# Animatable Neural Radiance Fields from Monocular RGB Video

Jianchuan Chen[1]    Ying Zhang[2]    Di Kang[2]    Xuefei Zhe[2]

Linchao Bao[2]    Huchuan Lu[1]

[1]Dalian University of Technology    [2]Tencent AI Lab

## Abstract

*We present animatable neural radiance fields for detailed human avatar creation from monocular videos. Our approach extends neural radiance fields (NeRF) to the dynamic scenes with human movements via introducing explicit pose-guided deformation while learning the scene representation network. In particular, we estimate the human pose for each frame and learn a constant canonical space for the detailed human template, which enables natural shape deformation from the observation space to the canonical space under the explicit control of the pose parameters. To compensate for inaccurate pose estimation, we introduce the pose refinement strategy that updates the initial pose during the learning process, which not only helps to learn more accurate human reconstruction but also accelerates the convergence. In experiments we show that the proposed approach achieves 1) implicit human geometry and appearance reconstruction with high-quality details, 2) photo-realistic rendering of the human from arbitrary views, and 3) animation of the human with arbitrary poses.*

## 1. Introduction

3D human digitization has a wide range of applications in industries such as film, animation, games, and virtual try-on. Existing approaches to obtain high-quality 3D human reconstruction often require expensive equipment such as multiple synchronized cameras [31], and RGB-D sensor [17], limiting their applications in practical scenarios. On the other hand, for various 3D human reconstruction approaches [29, 30, 2, 3, 38], modeling complex geometric details such as hair, glasses, and cloth wrinkles of real humans remains a challenging problem.

In this paper, we target to obtain photo-realistic 3D human avatars from monocular RGB videos. Different from existing approaches based on pre-scanned human models [36] or parametric body models [2, 38, 3], our approach implicitly reconstructs the human geometry and appear-

ance via generalizing Neural Radiance Fields (NeRF) [24], which uses a neural network to encode color and density as a function of location and viewing angle, and generates photo-realistic images by volume rendering. NeRF [24] has shown impressive ability in reconstructing a static scene from multi-view images, and inspires many researchers in extending NeRF to scenes with severe lighting changes [23] and non-rigid deformations [25, 35]. However, these approaches are uncontrollable and limited to nearly static scenes with small movements, failing to deal with human subjects with large movements.

To handle the dynamic human from monocular videos, we combine neural radiance fields with a parametric body model of SMPL [21], which enables more precise human geometry and appearance modeling, and further makes the neural radiance fields controllable. In particular, our approach extends NeRF via introducing the pose-guided deformation, which unwarps the observation space near the human body to a constant canonical space through the deformation of the SMPL. We observe that even the state-of-the-art SMPL estimator from monocular videos cannot obtain accurate parameters, which inevitably leads to blurry results. To address this problem, we propose to jointly optimize the NeRF and SMPL parameters via analysis-by-synthesis, which not only obtains better results but also accelerates the convergence of training. We demonstrate the superiority of the proposed method on multiple datasets, with both quantitative and qualitative results on novel view synthesis, 3D human reconstruction, and novel pose synthesis.

In summary, our work has the following contributions:

- We propose a method explicitly deforming the points according to SMPL pose to reconstruct a canonical NeRF model, relaxing the requirement of the static object and preserving details such as clothing and hair.

- We incorporate pose refinement into our analysis-by-synthesis approach to account for the inaccurate SMPL estimates, resulting in refined SMPL pose and greatly improved reconstruction quality.

- We achieve high-quality 3D human reconstruction from monocular RGB video, and can render photo-realistic images from arbitrary views.

- Due to our controllable SMPL-based geometry deformation, we can synthesize novel pose images, showing that our learned canonical space NeRF model is animatable.

## 2. Related work

**3D Human Reconstruction** Reconstructing 3D human has been more and more popular in recent years. Various approaches attempt to digitize a human from a single-view image [29, 30, 16, 4, 12], multi-view images [29, 30, 11, 41], RGB videos [3, 2, 1, 38, 40], or RGB-D videos [17, 42]. One stream of these approaches [2, 3, 38] utilize a parametric body model such as SMPL [21] to represent a human body with cloth deformations, which produces an animatable 3D model with high-quality textures but struggles with limited expressive ability in complex geometries such as hair and dresses. On the other hand, PIFu [29] and PIFuHD [30] based methods use an implicit representation to reconstruct a 3D surface and achieves impressive results in handling people with complex poses, hairstyles and clothing, which however, suffers from a blurry appearance and requires further registration for animation.

