

# Generalizable Neural Voxels for Fast Human Radiance Fields

Taoran Yi<sup>1\*</sup>, Jiemin Fang<sup>2,1\*</sup>, Xinggang Wang<sup>1†</sup>, Wenyu Liu<sup>1</sup>

<sup>1</sup>School of EIC, Huazhong University of Science & Technology

<sup>2</sup>Institute of Artificial Intelligence, Huazhong University of Science & Technology

{taoranyi, jaminfang, xgwang, liuwuy}@hust.edu.cn

## Abstract

*Rendering moving human bodies at free viewpoints only from a monocular video is quite a challenging problem. The information is too sparse to model complicated human body structures and motions from both view and pose dimensions. Neural radiance fields (NeRF) have shown great power in novel view synthesis and have been applied to human body rendering. However, most current NeRF-based methods bear huge costs for both training and rendering, which impedes the wide applications in real-life scenarios. In this paper, we propose a rendering framework that can learn moving human body structures extremely quickly from a monocular video. The framework is built by integrating both neural fields and neural voxels. Especially, a set of generalizable neural voxels are constructed. With pretrained on various human bodies, these general voxels represent a basic skeleton and can provide strong geometric priors. For the fine-tuning process, individual voxels are constructed for learning differential textures, complementary to general voxels. Thus learning a novel body can be further accelerated, taking only a few minutes. Our method shows significantly higher training efficiency compared with previous methods, while maintaining similar rendering quality. The project page is at <https://taoranyi.com/gneuvox>.*

## 1. Introduction

Rendering human bodies [84, 85, 67, 25, 82, 81, 68, 52, 44, 9, 4] is a longstanding research topic and plays important roles in varieties of applications, *e.g.*, virtual reality (VR), augmented reality (AR) and other interactive products. Recently, the emergence of neural radiance fields (NeRF) [50] has significantly facilitated the development of rendering techniques. Taking sparse-view images as input, NeRF models can generate images from novel views with

high qualities.

Some works [60, 86, 72, 12, 59, 58, 15, 33, 89] successfully apply NeRF methods to human body rendering frameworks. Benefiting from geometry priors, *e.g.* SMPL models [43] for predicting human poses, body structures can be learned accurately even for poses with complicated motions. However, NeRF-based models rely on volume rendering to connect 2D image pixels with 3D real points, which requires for inferring tons of points and takes huge cost. The training process for one single scene usually needs dozens of hours or even days to complete on one GPU. This will undoubtedly impede these methods from applications in real-life scenarios.

In this paper, a fast training framework is built for rendering free-view images of human bodies from a monocular video recording a moving person. We first propose to represent the human body under the canonical pose with neural voxel features, which can be optimized via gradient descent directly. This explicit representation significantly accelerates the optimization process. Equipped with deformation networks that transform human poses into canonical ones, human body information can be modeled accurately and quickly. More importantly, considering different human bodies share commonalities, *e.g.* the basic skeleton, we enable the generalization ability of the proposed framework and name our method as **GNeuVox**. Two types of neural voxels are constructed to achieve this goal. One as **general voxels** is pretrained across various human bodies where common properties/information are learned and stored. The other one as **individual voxels** is built for any novel scene, which learns the scene’s specific appearance and unique textures. With the collaboration of two types of neural voxels, the whole framework can learn a new human body extremely quickly and render high-quality images.

We summarize contributions as follows.

- A fast-training framework is designed by introducing explicit neural voxels, which renders free-view moving human bodies from a monocular video.
- As far as we know, we are the first to propose **gener-**

\*Equal contributions.

†Corresponding author.

**alizable explicit representations** to model common properties across different human bodies, which show significant effectiveness in accelerating novel human learning.

- The proposed framework is evaluated to show notably high efficiency with only **5-minute** training time. Even across different datasets, a strong generalization ability is still achieved. The code will be released.

## 2. Related Work

### 2.1. Neural Rendering for Human Bodies

Differentiable neural rendering has been widely adopted in human body scenarios which produces highly realistic images. Neural Body [60] proposes to compute latent code volumes by inferring mesh vertices but performs poorly on unseen poses. To solve this problem, Animatable NeRF [58, 59] maps human poses from the observation space to the predefined canonical space. Some other methods [83, 40] share similar ideas, which are not limited to the human body, but also apply to other dynamic scenes [61, 55, 54, 18, 32, 24]. In ARAH [80] and SNARF [15], a more accurate joint root-finding algorithm is used to obtain a more accurate mapping, which can better generalize to poses beyond the existing distribution. Neuman [33] separates the background and characters, which are then modeled separately. There are also other methods to separate the scene [71, 38], which can lead to better rendering results. A-NeRF [72] uses 3D skeleton poses to model the human body. In H-NeRF [89], the combination of NeRF and Signed Distance Field (SDF) also obtains good human rendering results. PIFu [65] and PI-HuHD [66] use the supervision of 3D mesh to model the human body. [14, 7, 79] achieve the effect of generating humans through additional supervision. ARCH [29], ARCH++ [26] and ICON [88] also require additional 3D supervision, but the cost of supervision data acquisition is expensive. Our method can achieve free-viewpoint rendering with only monocular videos, which greatly reduces the difficulty of dataset acquisition. Compared with other methods using monocular videos [38, 29, 31, 3, 2, 20, 75, 21, 87], ours also enjoys high convergence efficiency.

### 2.2. Generalizable NeRF

Generalization is an important but challenging problem in NeRF, as a NeRF model can usually only represent one specific scene. Making NeRF generalizable will greatly improve the efficiency of the representation. PixelNeRF [94] projects images onto generalizable feature volumes. Inspired by MVSNet [91], MVSNeRF [11] builds a voxel feature containing relative positions and directions for view synthesis. These methods, including some follow-up methods [76, 37, 64, 78], extract features from images of views

near the target view, so as to achieve the purpose of generalization. However, due to occlusion present in some poses, artifacts may occur in certain positions. For this reason, GeoNeRF [34] and NeuRay [41] gather information from consistent source views to reduce artifacts. Some works also explore NeRF generalization on human bodies [98, 36, 16, 22, 48, 13], but the convergence speed is slow.

### 2.3. NeRF Acceleration

NeRF [50] and its extensions [5, 96, 46, 28, 6, 49, 45, 42, 84] have achieved high rendering quality, but most of them bear long inference time and take large training cost to learn the scene. For inference speed, it can be promoted by improving the sampling strategy or baking related properties [93, 62, 27, 23, 77]. DVGO [73], Plenoxels [92], Instant-NGP [51] and other methods [53, 70, 90, 63] significantly speed up the convergence speed by introducing explicit representations. DS-NeRF [17] proposes to accelerate the convergence with depth supervision. TensorRF [10] uses tensor decomposition to achieve high storage efficiency while maintaining fast training speed. However, the aforementioned ones can only be applied to static scenes. [18, 39, 19, 57] extend the efficient NeRF framework to dynamic scenes. We not only introduce explicit representations for acceleration but also make them generalizable for extremely fast convergence.

## 3. Method

In this section, we first review the original NeRF [50] and HumanNeRF [86] methods in Sec. 3.1. Then in Sec. 3.2, the proposed framework for fast-learning human bodies is introduced. We discuss how to generalize the framework across different scenes in Sec. 3.3. Finally, in Sec. 3.4, we elaborate the optimization process.

### 3.1. Preliminaries

Neural radiance fields (NeRF) are first proposed in [50], aiming at connecting real 3D points with 2D pixels of images. NeRF is designed as a 5D function  $f$ , mapping 3D coordinates  $(x, y, z)$  along with 2D directions  $\mathbf{d} = (\theta, \phi)$  into the color  $c$  and density  $\sigma$ . This process can be formulated as:

$$c, \sigma = f(x, y, z, \mathbf{d}), \quad (1)$$

where  $f$  is commonly instantiated as multilayer perceptrons (MLPs).

