mario Fritz. Disentangled person image generation. *Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [34] Qianli Ma, Shunsuke Saito, Jinlong Yang, Siyu Tang, and Michael J. Black. Scale: Modeling clothed humans with a surface codec of articulated local elements. In *CVPR*, 2021. 2, 4
- [35] Qianli Ma, Jinlong Yang, Anurag Ranjan, Sergi Pujades, Gerard Pons-Moll, Siyu Tang, and Michael J. Black. Learning to dress 3d people in generative clothing. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6468–6477, 2020. 1
- [36] Qianli Ma, Jinlong Yang, Siyu Tang, and Michael J. Black. The power of points for modeling humans in clothing. *ArXiv*, abs/2109.01137, 2021. 2, 4
- [37] Lars M. Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4455–4465, 2019. 2
- [38] Mateusz Michalkiewicz, Jhony Kaesemodel Pontes, Dominic Jack, Mahsa Baktash, and Anders P. Eriksson. Deep level sets: Implicit surface representations for 3d shape inference. *ArXiv*, abs/1901.06802, 2019. 2
- [39] Marko Mihajlović, Yan Zhang, Michael J. Black, and Siyu Tang. Leap: Learning articulated occupancy of people. *ArXiv*, abs/2104.06849, 2021. 2
- [40] 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 *ECCV*, 2020. 2
- [41] 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 *ECCV*, 2020. 2
- [42] Natalia Neverova, Riza Alp Güler, and Iasonas Kokkinos. Dense pose transfer. *European Conference on Computer Vision (ECCV)*, 2018. 2
- [43] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11448–11459, 2021. 3
- [44] Atsuhiko Noguchi, Xiao Sun, Stephen Ching-Feng Lin, and Tatsuya Harada. Neural articulated radiance field. *ArXiv*, abs/2104.03110, 2021. 3
- [45] Pablo Rodríguez Palafox, Aljazz Bovzirc, Justus Thies, Matthias Nießner, and Angela Dai. Npms: Neural parametric models for 3d deformable shapes. *ArXiv*, abs/2104.00702, 2021. 2
- [46] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. *International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 13
- [47] Jeong Joon Park, Peter R. Florence, Julian Straub, Richard A. Newcombe, and S. Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 165–174, 2019. 2

- [48] Chaitanya Patel, Zhouyingcheng Liao, and Gerard Pons-Moll. Tailornet: Predicting clothing in 3d as a function of human pose, shape and garment style. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7363–7373, 2020. 1
- [49] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed 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)*, 2019. 3
- [50] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. 2021. 1, 2, 3, 6, 13
- [51] 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. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9050–9059, 2021. 1, 2, 3, 5, 6, 9, 13
- [52] Sergey Prokudin, Michael J. Black, and Javier Romero. Smpix: Neural avatars from 3d human models. *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1809–1818, 2021. 1, 2, 3, 6, 9, 13
- [53] Albert Pumarola, Antonio Agudo, Alberto Sanfeliu, and Francesc Moreno-Noguer. Unsupervised person image synthesis in arbitrary poses. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2
- [54] Amit Raj, Julian Tanke, James Hays, Minh Vo, Carsten Stoll, and Christoph Lassner. Anr: Articulated neural rendering for virtual avatars. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3721–3730, 2021. 1, 2, 3, 4, 6, 9, 13
- [55] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501*, 2020. 8, 13
- [56] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2304–2314, 2019. 2
- [57] Shunsuke Saito, Tomas Simon, Jason M. Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 81–90, 2020. 2
- [58] Shunsuke Saito, Jinlong Yang, Qianli Ma, and Michael J. Black. Scanimate: Weakly supervised learning of skinned clothed avatar networks. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2885–2896, 2021. 2
- [59] Kripasindhu Sarkar, Dushyant Mehta, Weipeng Xu, Vladislav Golyanik, and Christian Theobalt. Neural re-rendering of humans from a single image. In *European Conference on Computer Vision (ECCV)*, 2020. 2, 14
- [60] Aliaksandr Siarohin, Enver Sangineto, Stephane Lathuiliere, and Nicu Sebe. Deformable GANs for pose-based human image generation. In *CVPR 2018*, 2018. 2
- [61] K. Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2015. 5
- [62] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhöfer. Deepvoxels: Learning persistent 3d feature embeddings. In *Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [63] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2
- [64] Shih-Yang Su, Frank Yu, Michael Zollhoefer, and Helge Rhodin. A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose. 2021. 3
- [65] Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason M. Saragih, Matthias Nießner, Rohit Pandey, S. Fanello, Gordon Wetzstein, Jun-Yan Zhu, Christian Theobalt, Maneesh Agrawala, Eli Shechtman, Dan B. Goldman, and Michael Zollhofer. State of the art on neural rendering. *Computer Graphics Forum*, 39, 2020. 1
- [66] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 38, 2019. 2, 3, 4, 6, 9, 13
- [67] Garvita Tiwari, Nikolaos Sarafianos, Tony Tung, and Gerard Pons-Moll. Neural-gif: Neural generalized implicit functions for animating people in clothing. *ArXiv*, abs/2108.08807, 2021. 2
- [68] Shaofei Wang, Marko Mihajlović, Qianli Ma, Andreas Geiger, and Siyu Tang. Metaavatar: Learning animatable clothed human models from few depth images. *ArXiv*, abs/2106.11944, 2021. 2
- [69] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *NeurIPS*, 2018. 2
- [70] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 2
- [71] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. pages 8798–8807, 06 2018. 5, 13
- [72] Zhou Wang, A. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13:600–612, 2004. 6
- [73] Feng Xu, Yebin Liu, Carsten Stoll, James Tompkin, Gaurav Bharaj, Qionghai Dai, Hans-Peter Seidel, Jan Kautz, and Christian Theobalt. Video-based characters: creating new human performances from a multi-view video database. *ACM SIGGRAPH 2011 papers*, 2011. 1