**Neural Representations** Representing a scene with neural networks has achieved stunning success in recent years. SRN [32] proposes an implicit neural representation that assigns feature vectors to 3D positions, and uses a differentiable ray marching algorithm for image generation. NeRF [24] establishes a static scene that maps 3D coordinates and view direction to density and color using a neural network. These methods [32, 24, 18] can render very realistic images, but they are all limited to static scenes. Dynamic NeRFs [25, 28, 18, 35] extend NeRF to dynamic scenes by introducing the latent deformation field. However, it is difficult to implicitly control the complex non-rigid deformation of human body motion. These works [10, 27, 33] combine scene representation network with parametric models [7, 21] to reconstruct dynamic humans. Instead of using latent codes or expression parameters as input, we use the human body model SMPL to explicitly deform over different poses and shapes. At the same time, this explicit method allows us to fine-tune the parameters of the SMPL.

**Human Motion Transfer** Human motion transfer aims to synthesize an image of a person with the appearance from a source human and the motion from a reference image. Recent advances using Generative Adversarial Networks (GAN) have shown convincing performance without recovering detailed 3D geometry. These works [5, 8, 19, 20] use image-to-image translation [13, 37] to map 2D skeleton images to rendering output. Due to the lack of 3D reasoning, the geometry of the generated humans is usually not consistent across multiple views and motions.

## 3. Method

In this section, we will describe the method to create a human avatar from a single portrait video of a person as shown in Figure 1. Given a video sequence with $n$ frames $\{I_t\}_{t=1}^n$ of a single human subject, we estimate the SMPL [21, 26] parameters $M(\theta_t, \beta_t)$ and camera intrinsics $K_t$ of each frame using existing human body shape and pose estimation models [15]. In order to avoid the influence of background changes caused by camera movement, we first use segmentation network [39] to obtain the foreground human mask, and set the background color to white uniformly. Our animatable neural radiance fields (Section 3.1) can be decomposed into pose-guided deformation (Section 3.2) and a neural radiance field (NeRF) defined in the canonical space. We can use the volumetric rendering (Section 3.3) to render our neural radiance field. In order to avoid the negative effects of inaccurate SMPL parameters, we propose to jointly optimize the neural radiance field and SMPL parameters (Section 3.4). We also introduce background regularization and pose regularization to improve the robustness of optimization (Section 3.5).

### 3.1. Animatable Neural Radiance Fields

For modeling the human appearance and geometry with complex non-rigid deformation, we introduce the parameterized human model SMPL [21] into the neural radiance field and present the animatable neural radiance fields $F$ which maps the 3D position $\mathbf{x} = (x, y, z)$, shape $\beta_t$ and pose $\theta_t$ into color $\mathbf{c} = (r, g, b)$ and density $\sigma$:

$$F\left(D\left(\mathbf{x}, \theta_t, \beta_t\right)\right) = (\mathbf{c}, \sigma) \tag{1}$$

where $D\left(\mathbf{x}, \theta_t, \beta_t\right)$ transforms the 3D position $\mathbf{x} = (x, y, z)$ in the observation space to $\mathbf{x}^0 = (x^0, y^0, z^0)$ in canonical space, aiming to handle human movements between different frames. The view dependence in NeRF is mainly for dealing with specular reflections of materials such as metal and glass. But skin and clothes of human are mainly diffuse reflective materials, so we remove the view direction input.