In order to obtain the pixel color  $C(\mathbf{r})$ , points are sampled along the ray  $\mathbf{r} = \mathbf{o} + t\mathbf{d}$  emitted from the camera, where  $\mathbf{o}$  denotes the origin of the ray, and  $\mathbf{d}$  denotes the ray direction. Then the differentiable volume rendering [47] is performed to accumulate all 3D point colors and densities

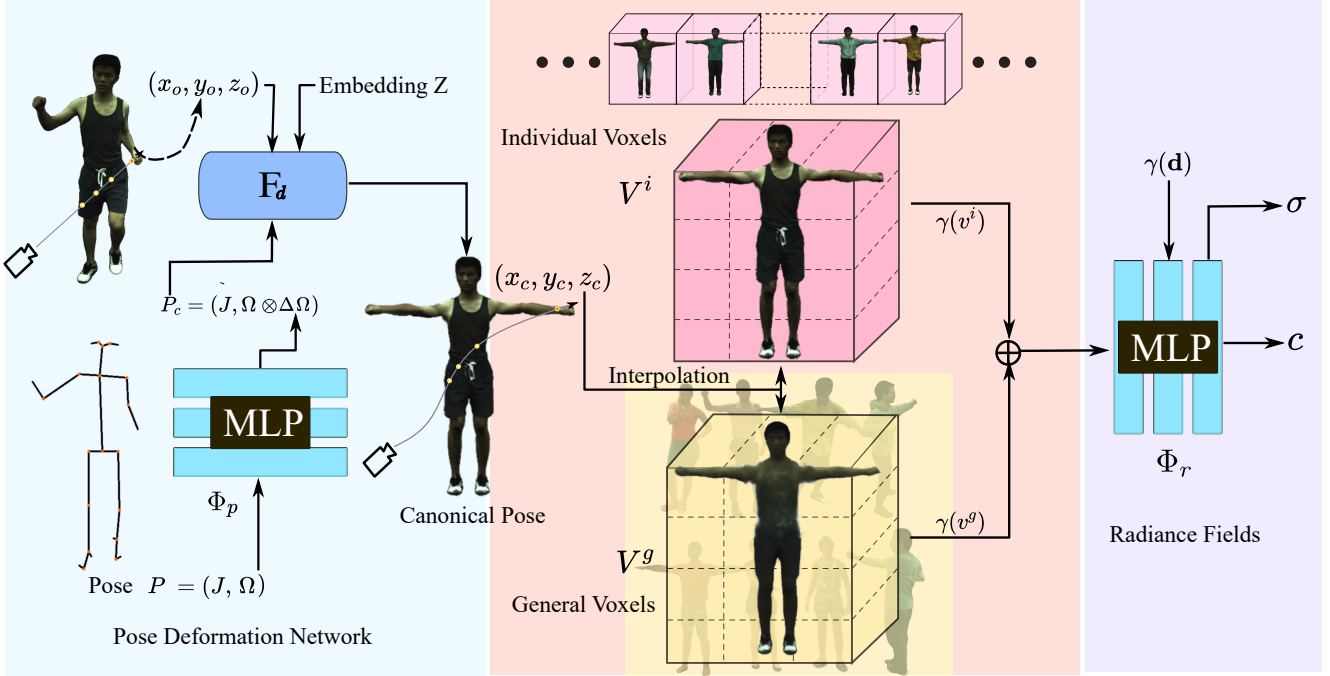


Figure 1. Overall framework of GNeuVox. First, we use the pose refinement module to adjust the obtained pose  $P$ , and then use the corrected pose  $P_c$  and the embedding  $Z$  to deform the coordinates  $(x_o, y_o, z_o)$  of the human body in the observation space to the T-pose in the canonical space. Using the shifted coordinates  $(x_c, y_c, z_c)$  to interpolate the individual voxels  $V^i$  and general voxels  $V^g$  to obtain their corresponding feature vectors  $v^i$  and  $v^g$ . Finally,  $v^i$  and  $v^g$  are fed into the radiance fields to obtain the color  $c$  and density  $\sigma$ .

along the ray and produce the final pixel color:

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) c_i, \quad (2)$$

where  $T_i = \exp(-\sum_{j=1}^{i-1} \sigma_j \delta_j)$ ,  $\delta_i = t_{i+1} - t_i$  denotes the distance between two adjacent points,  $N$  denotes the number of sampled points. To represent high-frequency details [74],  $(x, y, z)$  and  $\mathbf{d}$  in Eq. 1 are transformed with the positional encoding  $\gamma$ , which is formulated as

$$\gamma(x) = (\sin(2^0 x), \cos(2^0 x), \dots, \sin(2^{L-1} x), \cos(2^{L-1} x)), \quad (3)$$

where  $L$  is used to control the maximum frequency.

Parameters of the field model are optimized to make the rendered color  $\hat{C}(\mathbf{r})$  approach the ground truth pixel color  $C(\mathbf{r})$ . The loss function is defined as

$$\mathcal{L} = \|\hat{C}(\mathbf{r}) - C(\mathbf{r})\|_2^2. \quad (4)$$

HumanNeRF [86] is a representative work of extending NeRF to human bodies. It uses the NeRF model to represent a human body with the canonical pose. To achieve this, a pose refinement module and a coordinate deformation network are introduced to transform the human pose into the canonical one. The pose refinement module corrects the pose  $P$  obtained from the off-the-shelf technique [43] to

obtain a more accurate one  $P_c$ . The coordinate deformation network maps the observation pose  $(x_o, y_o, z_o)$  to the canonical one  $(x_c, y_c, z_c)$  with the corrected pose. The human body pose can be represented as  $P = (J, \Omega)$ ,  $J$  denotes the positions of  $K$  human body key points and  $\Omega$  denotes the rotation angles of local joints. The pose refinement network is denoted as  $\Delta\Omega = \Phi_p(\Omega)$ , which is performed to get a more accurate pose  $P_c = (J_c, \Omega_c)$ . It can be formulated as

$$P_c = (J_c, \Omega_c) = (J, \Omega \otimes \Delta\Omega), \quad (5)$$

where  $\otimes$  denotes  $(R^i \cdot \Delta R^i, T^i + \Delta T^i)$ ,  $R^i$  and  $T^i$  denote the rotation matrix and translation matrix to the canonical space corresponding to the  $i$ -th bone.  $(R^i, T^i)$  and  $(\Delta R^i, \Delta T^i)$  are derived from  $\Omega$  and  $\Delta\Omega$  respectively.

The coordinate deformation network<sup>1</sup> is formulated as:

$$(x_{skel}, y_{skel}, z_{skel}) = T_{skel}((x_o, y_o, z_o), P_c), \quad (6)$$

$$T_{skel}(x, P_c) = \sum_{i=1}^K w^i(x) (R_c^i x + T_c^i),$$

where  $(x_o, y_o, z_o)$  is denoted as  $x$  for abbreviation and  $(R_c^i, T_c^i)$  is derived from  $\Omega_c$ ,  $w^i$  is a weight parameter that controls the blend weight for the  $i$ -th bone.  $w^i$  is obtained

<sup>1</sup>The deformation network in HumanNeRF [86] contains two parts of rigid skeletal offsets and non-rigid offsets. Here we only refer to skeletal offsets for simplicity.

by mapping a learnable embedding vector  $Z$  into a volume via several 3D convolutional layers, which is further trilinearly interpolated.

### 3.2. Fast-training Framework for Rendering Human Bodies

Our framework is constructed with the following modules, *i.e.* the pose deformation network, neural voxels, and radiance fields. We show the overall framework of our method in Fig. 1.

**Overall Pipeline** We first use the pose refinement module to correct the obtained pose  $P$ , and use the corrected pose  $P_c$  with the embedding  $Z$  to deform the coordinates  $(x_o, y_o, z_o)$  of the human body in the observation space to the T-pose in the canonical space through the pose deformation network  $\mathbf{F}_d$ . Then we use the deformed coordinates  $(x_c, y_c, z_c)$  in the canonical space to interpolate the neural voxels  $V$  to get the voxel feature  $\mathbf{v}$ . Finally,  $\mathbf{v}$  together with the coordinates  $(x_o, y_o, z_o)$ , direction  $\mathbf{d}$ , and timestamp  $t$  are fed into the radiation field to obtain the color  $\mathbf{c}$  and density  $\sigma$ .