- [74] Hongyi Xu, Thiemo Alldieck, and Cristian Sminchisescu. H-nerf: Neural radiance fields for rendering and temporal reconstruction of humans in motion. *ArXiv*, abs/2110.13746, 2021. 6
- [75] Jinlong Yang, Jean-Sébastien Franco, Franck Hétroy-Wheeler, and Stefanie Wuhler. Analyzing clothing layer deformation statistics of 3d human motions. In *ECCV*, 2018. 1
- [76] Richard Zhang, Phillip Isola, Alexei A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 586–595, 2018. 6, 15
- [77] Yang Zheng, Ruizhi Shao, Yuxiang Zhang, Tao Yu, Zerong Zheng, Qionghai Dai, and Yebin Liu. Deepmulticap: Performance capture of multiple characters using sparse multiview cameras. *ArXiv*, abs/2105.00261, 2021. 2
- [78] Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai. Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. *IEEE transactions on pattern analysis and machine intelligence*, PP, 2021. 2
- [79] Zhen Zhu, Tengpeng Huang, Baoguang Shi, Miao Yu, Bofei Wang, and Xiang Bai. Progressive pose attention transfer for person image generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2347–2356, 2019. 2

## Appendices

### A. Implementation Details

**Optimization.** The networks were trained using the Adam optimizer [23] with an initial learning rate of  $2 \times 10^{-4}$ ,  $\beta_1 = 0.5$ . The loss weights  $\{\lambda_{vol}, \lambda_n, \lambda_{feat}, \lambda_{mask}, \lambda_{pix}, \lambda_{adv}, \lambda_{face}\}$  are set empirically to  $\{15, 1, 10, 5, 1, 1, 5\}$ . We train DNR[66], SMPLpix[52], ANR[54], and our method for 50,000 iterations, and 180,000 iterations for Neural Body [51], and 250,000 iterations for Animatable NeRF (AniNeRF [50]). We train the networks with a Nvidia P6000 GPU, and it generally takes 28 hours for DNR and SMPLpix, and 40 hours for our method. Note that Neural Body [51] cannot converge to a detailed generation as seen in Fig. 11.