### 3.2. Pose-guided Deformation

In contrast to [35, 25] that implicitly control of the deformation of space, we use the parametric body model - SMPL, to guide the deformation of space explicitly. Here we define the observed image as from the observation space and attempt to learn a template human in the canonical space, and the explicit SMPL model enables the appearance transformation from one space to another, which facilitates the learning of a meaningful canonical space, and reduces
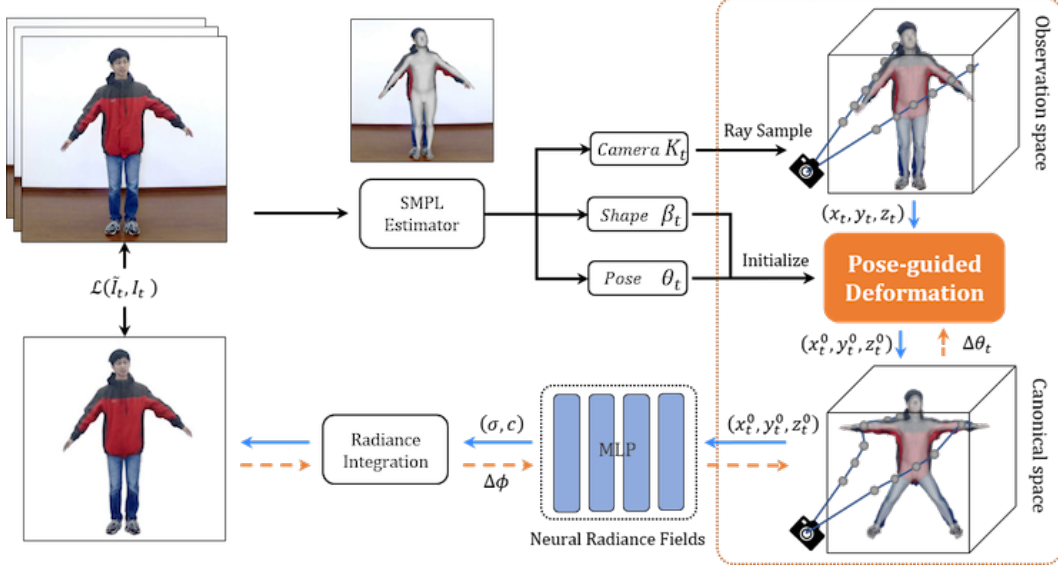
Figure 1. Overview of the proposed Animatable Neural Radiance Fields. Given a video sequence, we estimate the camera $K_t$ and SMPL parameters $M(\theta_t, \beta_t)$ of the human subject for initialization. We use volume rendering to sample points $(x_t, y_t, z_t)$ along the camera ray in observation space, and transform these points to canonical space according to pose-guided deformation. Then we input these points $(x_t^0, y_t^0, z_t^0)$ into the neural radiance field to get densities $\sigma$ and colors $\mathbf{c}$. Then we use the integral equation Eq. (5) to render the image, and jointly optimize the neural radiance field parameters $\phi$ and SMPL parameters $\theta_t$ by minimizing the error $\mathcal{L}\left(\tilde{I}_t, I_t\right)$ between the rendered image $\tilde{I}_t$ and the ground truth image $I_t$ with the mask.

the reliance of diverse input poses to generalize to unseen poses. The template pose in the canonical space is defined as X-pose $\theta_0$ (as shown in Figure 1), due to its good visibility and separability of each body part. By using the inverse transformation of the linear skinning of SMPL, the pose $\theta_t$ in observation space can be transformed into the X-pose $\theta_0$ in canonical space. Consider that the transformation functions are only defined on the surface vertices of the body mesh, we extend them to the space near the mesh surface based on the intuition that points in space near the mesh should move along with neighboring vertices. Following PaMIR [41] we define the transformation of $\mathbf{x}$ from observation space to canonical space as

$$
\begin{bmatrix} \mathbf{x}^0 \\ 1 \end{bmatrix} = \mathbf{M}(\mathbf{x}, \beta, \theta_t, \theta_0) \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}
$$
$$
\mathbf{M}(\mathbf{x}, \beta, \theta_t, \theta_0) = \sum_{v_i \in \mathcal{N}(\mathbf{x})} \frac{\omega_i}{\omega} \mathbf{M}_i (\beta, \theta_0) \left(\mathbf{M}_i(\beta, \theta_t)\right)^{-1}
$$

(2)

and the transformation matrix $\mathbf{M}_i (\beta, \theta)$ of vertex $v_i$ from T-pose to $\theta$-pose is computed by

$$
\mathbf{M}_i(\beta, \theta) = \left( \sum_{j=1}^{K} b_{i,j} \mathbf{G}_j \right) \begin{bmatrix} \mathbf{I} & \mathbf{B}_{S,i}(\beta) + \mathbf{B}_{P,i}(\theta) \\ \mathbf{0}^T & 1 \end{bmatrix}
$$

(3)

where $\mathbf{G}_j \in \mathbb{R}^{4 \times 4}$ is the world transformation of joint $j$, $b_{i,j}$ is the blend weight representing how much the rotation of part $j$ effects the vertex $v_i$, $\mathbf{B}_{P,i}(\theta) \in \mathbb{R}^3$ and $\mathbf{B}_{S,i}(\beta) \in \mathbb{R}^3$ are the pose blendshape and shape blendshape of vertex $v_i$ respectively. $\mathcal{N}(\mathbf{x})$ denotes the SMPL vertex set near $\mathbf{x}$ in the observation space, and Eq. (2) indicates that the movement of $\mathbf{x}$ relies on the movement of neighboring vertices.