**Pose Deformation Network** We construct the pose refinement network and coordinate deformation network, following HumanNeRF [86], to model human motions. The process of mapping the pose in the observation space  $(x_o, y_o, z_o)$  to the canonical space  $(x_c, y_c, z_c)$  is formulated as

$$(x_c, y_c, z_c) = \mathbf{F}_d((x_o, y_o, z_o), P_c), \quad (7)$$

where  $P_c$  is the human body pose after pose refinement network  $\Phi_p$ , which is the same as Eq.5.  $\mathbf{F}_d$  has the same structure as  $T_{skel}$  in Eq.6, which is achieved via 3D convolutional layers with corrected poses  $P_c$  and zero-initialized embedding  $Z$ .

#### Representing Canonical Bodies with Neural Voxels

Most previous NeRF-based methods for human bodies [60, 86, 72, 12, 59, 58, 15, 33, 89] adopt purely implicit representations. Although good rendering quality is obtained, it takes a long time to complete the training phase. To speed up the training on human bodies, an explicit data structure is introduced, namely *neural voxels*. We construct a set of optimizable feature parameters which are organized into a voxel grid structure. These voxel features are designed for representing the canonical T-pose human body, which ease the difficulty of optimization. Then trilinear interpolation is performed on *neural voxels*  $V \in \mathbb{R}^{C_v \times N_x \times N_y \times N_z}$  to get the features  $\mathbf{v}$ , which is performed as multi-distance interpolation (MDI) [18].  $C_v$  denotes the channel number of each voxel feature, while  $N_x$ ,  $N_y$ , and  $N_z$  denote the length of each dimension in 3D space.

$$\mathbf{v} = \text{interp}\{(x_c, y_c, z_c), V\}, \quad (8)$$

where  $\text{interp}$  denotes trilinear interpolation operation. Besides, the timestamp  $t$ , direction  $\mathbf{d}$ , original coordinates  $(x_o, y_o, z_o)$  along with interpolated voxel features  $\mathbf{v}$  are all fed into the radiance field  $\Phi_r$  to predict the final color  $\mathbf{c}$  and density  $\sigma$ .

$$\mathbf{c}, \sigma = \Phi_r(\gamma(\mathbf{v}), \gamma(t), \gamma(x_o, y_o, z_o), \gamma(\mathbf{d})), \quad (9)$$

where  $\gamma$  is the positional encoding function as defined in Eq. 3.

### 3.3. Generalization across Human Bodies

Rendering human bodies is a domain-specific task. Though different human bodies own various appearances or wear diverse clothes, they share commonalities in the basic skeleton. To take full advantage of the geometry priors, we propose to enable the generalization ability of the framework. This is achieved by constructing two types of neural voxels introduced in Sec. 3.2. One is designed as **general voxels**  $V^g$  which are pretrained across various human bodies. General voxels extract substantial information from different human bodies, representing a basic skeleton or template applicable for any human body. The other one is constructed instantly as **individual voxels**  $V^i$ , which are optimized to represent unique appearances for a specific human body. The interpolation process, which is same as Eq. 8, is performed on both voxel grids to obtain voxel features  $\mathbf{v}^g$  and  $\mathbf{v}^i$ , coming from  $V^g$  and  $V^i$  respectively. Two voxel features  $\mathbf{v}^g$  and  $\mathbf{v}^i$  are concatenated and fed into the radiance network. The radiance field function in Eq. 9 is re-defined as:  $\mathbf{c}, \sigma = \Phi_r(\gamma(\mathbf{v}^g), \gamma(\mathbf{v}^i), \gamma(t), \gamma(x_o, y_o, z_o), \gamma(\mathbf{d}))$ . With pretrained general voxels, optimization of the rendering framework directly starts from a body template, not from nothing or chaos. Then the training phase can be further accelerated, benefiting from structure priors obtained from various human bodies.

The training process is divided into two phases, *i.e.* one pretraining phase across multiple different human bodies and fine-tuning for a specific human. During the pretraining, besides general voxels, the deformation network and radiance field are also shared for all humans. The two networks can be treated as decoders and the decoding function becomes stronger and more accurate with seeing substantially different humans. On the other hand, the individual voxels and deformation embedding  $Z$  are built separately for each specific human. For the fine-tuning phase, a new set of individual voxels and embedding  $Z$  are built with a zero initialization. Parameters of all modules in the framework are updated to fit the new human body. The fine-tuning phase can be completed with quite few iterations, where mainly differential textures need to be captured.

Table 1. Quantitative comparisons on the ZJU-MoCap dataset [60]. GNeuVox represents training from scratch. GNeuVox-ZJU and GNeuVox-H36m represent the fine-tuning results after pretraining on the ZJU-MoCap and Human3.6M [30] dataset. We color cells with the best values in orange and the second best values in blue.

	Subject 377			Subject 386			Subject 387		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS ( $\times 10^{-2}$ ) $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS ( $\times 10^{-2}$ ) $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS ( $\times 10^{-2}$ ) $\downarrow$
Neural Body [60]	29.11	0.9674	4.095	30.54	0.9678	4.643	27.00	0.9518	5.947
HumanNeRF [86]	30.41	0.9743	2.406	33.20	0.9752	2.899	28.18	0.9632	3.558
GNeuVox	29.77	0.9754	2.673	33.17	0.9755	3.070	28.05	0.9595	4.270
GNeuVox-ZJU	29.93	0.9763	2.408	33.45	0.9769	2.771	28.43	0.9623	3.790
GNeuVox-H36m	30.03	0.9761	2.500	33.02	0.9758	2.840	28.24	0.9613	4.028
	Subject 392			Subject 393			Subject 394		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS ( $\times 10^{-2}$ ) $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS ( $\times 10^{-2}$ ) $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS ( $\times 10^{-2}$ ) $\downarrow$
Neural Body [60]	30.10	0.9642	5.327	28.61	0.9590	5.905	29.10	0.9593	5.455
HumanNeRF [86]	31.04	0.9705	3.212	28.31	0.9603	3.672	30.31	0.9642	3.289
GNeuVox	30.81	0.9687	3.536	28.53	0.9594	4.072	29.83	0.9620	3.822
GNeuVox-ZJU	31.07	0.9705	3.376	28.53	0.9612	3.794	30.49	0.9648	3.439
GNeuVox-H36m	30.78	0.9694	3.469	28.38	0.9604	3.962	30.21	0.9635	3.593

Table 2. Comparisons about training cost and rendering quality on ZJU-MoCap dataset [60]. We list the number of training iterations on each scene and the corresponding training time. We color cells with the best values in orange and the second best values in blue.

Method	Pretrain Dataset	Perscene Iterations	Time	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS ( $\times 10^{-2}$ ) $\downarrow$
Neural Body [60]	$\times$	–	–	29.08	0.9616	5.229
HumanNeRF [86]	$\times$	400k	288 hours	30.24	0.9680	3.173
GNeuVox	$\times$	10k	50 mins	30.26	0.9678	3.450
GNeuVox	ZJU-MoCap	1k	5 mins	30.26	0.9682	3.420
GNeuVox	Human3.6M [30]	3k	15 mins	30.11	0.9677	3.399

### 3.4. Optimization

The mean square error (MSE) between the rendered color and the ground truth is used as one term of our loss function as defined in Eq. 4. Following HumanNeRF [86], we add the learned perceptual image patch similarity (LPIPS) [97] loss with the VGG [69] network to improve rendering details. Our final loss function can be formulated as

$$\mathcal{L} = \lambda_M \mathcal{L}_{MSE} + \lambda_L \mathcal{L}_{LPIPS}, \quad (10)$$

where  $\lambda_M$  and  $\lambda_L$  control the magnitudes of two loss terms respectively.

In each iteration, we sample rays of  $G$  patches with a size of  $H \times H$  on one image for computing the LPIPS loss. To speed up training, we only use the MSE loss in the first  $I_n$  iterations during fine-tuning. In the left iterations, the LPIPS loss is added to improve the details.