**Network Architectures and Optimizable Latents.**  $Z_G$  and  $Z_T$  both have a size of  $128 \times 128 \times 16$ .  $F_P$  is based on Pix2PixHD [71] architecture with Encoder blocks of [Conv2d, Batch-Norm, ReLU], ResNet [15] blocks, and Decoder blocks of [ReLU, ConvTranspose2d, BatchNorm].  $F_P$  has 3 Encoder and Decode blocks, and 2 ResNet blocks.  $F_N$  has 2 Decode blocks.  $F_T$  has  $n$  ( $n = 2$  or  $3$  or  $4$ ) Encoder blocks, and the exact number depends on the downsampling factor of PD-NeRF such that the textural features and volumetric features have the same size as discussed at Sec. 3.4.  $F_R$  has  $(4 - n)$  Encoder blocks, 4 Decoder blocks, and 5 ResNet blocks. For  $F_\theta$ , we use a 7-layer MLP with a skip connection from the input to the 4th layer as in DeepSDF [46]. From the 5th layer, the network branches out two heads, one to predict density with one fully-connected layer and the other one to predict color features with two fully-connected layers.

**Geometry-guided Ray Marching.** The success of our method depends on the efficient and effective training of the pose-conditioned downsampled NeRF (PD-NeRF). First, instead of sampling rays in the whole space, we utilize a geometry-guided ray marching mesh as illustrated in Fig. 8. Specifically, we only sample query points along the corresponding rays near the SMPL [31] mesh surface, which is determined by a dilated SMPL mesh. The SMPL mesh is dilated along the normal of each face with a radius of  $d$ , where  $d$  is about 12cm for general clothes and 20cm for loose clothing like skirts for M1 dataset (see Fig. 5 of the paper). We find the near and far points by querying the Z-buffer of the corresponding pixels after projecting the dilated SMPL mesh using Pytorch3D [55]. In addition, we sample more points to the near region, which is expected to contain visible contents. The geometry-guided ray marching algorithm and UV conditioned architecture enable us to train a PD-NeRF with  $45 \times 45$  resolution images and only 7 sampled point along each ray, as shown in Fig. 9. Though learned from low resolution images, the reconstructed geometry still preserves some pose-dependent features.

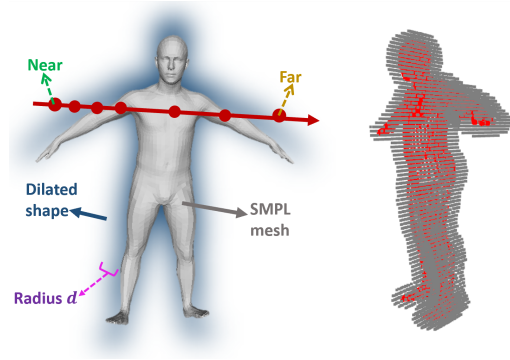


Figure 8: Geometry-guided ray marching. Left: sampling points by SMPL mesh dilation. Right: Red - SMPL model; Gray - rays and sampled points.



Figure 9: Construct PD-NeRF with  $45 \times 45$  resolution images and 7 sampled point along each ray: left (geometry), right (reference image).

### B. More Experimental Results

	LPIPS ↓	FID ↓	SSIM ↑	PSNR ↑
AniNeRF [50]	.271	196.44	.773	23.36
Ours	<b>.090</b>	<b>62.33</b>	<b>.842</b>	<b>26.22</b>

Table 5: Comparisons with AniNeRF on R1 dataset under novel poses. To reduce the influence of the background, all scores are calculated from images cropped to 2D bounding boxes.

#### B.1. Comparisons

**Comparisons with Animatable NeRF [50]** A quantitative comparison with AniNeRF on R1 dataset is shown in Tab. 5, and the results of the other methods are shown in Tab. 2 and Tab. 4 of the paper. Our method significantly outperforms AniNeRF on all the four metrics. The qualitative comparison is shown in Fig. 12 and the supplementary video.



Figure 10: Comparisons of methods trained with or without face identity loss. Ours(-F) indicates a variant of our method that is not trained with face identity loss.