The transformation weight $\omega_i$ is defined as

$$
\omega_i = \exp \left( -\frac{\|\mathbf{x} - v_i\| \left\| \hat{\mathbf{b}} - \mathbf{b}_i \right\|}{2\sigma^2} \right)
$$
$$
\omega = \sum_{v_i \in \mathcal{N}(\mathbf{x})} \omega_i
$$

(4)

where $\mathbf{b}_i$ is the blending skinning weight of $v_i$ and $\hat{\mathbf{b}}$ is the blending skinning weight of the nearest vertex, and $\|\mathbf{x} - v_i\|$ computes the L2 distance between $\mathbf{x}$ and $v_i$. Consider the fact that a point might be affected by different body parts, leading to ambiguous or even non-meaningful transformation, we adopt the blending skinning weight which characterizes the movement patterns of a vertex along with the SMPL joints [21], to strengthen the movement impact of the nearest neighbor.

### 3.3. Volumetric Rendering

We use the volume rendering techniques introduced in NeRF [24] to render the neural radiance field into a 2D image. For a given video frame $I_t$, we first convert the camera coordinate system to the SMPL coordinate system, that is, transform the SMPL global rotation and translation to the camera. Then the pixel colors are obtained by accumulating the colors and densities along the corresponding camera ray $\mathbf{r}$. In practice, the continuous integration is approximated by sampling $N$ points $\{\mathbf{x}_k\}_{k=1}^{N}$ between the near plane and the far plane along the camera ray $\mathbf{r}$ as

$$\tilde{C}_t(\mathbf{r}) = \sum_{k=1}^{N} T_k \left(1 - \exp\left(\eta_t(\mathbf{x}_k)\sigma_t(\mathbf{x}_k)\delta_k\right)\right)\mathbf{c}_t(\mathbf{x}_k)$$

$$\tilde{D}_t(\mathbf{r}) = \sum_{k=1}^{N} T_k \left(1 - \exp\left(\eta_t(\mathbf{x}_k)\sigma_t(\mathbf{x}_k)\delta_k\right)\right)$$

$$T_k = \exp\left(-\sum_{j=1}^{k-1}\eta_t(\mathbf{x}_k)\sigma_t(\mathbf{x}_j)\delta_j\right)$$

(5)

where $\delta_k = \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2$ is the distance between adjacent sampling points, and $\eta_t(\mathbf{x}_k)$ is a prior mask used to provide geometric prior guidance and deal with ambiguity during pose-guided deformation. In experiments we follow NeRF [24] to perform hierarchical volume sampling to obtain $\tilde{C}_t^c(\mathbf{r})$ and $\tilde{C}_t^f(\mathbf{r})$ with the coarse and fine network, respectively.

Since we only focus on modeling a single human subject, here we introduce an assumption for learning a more accurate neural radiance field: The densities should be zeros for points far from the surface of human mesh;

$$\eta_t(\mathbf{x}_k) = d(\mathbf{x}_k) \leq \delta$$

$$d(\mathbf{x}_k) = \sum_{v_i \in \mathcal{N}(\mathbf{x}_k)} \frac{\omega_j}{\omega} \|\mathbf{x}_k - v_i\|$$

(6)

where $d(\mathbf{x}_k)$ is the weighted distance from point $\mathbf{x_k}$ to the nearest neighbor vertices $\mathcal{N}(\mathbf{x}_k)$ in the observation space. $\delta$ is the distance threshold limiting distance between the sample point to the SMPL surface in the observation space.

### 3.4. Pose Refinement via Analysis-by-Synthesis

Our proposed method learns an animatable neural radiance field for human subjects via explicitly deforming the observation space from different frames to a constant canonical space, under the guidance of SMPL transformations. Although the current state-of-the-art pose and shape estimation methods [15] is adopted to obtain more stable SMPL parameters, in experiments we observe that results estimated by these methods do not align well with the ground truth, especially in depth. The inaccurate human body estimation could easily lead to fuzzy results and to address this problem we propose to fine-tune the

SMPL parameters during training. More specifically, we use VIBE [15] to estimate SMPL parameters $M(\theta_t, \beta_t)$ for each frame $I_t$ as initialization of a pose variable, which will be optimized during training. We use the mean shape parameters $\beta = \frac{1}{n}\sum_{t=1}^{n}\beta_t$ for different frames. It turns out that the refined pose can better fit the input image, and helps to obtain clearer and sharper results.