## 4. Experiments

In this section, we first describe the implementation details in Sec. 4.1. Then we show evaluation results and compare with other related methods in Sec. 4.2. In Sec. 4.3, ablation studies and a series of analysis are performed for better understanding the key components of GNeuVox.

### 4.1. Implementation Details

We implemented our framework using PyTorch [56]. The general voxels are constructed at a resolution of  $160^3$ . We build individual voxels separately for each scene with the same resolution as general voxels. The channel number of features in both types of voxel grids is set as 6. For voxel features  $v$ , we use positional encoding with frequencies of 2. Frequencies of positional encoding on the coordinates  $(x, y, z)$  and direction  $d$  are set as 10 and 4 respectively. We provide more details of the pose deformation network and radiation network in supplementary materials.

For training, we use an Adam [35] optimizer with  $(0.9, 0.99)$   $\beta$  values. Our training is divided into two

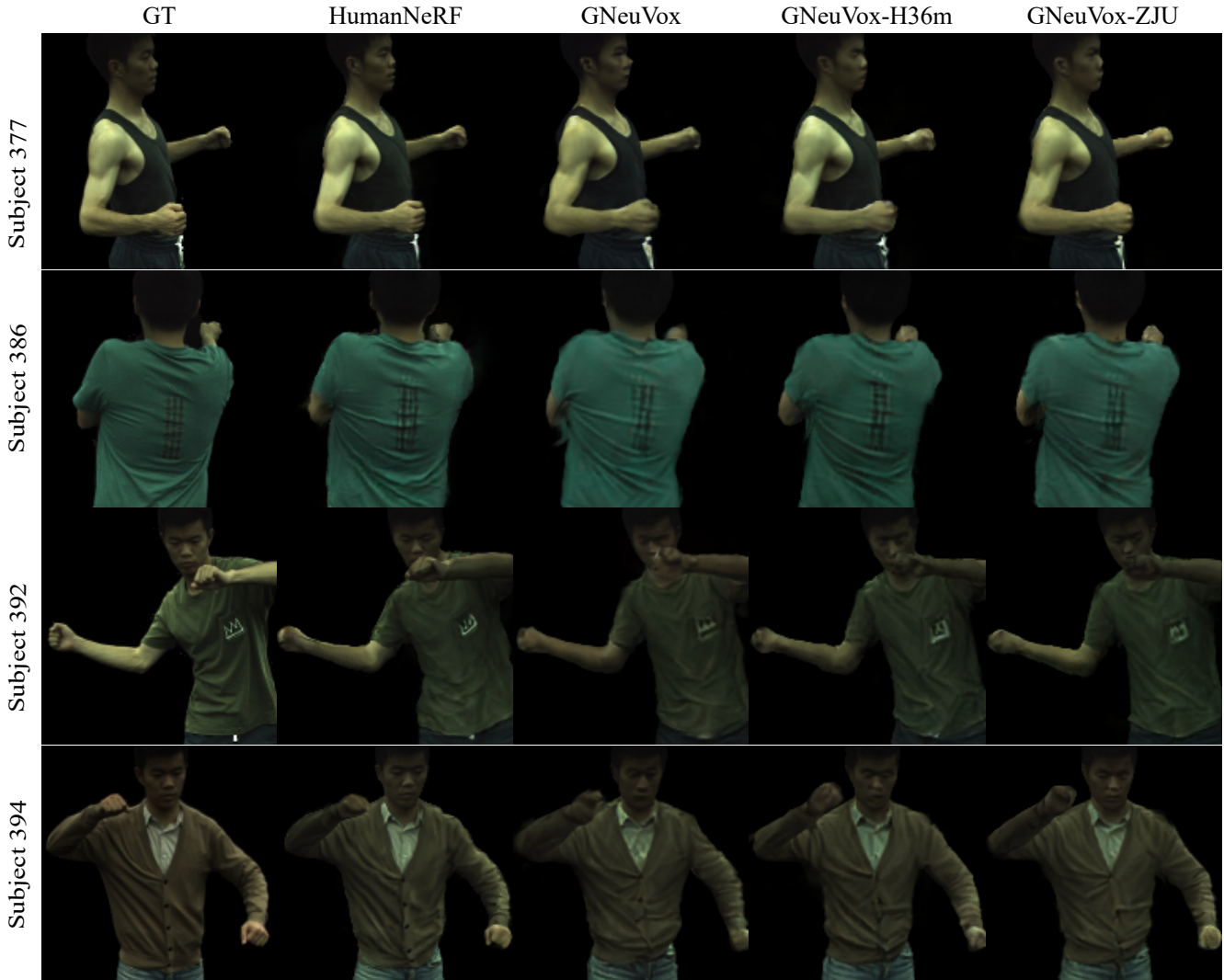


Figure 2. Qualitative comparisons between HumanNeRF [98] and our GNeuVox on the ZJU-MoCap [60] dataset. GNeuVox represents training from scratch, GNeuVox-ZJU and GNeuVox-H36m represent the fine-tuning results with pretraining on the ZJU-MoCap and Human3.6M [30] dataset respectively.

phases, pretraining and fine-tuning. During both stages, we initialize the learning rate as  $5 \times 10^{-5}$ , except voxel grids and  $\Phi_r$  are set as  $2 \times 10^{-2}$  and  $5 \times 10^{-4}$  respectively. The learning rate decays by a factor of 0.1 for every 500k iterations. The pretraining phase takes 50k iterations in total. The model is fine-tuned for 1k iterations when pretrained on the ZJU-MoCap dataset [60] and for 3k iterations when pretrained on Human3.6M [30]. We also provide a training-from-scratch version for 10k iterations. All experiments are performed on one single GeForce RTX 3090 GPU.

In practice, we get 6  $32 \times 32$  patches for sampling rays. 128 points are randomly sampled along each ray. For pretraining or training from scratch, we set  $\lambda_M$  as 0.2 and  $\lambda_L$  as 1 in Eq.10.  $\lambda_M$  is set as 10 and  $\lambda_L$  as 0 in the first 300 iterations during fine-tuning. For the rest iterations,  $\lambda_M$  is set as 0.2 and  $\lambda_L$  as 1.

## 4.2. Evaluation

**Training from Scratch** When training from scratch, We only use the images taken by “camera 1” in ZJU-MoCap [60] as training for the monocular setting. And the remaining 22 camera views, except for “camera 1” are used as evaluation. The selected images ensure that the human body rotates 360 degrees. We downsample the image to the half resolution of  $540 \times 540$ . These settings are consistent with HumanNeRF [86]. The comparison results of our method, Neural body [60], and HumanNeRF are shown in Tab. 1. We use three metrics for evaluation, namely peak signal-to-noise ratio (PSNR), structural similarity (SSIM), and LPIPS [97]. The rendering results from Neural body are worse than other methods, as it is designed for the multi-view setting. HumanNeRF and our method can achieve similar performance on the one-view setting. Neural body