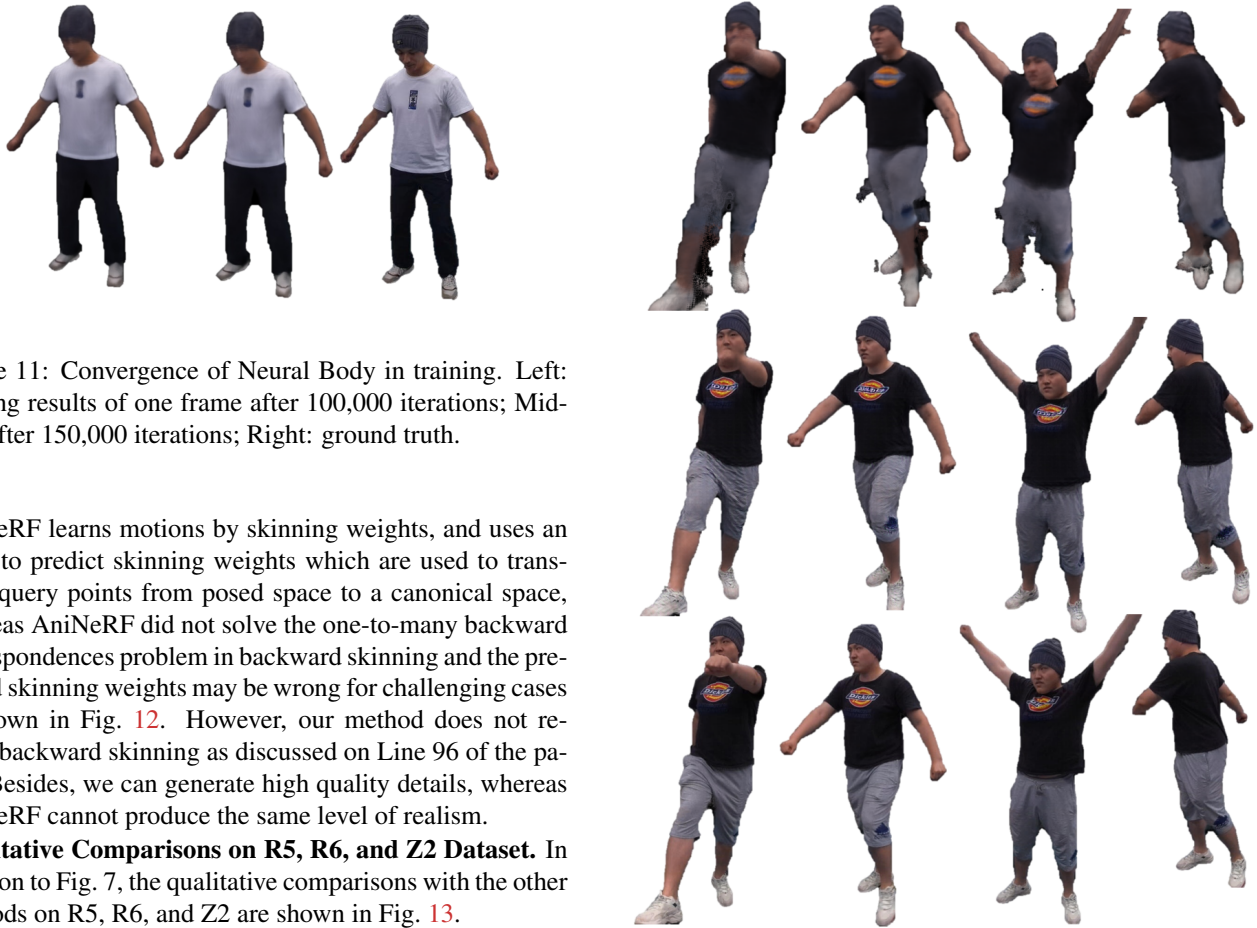


Figure 11: Convergence of Neural Body in training. Left: training results of one frame after 100,000 iterations; Middle: after 150,000 iterations; Right: ground truth.

AniNeRF learns motions by skinning weights, and uses an MLP to predict skinning weights which are used to transform query points from posed space to a canonical space, whereas AniNeRF did not solve the one-to-many backward correspondences problem in backward skinning and the predicted skinning weights may be wrong for challenging cases as shown in Fig. 12. However, our method does not require backward skinning as discussed on Line 96 of the paper. Besides, we can generate high quality details, whereas AniNeRF cannot produce the same level of realism.

**Qualitative Comparisons on R5, R6, and Z2 Dataset.** In addition to Fig. 7, the qualitative comparisons with the other methods on R5, R6, and Z2 are shown in Fig. 13.