### 3.5. Objective Functions

Given a monocular video sequence, we learn the animatable neural radiance field with the following objective function

$$\mathcal{L} = \mathcal{L}_c + \mathcal{L}_p + \lambda_d * \mathcal{L}_d \tag{7}$$

where $\mathcal{L}_c$, $\mathcal{L}_p$, and $\mathcal{L}_d$ are reconstruction loss, pose regularization, and background regularization, respectively. $\lambda_d$ aims to balance the importance of background regularization.

**Reconstruction Loss** The reconstruction loss aims to minimize the error between the rendered images and the corresponding observed frames, which is defined as

$$\mathcal{L}_c = \sum_t \sum_{\mathbf{r}} \left\|\tilde{C}_t^c(\mathbf{r}) - C_t(\mathbf{r})\right\|_2^2 + \left\|\tilde{C}_t^f(\mathbf{r}) - C_t(\mathbf{r})\right\|_2^2 \tag{8}$$

where $\mathbf{r}$ is a camera ray passing through the image $I_t$, and $\tilde{C}_t^c(\mathbf{r})$, $\tilde{C}_t^f(\mathbf{r})$ are the rendered colors using volume rendering introduced in Sec. 3.3. $C_t(\mathbf{r})$ is the ground truth color from observed image $I_t$.

**Pose Regularization** To obtain stable and smooth pose parameters, we add the pose regularization to constrain that the optimized pose parameter should not be far away from the initial pose, and the pose parameters between adjacent frames should be similar.

$$\mathcal{L}_p = \lambda_1 \left\|\tilde{\theta}_t - \theta_t\right\| + \lambda_2 \left\|\tilde{\theta}_t - \tilde{\theta}_{t+1}\right\| \tag{9}$$

where $\tilde{\theta}_t$ and $\tilde{\theta}_{t+1}$ are the pose parameters of frame $t$ and $t + 1$. $\lambda_1$ and $\lambda_2$ are the corresponding penalty weights.

**Background Regularization** In this paper we focus on the reconstruction of the subject, so we uniformly set the background color to white. We want the density to be concentrated in the foreground, and there should be no density in the background. In order to avoid the ambiguity between the color of clothes and the background, we minimize the error between the rendered density and the mask obtained by segmentation.

$$\mathcal{L}_d = \sum_t \sum_{\mathbf{r}} \left\|\tilde{D}_t^c(\mathbf{r}) - D_t(\mathbf{r})\right\| + \left\|\tilde{D}_t^f(\mathbf{r}) - D_t(\mathbf{r})\right\| \tag{10}$$

where $\tilde{D}_t^c$ and $\tilde{D}_t^f$ is the rendered density of the coarse and fine network for the camera ray $\mathbf{r}$ from the image $I_t$, and $D_t(\mathbf{r})$ is the corresponding segmentation mask.
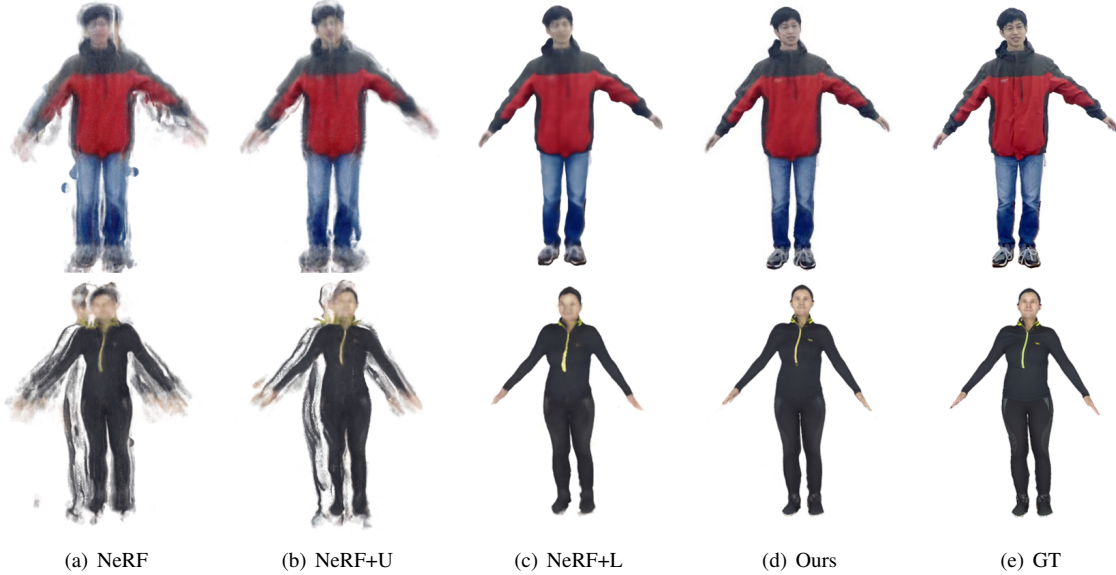
(a) NeRF  (b) NeRF+U  (c) NeRF+L  (d) Ours  (e) GT