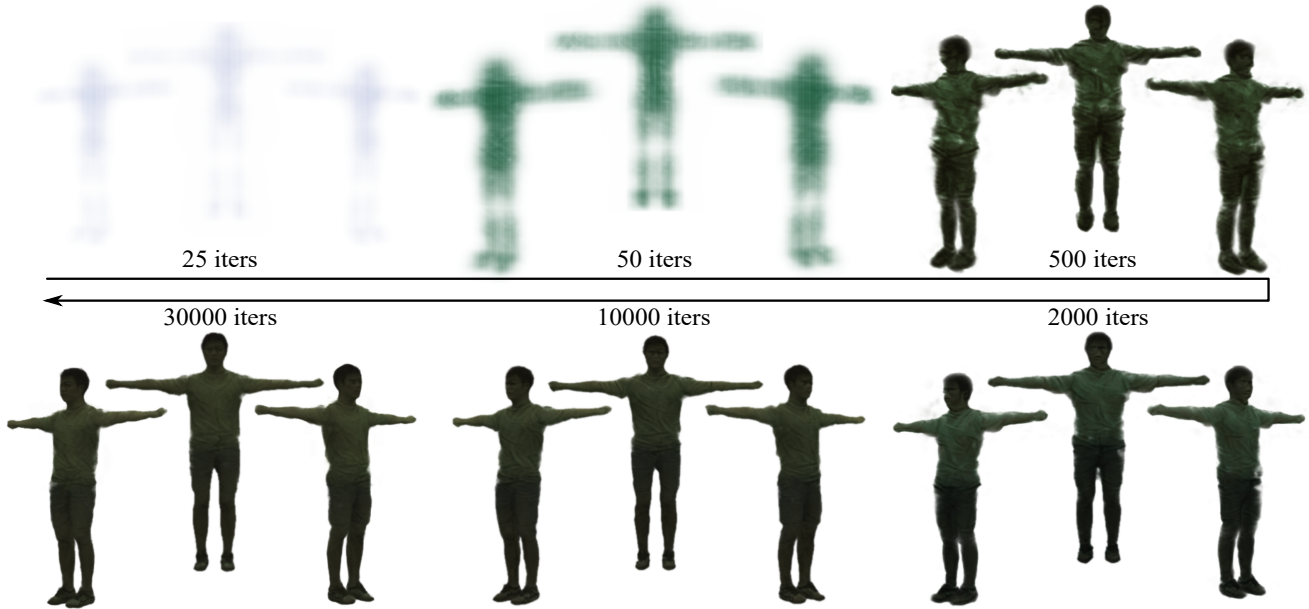


Figure 3. Change process of general voxels during pretraining on the ZJU-MoCap [60] dataset.

Table 3. Comparisons about rendering quality on PeopleSnapshot [2]. GNeuVox-ZJU and GNeuVox-H36m represent the fine-tuning results after pretraining on the ZJU-MoCap [60] and Human3.6M [30] dataset.

Method	Time	PSNR $\uparrow$
Neural Body [60]	$\sim$ 14 hours	24.47
Anim-NeRF [12]	$\sim$ 13 hours	28.89
GNeuVox-ZJU	5 mins	27.73
GNeuVox-H36m	5 mins	27.94

and HumanNeRF take days to converge. In Tab. 2, our method trained from scratch can get comparable results on ZJU-MoCap in less than 50 minutes. Here we only show the evaluation results for 6 subjects.

**Pretraining on ZJU-MoCap** We adopt the leave-one-out evaluation approach for ZJU-MoCap pertaining, which is also used in MPS-NeRF [22]. 7 out of 8 subjects (313, 377, 386, 387, 390, 392, 393, 394) are selected for pretraining while the left one is used for fine-tuning and evaluation, so we need to pretrain more than once on ZJU-MoCap to get pretrained models. This manner guarantees no overlap between datasets for pertaining and fine-tuning. When evaluation, datasets are no different from training from scratch. As shown in Tab. 2, training time can be shortened to 5 minutes if we load a pretrained model.

**Pretraining on Human3.6M [30]** For pretraining on Human3.6M, 7 subjects (S1, S5, S6, S7, S8, S9, and S11) pre-

Table 4. Ablation study about the generalization effects of several key components.

General Voxels	Radiance Fields	Deform CNN	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
$\checkmark$	$\checkmark$	$\checkmark$	30.17	0.9678	3.725
$\times$	$\checkmark$	$\checkmark$	30.04	0.9663	4.493
$\checkmark$	$\times$	$\checkmark$	29.44	0.9625	4.955
$\checkmark$	$\checkmark$	$\times$	30.14	0.9669	4.064
$\times$	$\times$	$\times$	28.89	0.9540	6.561

processed by Animatable NeRF [12] are used. Then the pre-trained model is evaluated across datasets on ZJU-MoCap. For each subject in Human3.6M, all four views are used. Images are sampled for every five frames. The image is downsampled to half the original size to  $500 \times 500$ . Note that the pretrain process only needs to be performed once on Human3.6M and applied to all subjects in ZJU-MoCap. Tab. 2 shows it only takes 15 minutes to converge with a Human3.6M-pretrained model and the PSNR is comparable to HumanNeRF. It is observed that the quantitative results have big variances when the model is trained from scratch. And sometimes the training may go collapse. Not only with our method, but the same problem also occurs with HumanNeRF. However, with pretrained general neural voxels and networks, the above situation is alleviated, and more stable training results can be achieved. The prior knowledge stored in general voxels can help our method to fit stably on new human bodies.

**Evaluation on PeopleSnapshot [2]** To explore the model’s ability to synthesize free-viewpoint of loose clothing, we conduct experiments on PeopleSnapshot using the same setting as [12]. Due to the different human body pose definitions used, we recalculate the poses using [1]. We show the average PSNR of four scenes in in Tab. 3, and the results of Neural Body [60] and Anim-NeRF [12] are from Anim-NeRF. Our method can achieve results comparable to the current state-of-the-art method (Anim-NeRF) in just 5 minutes by loading the pretrained model. Even on human bodies wearing loose clothing, our method converges quickly, demonstrating strong generalization performance. Our method has indeed learned a good human template in general voxels, leading to fast convergence.

**Results of Visualization** In Fig. 2 we show rendered images of training from scratch for 3k iterations and fine-tuning the same iterations for each scene. It can be observed that the rendering results are significantly improved with the pretrained model, which are comparable to HumanNeRF with far less training cost. Pretraining on the same dataset and across datasets can both speed up the convergence, especially in terms of better facial expressions and clothing details.

### 4.3. Ablation Study and Analysis

**Generalization Effects of Separate Components** We selectively and partially load the pretrained model to explore the generalization effectiveness of several components in the fine-tuning phase. We conduct several experiments, and at each time we choose not to load a part of the pretrained network and fine-tune the model for 500 iterations. Three evaluated parts of the network include the convolutional layers in  $F_d$ , general voxels  $V^g$  and the radiance field  $\Phi_r$ . As shown in Tab. 4, the three components all play important roles for short-time training. Promotions brought by general voxels and  $\Phi_r$  are particularly prominent, 0.13dB and 0.73dB PSNR respectively. When loading the pretrained general voxels, the template features stored in it assist individual voxels to learn the unique features of the scene, and the two together represent the information of the entire human. Loading the pretrained radiation field can well decode the features in the two types of voxels. The pretrained CNN in  $F_d$  can help quickly find out how to convert the coordinates of the observation space to the canonical space at the beginning of training. Making all these components generalizable guarantees the best acceleration performance.

**Changes of General Voxels in Pretraining** In Fig. 3, we show the change process of general voxels in the pretraining process. It can be observed that general voxels first gradually learn the outline of the human body in T-pose, but they learn some detailed information such as clothes folds in the

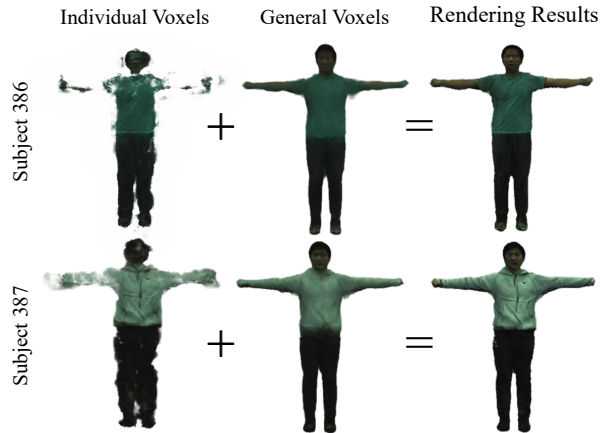


Figure 4. Rendering results by inferring general voxels and individual voxels separately.

middle pretraining process. This fold information is specific and different for each human. In the later pretraining process, these contents are degraded and the learned human becomes more general. Besides, the approximate positions of facial features are gradually learned. This information is beneficial for individual voxels, at the fine-tuning phase, to fill the details of facial features into general voxels, which results in rapid convergence.

**Roles of Two Neural Voxels** We analyze how the learned features of both voxels look like after fine-tuning. To achieve this, we set the parameters of the two voxels to zero, respectively, and obtain the result of rendering the two voxels separately after fine-tuning. The rendering result of individual voxels is very similar to the texture map, including details such as clothes folds and facial expressions. Although it consists of feature vectors that can be optimized, it exhibits obvious physical meaning. The rendered result of general voxels supplements the approximate shapes and colors of the clothes, compared to that before fine-tuning as in Fig. 4. The final rendering result is obtained by combining the differential textures contained in individual voxels and the contour features of general voxels.