**Accuracy and Inference Time.** The accuracy and inference time of each method are shown in Tab. 6.

## B.2. Ablation Study

**Face Identity Loss.** We use the face identity loss to improve the qualitative results as shown in Fig. 10 (also used in [17, 59]), whereas the improvements of faces do not im-

Figure 12: Qualitative comparisons with AniNeRF. Top: AniNeRF; Middle: ours; Bottom: ground truth.

prove the overall quantitative results of each method, as listed in Tab. 8.

Models	LPIPS↓	FID↓	SSIM↑	PSNR↑	Time (s)	VR.T(%)
DNR	0.1023	75.0152	0.8310	25.7303	0.184	-
SMPLpix	0.1002	69.8119	0.8350	25.9295	0.198	-
ANR	0.1172	78.5012	0.8301	26.0168	0.224	-
Neuray Body	0.2124	155.8382	0.8328	26.1718	18.200	-
Ours_ng	0.0991	70.5282	0.8370	25.9842	0.257	-
Ours_16(7)	0.0966	64.8711	0.8489	26.4356	0.292	11.99
Ours_16(12)	0.0959	63.7922	0.8494	26.4822	0.295	12.88
Ours_16(20)	0.0959	63.8337	<b>0.8495</b>	<b>26.4841</b>	0.305	15.41
Ours_8(12)	0.0901	62.3330	0.8415	26.2165	0.349	26.36
Ours_4(20)	<b>0.0861</b>	<b>60.7884</b>	0.8461	26.2465	0.464	44.61

Table 6: Performance, inference time of each methods. VR.T(%) indicates the percentages of the volume rendering time. Compared with Tab. 4 of the paper, the other two metrics SSIM and PSNR are included.