Figure 2. Comparison of different methods on novel view synthesis on iPER(top) and Mutil-Garment(bottom).

| Subject ID | PSNR↑ | | | | | SSIM↑ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | NeRF | NeRF+L | NeRF+U | OURS | NeRF+U(GT) | NeRF | NeRF+L | NeRF+U | OURS | NeRF+U(GT) |
| people1 | 15.82 | 17.06 | 15.93 | **17.99** | 32.78 | .8332 | .9062 | .8536 | **.9100** | .9815 |
| people2 | 18.74 | 19.58 | 18.91 | **20.98** | 34.40 | .8510 | .9007 | .8632 | **.9028** | .9734 |
| people3 | 16.58 | 17.58 | 16.29 | **19.05** | 33.05 | .8340 | .9023 | .8464 | **.9094** | .9772 |
| people4 | 15.97 | 16.54 | 16.14 | **18.70** | 32.79 | .8926 | .8261 | .8406 | **.9037** | .9799 |
| Average | 16.78 | 17.69 | 16.82 | **19.18** | 33.26 | .8361 | .9005 | .8510 | **.9065** | .9780 |

Table 1. Ablation study on Multi-Garment dataset. where L represents latent code, U represents pose-guided deformation, and U(GT) represents pose-guided deformation using the ground truth SMPL.

## 3.6. Applications

The proposed approach learns an animatable neural radiance field, allowing us to reconstruct the implicit neural representation of the geometry and appearance of the human body, from a monocular video of a person turning around before a camera while holding the A-pose. After the model is trained, we can use it for novel view synthesis through volume rendering, and we can use the Marching Cubes algorithm [22] to extract the surface to get a 3D human mesh with high-quality details. In addition, we can also deform the neural radiance field to our desired pose for rendering by the pose-guided deformation.

## 4. Experiments

### 4.1. Implementation Details

Following NeRF [24] we implement the neural radiance field, and we use coarse and fine networks to represent the human body same as NeRF [24], and 64 coarse and $64+16$ fine rays samples for all experiments. Focusing on the foreground subject, we set $80\%$ of the rays to be sampled from the foreground, and the remaining $20\%$ to be sampled from the background. We set the hyper-parameters as $|\mathcal{N}(i)| = 4$, $\delta = 0.2$, $\lambda_1 = 0.001$, $\lambda_2 = 0.01$ and $\lambda_d = 0.1$. We use the image resolution of $512 \times 512$ in all experiments. For training the model we adopt the Adam optimizer [14], and it spends about 26 hours on 2 Nvidia Tesla V100 32GB GPUs.