## 5. Discussion

**Limitations** It is observed that human face rendering shows vagueness and distortion to some degree. We analyze the reason as follow. Geometry priors used in the proposed generalization framework mainly comes from the human pose (SMPL) and skeleton (general voxels). As the fine-tuning iterations are quite few ( $\leq 3k$ ), it is quite hard to learn the face contour accurately which may include complex non-rigid deformation. Besides, the clothes wrinkles do not match the ground truth image exactly. Improving the method towards better non-rigid deformation will be an important and interesting direction.



**Future Works & Potential** Besides the aforementioned orthogonal direction, we would like to further demonstrate the potential value our method carries. The generalization ability of the proposed general neural voxels is proved to be effective on human bodies in this paper. We plan to extend the pretraining process to some large-scale human body datasets [95, 8] and explore its effects on learning downstream human structures. More importantly, we believe the general voxels can be a promising data structure for representing scenes/objects, not limited to human bodies. General neural voxels, as **generalizable explicit representations**, enjoy high efficiency in optimization while they are more explainable than implicit representations. It is expected that general voxels can be applied to more categories of real-life objects, *e.g.* vehicles, animals and plants *etc.*

## 6. Conclusion

In this paper, we build a fast-training framework for rendering free-view moving human bodies from a monocular video. Neural voxels are introduced to represent humans and accelerate the optimization phase. Besides, we propose to construct two types of neural voxels, one pretrained on various human bodies to extract a basic skeleton while the other one targeted at a specific human. With two neural voxels collaborating, the training phase can be completed in an extremely short time. The concept of general voxels, enjoying high optimization efficiency, is expected to be extended for modeling more categories of objects.

## References

- [1] Easymocap - make human motion capture easier. Github, 2021. [8](#)
- [2] Thiemo Alldieck, Marcus Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. Detailed human avatars from monocular video. In *3DV*, 2018. [2](#), [7](#), [8](#), [13](#)
- [3] Thiemo Alldieck, Mihai Zanfir, and Cristian Sminchisescu. Photorealistic monocular 3d reconstruction of humans wearing clothing. In *CVPR*, 2022. [2](#)
- [4] Guha Balakrishnan, Amy Zhao, Adrian V Dalca, Fredo Durand, and John Guttag. Synthesizing images of humans in unseen poses. In *CVPR*, 2018. [1](#)
- [5] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, 2021. [2](#), [12](#)
- [6] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *CVPR*, pages 5470–5479, 2022. [2](#)
- [7] Alexander Bergman, Petr Kellnhofer, Wang Yifan, Eric Chan, David Lindell, and Gordon Wetzstein. Generative neural articulated radiance fields. *NeurIPS*, 2022. [2](#)
- [8] Zhongang Cai, Daxuan Ren, Ailing Zeng, Zhengyu Lin, Tao Yu, Wenjia Wang, Xiangyu Fan, Yang Gao, Yifan Yu, Liang Pan, Fangzhou Hong, Mingyuan Zhang, Chen Change Loy, Lei Yang, and Ziwei Liu. Humman: Multi-modal 4d human dataset for versatile sensing and modeling. In *ECCV*, 2022. [9](#)
- [9] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. Everybody dance now. In *ICCV*, 2019. [1](#)
- [10] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *ECCV*, 2022. [2](#)
- [11] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *ICCV*, 2021. [2](#)
- [12] Jianchuan Chen, Ying Zhang, Di Kang, Xuefei Zhe, Linchao Bao, Xu Jia, and Huchuan Lu. Animatable neural radiance fields from monocular rgb videos, 2021. [1](#), [4](#), [7](#), [8](#)
- [13] Mingfei Chen, Jianfeng Zhang, Xiangyu Xu, Lijuan Liu, Yujun Cai, Jiashi Feng, and Shuicheng Yan. Geometry-guided progressive nerf for generalizable and efficient neural human rendering. In *ECCV*, 2022. [2](#)
- [14] Xu Chen, Tianjian Jiang, Jie Song, Jinlong Yang, Michael J Black, Andreas Geiger, and Otmar Hilliges. gdna: Towards generative detailed neural avatars. In *CVPR*, 2022. [2](#)
- [15] Xu Chen, Yufeng Zheng, Michael J Black, Otmar Hilliges, and Andreas Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. In *ICCV*, 2021. [1](#), [2](#), [4](#)
- [16] Wei Cheng, Su Xu, Jingtian Piao, Chen Qian, Wayne Wu, Kwan-Yee Lin, and Hongsheng Li. Generalizable neural performer: Learning robust radiance fields for human novel view synthesis. *arXiv:2204.11798*, 2022. [2](#)
- [17] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised NeRF: Fewer views and faster training for free. In *CVPR*, June 2022. [2](#)
- [18] Jiemin Fang, Taoran Yi, Xinggang Wang, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Matthias Nießner, and Qi Tian. Fast dynamic radiance fields with time-aware neural voxels. *arxiv:2205.15285*, 2022. [2](#), [4](#)
- [19] Wanshui Gan, Hongbin Xu, Yi Huang, Shifeng Chen, and Naoto Yokoya. V4d: Voxel for 4d novel view synthesis. *arXiv:2205.14332*, 2022. [2](#)
- [20] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *ICCV*, 2021. [2](#)
- [21] Hang Gao, Ruilong Li, Shubham Tulsiani, Bryan Russell, and Angjoo Kanazawa. Monocular dynamic view synthesis: A reality check. In *NeurIPS*, 2022. [2](#)
- [22] Xiangjun Gao, Jiaolong Yang, Jongyoo Kim, Sida Peng, Zicheng Liu, and Xin Tong. Mps-nerf: Generalizable 3d human rendering from multiview images. *PAMI*, 2022. [2](#), [7](#)
- [23] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *ICCV*, 2021. [2](#)
- [24] Chen Geng, Sida Peng, Zhen Xu, Hujun Bao, and Xiaowei Zhou. Learning neural volumetric representations of dynamic humans in minutes. In *arXiv:2302.12237*, 2023. [2](#)