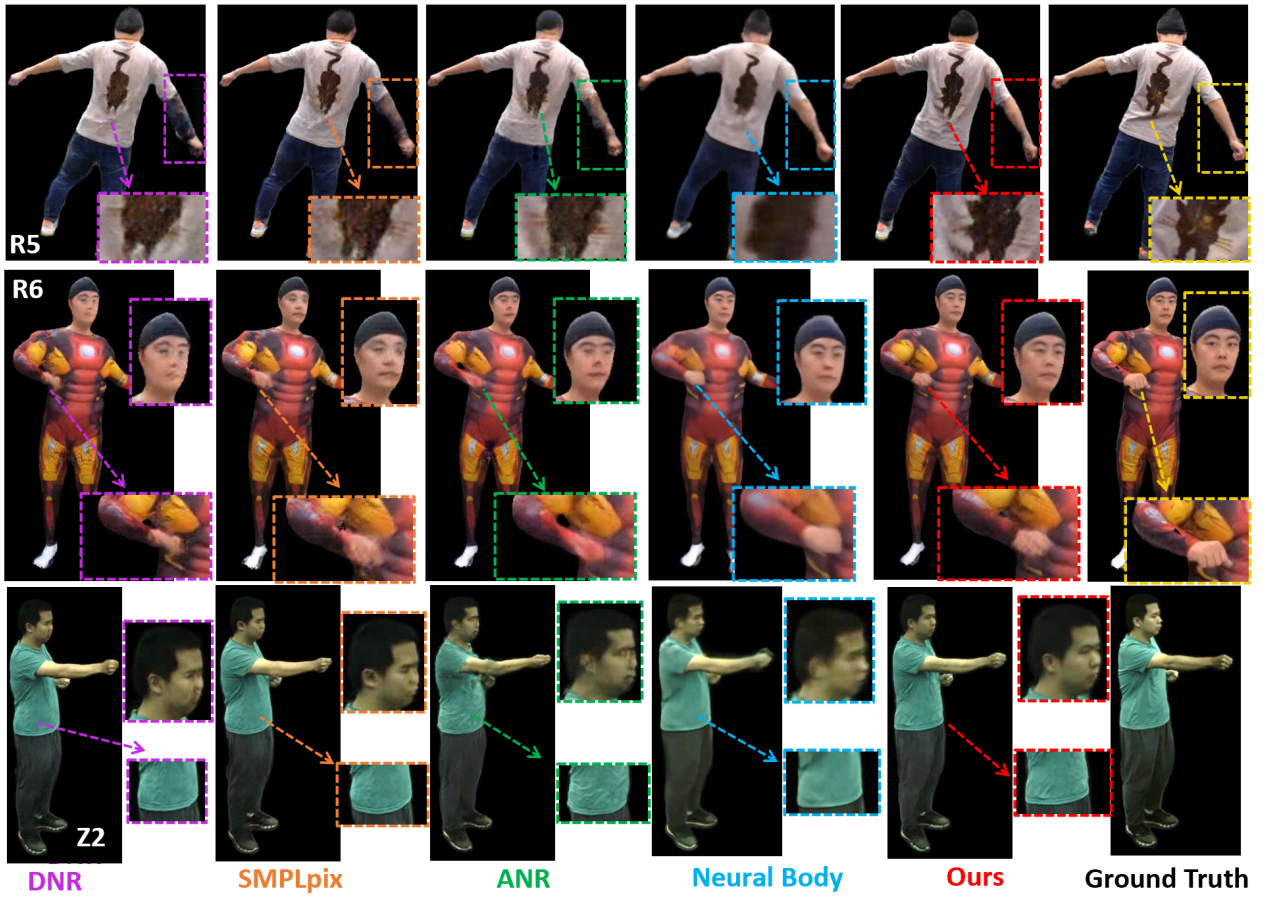


Figure 13: Comparisons with the other methods on R5, R6, and Z2 Dataset.

**Feature Fusion.** We compare two methods to fuse the volumetric and textural features as discussed at Sec. 3.4 by concatenation (Concat) and Attentional Feature Fusion (AFF [8]) on two datasets, R1 and R2 (about 12,000 frames in training, 3,000 frames in testing). We test the performances on novel poses. The quantitative results show that AFF can

improve the LPIPS [76] and FID [16] results.

<b>R2</b>	LPIPS ↓	FID ↓	SSIM↑	PSNR↑
Concat	.117	90.43	<b>.838</b>	28.55
AFF	<b>.108</b>	<b>84.43</b>	.833	<b>28.62</b>
<b>R1</b>	LPIPS ↓	FID ↓	SSIM↑	PSNR↑
Concat	.099	64.54	<b>.856</b>	<b>26.78</b>
AFF	<b>.090</b>	<b>62.33</b>	.842	26.22

Table 7: Comparisons of fusing volumetric and textural features by concatenation (Concat) and Attentional Feature Fusion (AFF [8]) on R1 and R2 dataset.

<b>R3</b>	LPIPS ↓	FID ↓	SSIM↑	PSNR↑
DNR	.108	80.33	.809	24.05
DNR + F	<b>.103</b>	<b>75.42</b>	<b>.812</b>	<b>24.13</b>
SMPLPix	<b>.104</b>	<b>74.57</b>	.810	<b>24.16</b>
SMPLPix + F	.109	78.59	<b>.811</b>	23.93
<b>R2</b>	LPIPS ↓	FID ↓	SSIM↑	PSNR↑
DNR	.128	105.63	<b>.820</b>	27.82
DNR + F	<b>.125</b>	<b>98.86</b>	<b>.820</b>	<b>27.92</b>
SMPLPix	<b>.124</b>	99.81	.822	27.92
SMPLPix + F	<b>.124</b>	<b>94.81</b>	<b>.826</b>	<b>27.97</b>
<b>R1</b>	LPIPS ↓	FID ↓	SSIM↑	PSNR↑
DNR	<b>.102</b>	<b>75.02</b>	.831	25.73
DNR + F	.103	75.35	<b>.832</b>	<b>25.85</b>
SMPLPix	<b>.100</b>	<b>69.81</b>	<b>.835</b>	<b>25.93</b>
SMPLPix + F	.104	75.34	.833	25.81
Ours(-F)	<b>.088</b>	<b>61.03</b>	.841	<b>26.42</b>
Ours	.090	62.33	<b>.842</b>	26.22

Table 8: Quantitative results of each method trained with or without face identity loss. Ours(-F) indicates a variant of our method that is not trained with face identity loss. (DNR + F) and (SMPLpix + F) are trained with face identity loss.