### 4.2. Datasets and Evaluation

**Datasets** To evaluate the effectiveness of the proposed method, we conducted experiments on 3 different datasets. **Multi-Garment** [6] The dataset contains 3D scanned human body models, and registered SMPLD models and textures can be used for animation. We selected 4 human body models to synthesize the data. According to motion sequences which the subjects rotate while holding an A-pose in People-Snapshot dataset[3]. **People-Snapshot** [3] and **iPER** [19] datasets contain monocular RGB vides captured in different real scenes. The subjects are turning around
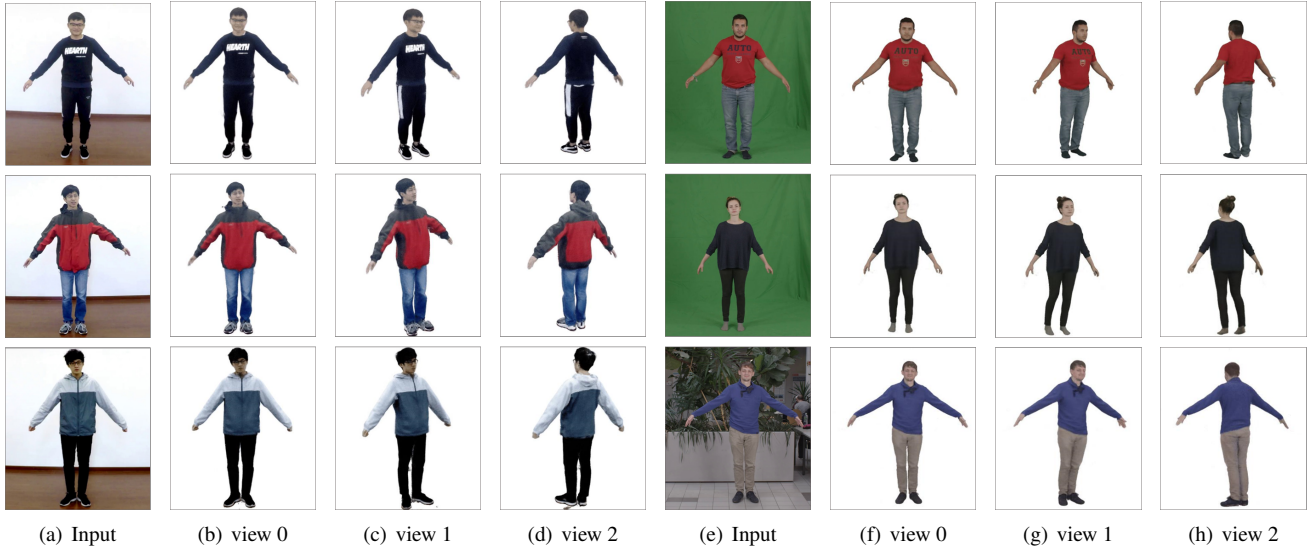
(a) Input    (b) view 0    (c) view 1    (d) view 2    (e) Input    (f) view 0    (g) view 1    (h) view 2

Figure 3. Results of Novel View Synthesis on iPER (a-d) and People-Snapshot (e-h).



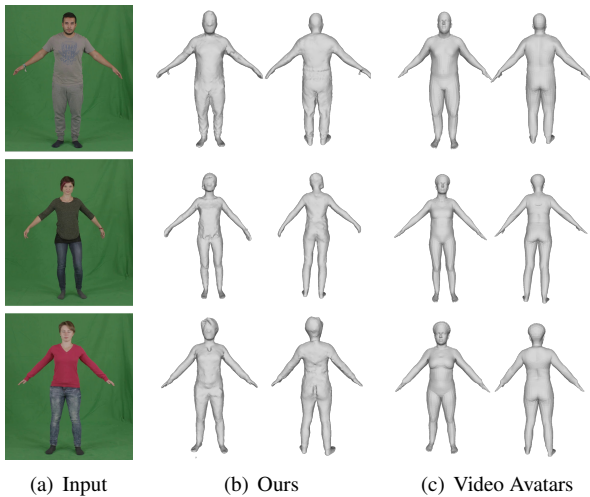(a) Input    (b) Ours    (c) Video Avatars

Figure 4. Comparison with video avatars [3] of 3D reconstruction results on People-Snapshot

from a fixed camera while holding an A-pose.

**Evaluation** In our experiment, we uniformly select about 300 frames from each video, where the subject rotates about 3 circles. The first 200 frames are used for training, and the remaining frames are used for testing. For quantitative evaluation, we evaluate our method for novel view synthesis using two metrics: peak signal-to-noise ratio (PSNR) and structural similarity index (SSIM[34]). For 3D reconstruction, we use point-to-surface Euclidean distance (P2S) and Chamfer distance (Chamfer[9]) in cm between the reconstructed surface and the ground truth surfaces. Since the datasets from the real scenes don't have the corresponding ground truth geometry, we will provide qualitative results.

**Comparison** We compare against the original NeRF [24] and NeRF+L baseline, where NeRF is conditional on per-frame learnable latent deformation code to modulate dynamic scenes. We transform the global rotation and translation of SMPL to the camera so that it is almost equivalent to a multi-view scene because the subject always keeps in A-pose, In order to verify the effectiveness of our pose fine-tuning strategy, we first use VIBE [15] to estimate the parameters of SMPL, and conduct comparative experiments between pose refinement (NeRF+U+P) and without pose refinement (NeRF+U). For synthetic data, we use the ground truth SMPL parameters (NeRF+U(GT)) for pose-guided deformation as the upper bound. For testing NeRF+L, we select the average of latent codes of $k = 5$ frames that are most similar to poses in training frames. Since there are scale and depth ambiguities in 3D reconstruction from monocular videos, we register our meshes to ground truth geometry for comparison.