- [25] Marc Habermann, Lingjie Liu, Weipeng Xu, Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. Real-time deep dynamic characters. *ACM Transactions on Graphics (TOG)*, 2021. 1
- [26] Tong He, Yuanlu Xu, Shunsuke Saito, Stefano Soatto, and Tony Tung. Arch++: Animation-ready clothed human reconstruction revisited. In *ICCV*, 2021. 2
- [27] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. In *ICCV*, 2021. 2
- [28] Xin Huang, Qi Zhang, Feng Ying, Hongdong Li, Xuan Wang, and Qing Wang. Hdr-nerf: High dynamic range neural radiance fields. *arXiv:2111.14451*, 2021. 2
- [29] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. Arch: Animatable reconstruction of clothed humans. In *CVPR*, 2020. 2
- [30] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *PAMI*, 2013. 5, 6, 7, 13
- [31] Boyi Jiang, Yang Hong, Hujun Bao, and Juyong Zhang. Selfrecon: Self reconstruction your digital avatar from monocular video. In *CVPR*, 2022. 2
- [32] Tianjian Jiang, Xu Chen, Jie Song, and Otmar Hilliges. Instantavatar: Learning avatars from monocular video in 60 seconds. *arXiv:2212.10550*, 2022. 2
- [33] Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel, and Anurag Ranjan. Neuman: Neural human radiance field from a single video. In *ECCV*, 2022. 1, 2, 4
- [34] M. Johari, Y. Lepoittevin, and F. Fleuret. Geonerf: Generalizing nerf with geometry priors. In *CVPR*, 2022. 2
- [35] A KingaD. A methodforstochasticoptimization. *ICLR*, 2015. 5
- [36] Youngjoong Kwon, Dahun Kim, Duygu Ceylan, and Henry Fuchs. Neural human performer: Learning generalizable radiance fields for human performance rendering. *NeurIPS*, 2021. 2
- [37] Jiaxin Li, Zijian Feng, Qi She, Henghui Ding, Changhu Wang, and Gim Hee Lee. Mine: Towards continuous depth mpi with nerf for novel view synthesis. In *ICCV*, 2021. 2
- [38] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *CVPR*, 2021. 2
- [39] Jia-Wei Liu, Yan-Pei Cao, Weijia Mao, Wenqiao Zhang, David Junhao Zhang, Jussi Keppo, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Devrf: Fast deformable voxel radiance fields for dynamic scenes. *arXiv:2205.15723*, 2022. 2
- [40] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. *ACM Trans. Graph.(ACM SIGGRAPH Asia)*, 2021. 2
- [41] Yuan Liu, Sida Peng, Lingjie Liu, Qianqian Wang, Peng Wang, Christian Theobalt, Xiaowei Zhou, and Wenping Wang. Neural rays for occlusion-aware image-based rendering. In *CVPR*, 2022. 2
- [42] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: learning dynamic renderable volumes from images. *ACM Transactions on Graphics*, 2019. 2
- [43] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graph.(ACM SIGGRAPH Asia)*, 2015. 1, 3
- [44] Liqian Ma, Xu Jia, Qianru Sun, Bernt Schiele, Tinne Tuytelaars, and Luc Van Gool. Pose guided person image generation. *NeurIPS*, 2017. 1
- [45] Li Ma, Xiaoyu Li, Jing Liao, Qi Zhang, Xuan Wang, Jue Wang, and Pedro V Sander. Deblur-nerf: Neural radiance fields from blurry images. In *CVPR*, 2022. 2
- [46] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*, 2021. 2
- [47] Nelson Max. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics*, 1995. 2
- [48] Marko Mihajlovic, Aayush Bansal, Michael Zollhoefer, Siyu Tang, and Shunsuke Saito. KeypointNeRF: Generalizing image-based volumetric avatars using relative spatial encoding of keypoints. In *ECCV*, 2022. 2
- [49] Ben Mildenhall, Peter Hedman, Ricardo Martin-Brualla, Pratul P Srinivasan, and Jonathan T Barron. Nerf in the dark: High dynamic range view synthesis from noisy raw images. In *CVPR*, pages 16190–16199, 2022. 2
- [50] 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. 1, 2
- [51] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv:2201.05989*, 2022. 2
- [52] Natalia Neverova, Riza Alp Guler, and Iasonas Kokkinos. Dense pose transfer. In *ECCV*, 2018. 1
- [53] Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Marc Stamminger. Real-time 3d reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (ToG)*, 2013. 2
- [54] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021. 2
- [55] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *ACM Trans. Graph.*, 2021. 2
- [56] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *NeurIPS*, 32, 2019. 5

- [57] Bo Peng, Jun Hu, Jingtao Zhou, Xuan Gao, and Juyong Zhang. Intrinsicngp: Intrinsic coordinate based hash encoding for human nerf. *arXiv:2302.14683*, 2023. 2
- [58] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *ICCV*, 2021. 1, 2, 4
- [59] Sida Peng, Shangzhan Zhang, Zhen Xu, Chen Geng, Boyi Jiang, Hujun Bao, and Xiaowei Zhou. Animatable neural implicit surfaces for creating avatars from videos. *arXiv:2203.08133*, 2022. 1, 2, 4
- [60] 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. 1, 2, 4, 5, 6, 7, 8, 12, 13
- [61] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *CVPR*, 2021. 2
- [62] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *ICCV*, 2021. 2
- [63] Christian Reiser, Richard Szeliski, Dor Verbin, Pratul P Srinivasan, Ben Mildenhall, Andreas Geiger, Jonathan T Barron, and Peter Hedman. Merf: Memory-efficient radiance fields for real-time view synthesis in unbounded scenes. *arXiv:2302.12249*, 2023. 2
- [64] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *ICCV*, 2021. 2
- [65] 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 *ICCV*, 2019. 2
- [66] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *CVPR*, 2020. 2
- [67] Soubhik Sanyal, Alex Vorobiov, Timo Bolkart, Matthew Loper, Betty Mohler, Larry S Davis, Javier Romero, and Michael J Black. Learning realistic human posing using cyclic self-supervision with 3d shape, pose, and appearance consistency. In *ICCV*, 2021. 1
- [68] Kripasindhu Sarkar, Vladislav Golyanik, Lingjie Liu, and Christian Theobalt. Style and pose control for image synthesis of humans from a single monocular view. *arXiv:2102.11263*, 2021. 1
- [69] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv:1409.1556*, 2014. 5
- [70] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. Deepvoxels: Learning persistent 3d feature embeddings. In *CVPR*, 2019. 2
- [71] Liangchen Song, Anpei Chen, Zhong Li, Zhang Chen, Lele Chen, Junsong Yuan, Yi Xu, and Andreas Geiger. Nerf-player: A streamable dynamic scene representation with decomposed neural radiance fields. *arXiv:2210.15947*, 2022. 2
- [72] Shih-Yang Su, Frank Yu, Michael Zollhöfer, and Helge Rhodin. A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose. In *NeurIPS*, 2021. 1, 2, 4
- [73] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, 2022. 2, 12
- [74] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *NeurIPS*, 2020. 3
- [75] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *ICCV*, 2021. 2
- [76] Alex Trevithick and Bo Yang. Grf: Learning a general radiance field for 3d scene representation and rendering. In *arXiv:2010.04595*, 2020. 2
- [77] Liao Wang, Jiakai Zhang, Xinhang Liu, Fuqiang Zhao, Yanshun Zhang, Yingliang Zhang, Minye Wu, Lan Xu, and Jingyi Yu. Fourier plencotrees for dynamic radiance field rendering in real-time. *arXiv:2202.08614*, 2022. 2
- [78] 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. 2
- [79] Shaofei Wang, Marko Mihajlovic, Qianli Ma, Andreas Geiger, and Siyu Tang. Metaavatar: Learning animatable clothed human models from few depth images. *NeurIPS*, 2021. 2
- [80] Shaofei Wang, Katja Schwarz, Andreas Geiger, and Siyu Tang. Arah: Animatable volume rendering of articulated human sdfs. In *ECCV*, 2022. 2
- [81] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. *arXiv:1808.06601*, 2018. 1
- [82] Tuanfeng Y Wang, Duygu Ceylan, Krishna Kumar Singh, and Niloy J Mitra. Dance in the wild: Monocular human animation with neural dynamic appearance synthesis. In *3DV*, 2021. 1
- [83] Yiming Wang, Qingzhe Gao, Libin Liu, Lingjie Liu, Christian Theobalt, and Baoquan Chen. Neural novel actor: Learning a generalized animatable neural representation for human actors. *arXiv:2208.11905*, 2022. 2
- [84] Ziyang Wang, Timur Bagautdinov, Stephen Lombardi, Tomas Simon, Jason Saragih, Jessica Hodgins, and Michael Zollhofer. Learning compositional radiance fields of dynamic human heads. In *CVPR*, 2021 MonocularReal. 1, 2
- [85] Chung-Yi Weng, Brian Curless, and Ira Kemelmacher-Shlizerman. Vid2actor: Free-viewpoint animatable person synthesis from video in the wild. *arXiv:2012.12884*, 2020. 1

Table 5. Ablation study about the number of pretraining iterations on fine-tuning results.

Pretrain Iterations	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS ( $\times 10^{-2}$ ) $\downarrow$
50k	30.11	0.9677	3.399
200k	30.26	0.9681	3.489

- [86] Chung-Yi Weng, Brian Curless, Pratul P. Srinivasan, Jonathan T. Barron, and Ira Kemelmacher-Shlizerman. HumanNeRF: Free-viewpoint rendering of moving people from monocular video. In *CVPR*, 2022. 1, 2, 3, 4, 5, 6
- [87] Chung-Yi Weng, Pratul P. Srinivasan, Brian Curless, and Ira Kemelmacher-Shlizerman. PersonNeRF: Personalized reconstruction from photo collections. *CVPR*, 2023. 2
- [88] Yuliang Xiu, Jinlong Yang, Dimitrios Tzionas, and Michael J. Black. ICON: Implicit Clothed humans Obtained from Normals. In *CVPR*, 2022. 2
- [89] Hongyi Xu, Thiemo Alldieck, and Cristian Sminchisescu. H-nerf: Neural radiance fields for rendering and temporal reconstruction of humans in motion. In *NeurIPS*, 2021. 1, 2, 4
- [90] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Point-based neural radiance fields. In *CVPR*, 2022. 2
- [91] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. *ECCV*, 2018. 2
- [92] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. 2
- [93] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 2
- [94] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 2
- [95] Tao Yu, Zerong Zheng, Kaiwen Guo, Pengpeng Liu, Qionghai Dai, and Yebin Liu. Function4d: Real-time human volumetric capture from very sparse consumer rgbd sensors. In *CVPR*, 2021. 9
- [96] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv:2010.07492*, 2020. 2
- [97] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 5, 6
- [98] 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 *CVPR*, 2022. 2, 6, 13

Table 6. Ablation study about fine-tuning iterations with the pre-trained model on the ZJU-MoCap dataset [60].

Iterations	500	1000	3000
PSNR $\uparrow$	30.17	30.26	30.32

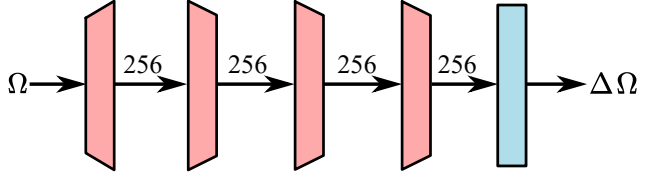


Figure 5. Architecture of the pose refinement network  $\Phi_p$ . The red trapezoidal squares denote that ReLU activation function is applied, and the blue square denotes no activation function.

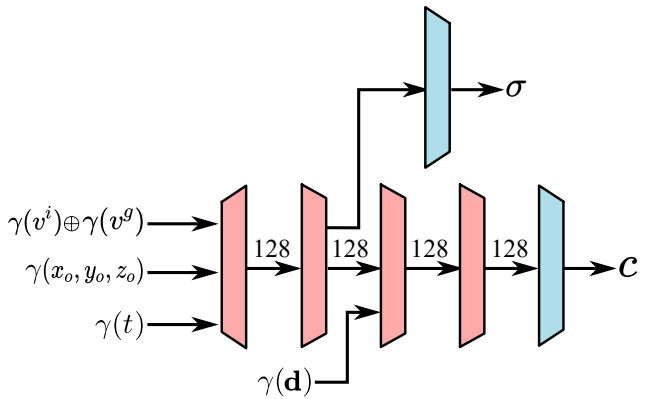


Figure 6. Architecture of the radiance network  $\Phi_r$ . The blue trapezoidal squares denote that other activation functions are applied. The activation functions applied to the density  $\sigma$  and color  $c$  are softplus [5, 73] and sigmoid respectively.  $\gamma$  denotes the positional encoding.

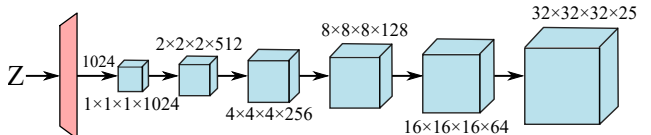


Figure 7. The convolutional layers in the pose deformation network. First, we use one layer MLP to reshape the embedding  $Z$  to  $1 \times 1 \times 1 \times 1024$ , and then generate a volume of  $32 \times 32 \times 32 \times 25$  by transposed convolutions.

## A. Appendix

### A.1. Neural Architectures

In this section, we provide detailed structures of the pose refinement network  $\Phi_p$ , the convolutional layers in the pose deformation network  $\mathbf{F}_d$ , and the radiance network  $\Phi_r$ . As shown in Fig. 5 and Fig. 6,  $\Phi_p$  and  $\Phi_r$  are both instantiated as MLPs. Fig. 7 shows that through one layer MLP and transposed convolutions, the embedding  $Z$  is transformed into a volume, which is trilinearly interpolated to obtain the

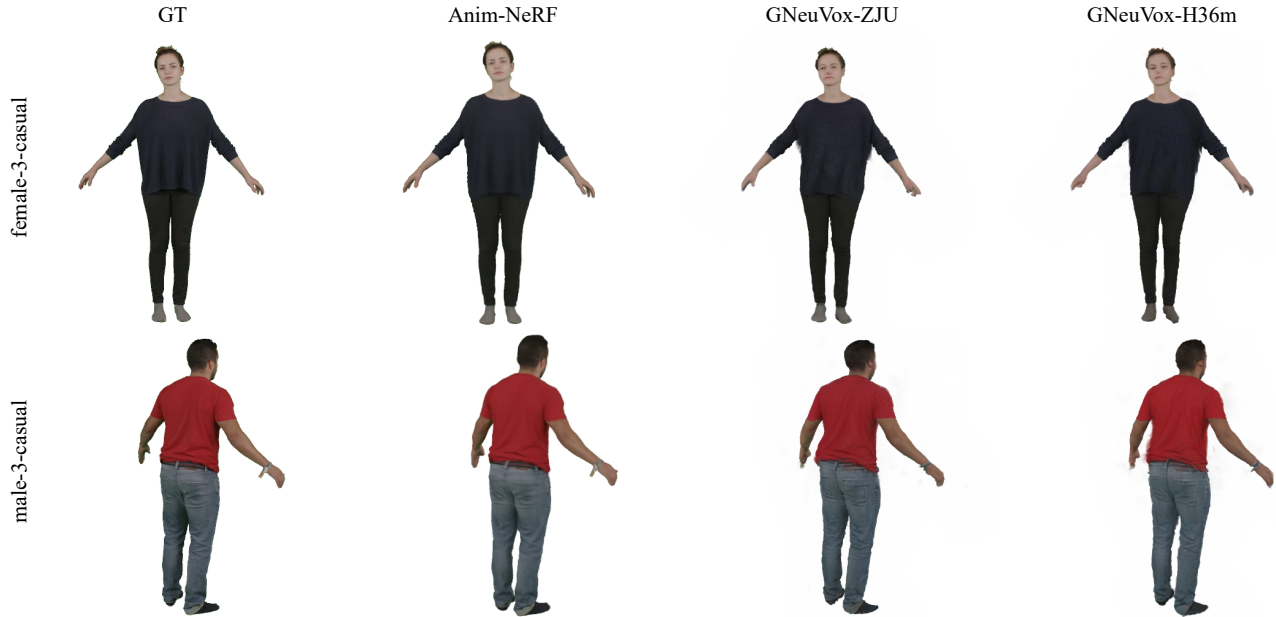


Figure 8. Qualitative comparisons between Anim-NeRF [98] and our GNeuVox on the PeopleSnapshot [2] dataset. GNeuVox-ZJU and GNeuVox-H36m denote the fine-tuning results with pretraining on the ZJU-MoCap [60] and Human3.6M [30] datasets respectively.

blend weigh  $w^i$  in Eq.6 for the coordinates deformation.

## A.2. More Ablation Studies

We pretrain the model on Human3.6M [30] for different iterations and then fine-tune it on the ZJU-MoCap [60] dataset for  $3k$  iterations to explore the impact of pretraining iterations on fine-tuning results. As shown in Tab. 5, more iterations on the pretraining stage show little effect on fine-tuning results. The Human3.6M dataset has only 7 human bodies, and  $50k$  iterations of pretraining are enough. We might get better fine-tuning results if we pretrain on more human bodies.

Moreover, we fine-tune the model with different iterations on ZJU-MoCap. Based on the experiment results as shown in Tab. 6, more fine-tuning iterations do not lead to a significant improvement in rendering quality.

## A.3. More Qualitative Comparisons

We conduct experiments on the PeopleSnapshot [2] dataset, where human bodies are wearing looser clothing. We show the visualization results in Fig. 8. Despite our method not constructing specialized modules for modeling non-rigid clothing or facial details, GNeuVox still achieves decent results. We leave building stronger modules for non-rigid deformations as an important future direction.