### 4.3. Novel view synthesis

Our animatable neural radiance field can be rendered from arbitrary views of the same pose. Table 1 compares the results of different approaches on novel view synthesis on the Multi-Garment dataset, and we can see that the introduction of a latent code (NeRF+L) can produce better results than NeRF, and the proposed approach achieves more reliable results than implicit control of deformation. Experiments also show that inaccurate SMPL parameters will have a very negative impact (NeRF+U), and with pose refinement, the quality of novel view synthesis can be significantly improved. It is interesting to see that with ground truth SMPL parameters, the PSNR and SSIM is much higher than all the approaches, demonstrating the necessity
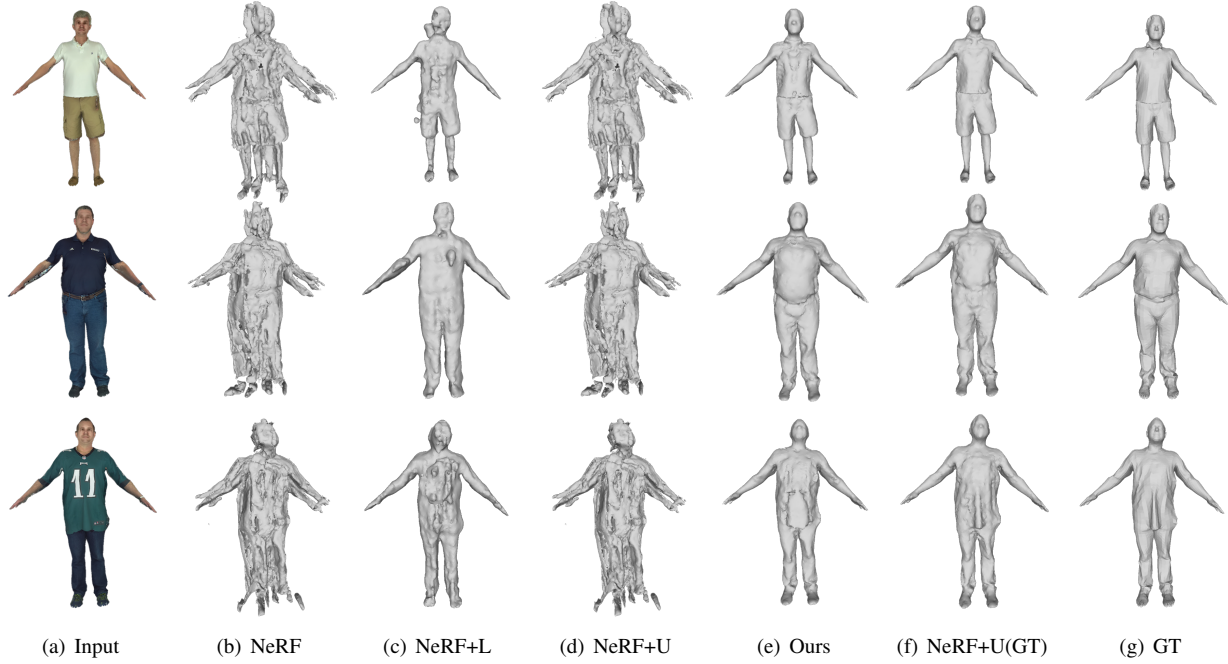
| (a) Input | (b) NeRF | (c) NeRF+L | (d) NeRF+U | (e) Ours | (f) NeRF+U(GT) | (g) GT |

Figure 5. Qualitative comparison of 3D reconstruction on the Mutil-Garment.

| Subject ID | P2S↓ | | | | | Chamfer↓ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | NeRF | NeRF+L | NeRF+U | OURS | NeRF+U(GT) | NeRF | NeRF+L | NeRF+U | OURS | NeRF+U(GT) |
| people1 | 65.53 | 13.57 | 33.51 | **4.09** | 0.86 | 89.32 | 13.96 | 41.81 | **4.25** | 0.25 |
| people2 | 36.26 | 11.67 | 28.50 | **1.55** | 0.85 | 34.95 | 10.78 | 28.86 | **0.96** | 0.25 |
| people3 | 34.78 | 16.01 | 36.40 | **4.17** | 1.17 | 33.62 | 13.83 | 38.36 | **3.30** | 0.43 |
| people4 | 33.29 | 26.84 | 32.74 | **3.53** | 1.06 | 33.70 | 26.59 | 32.08 | **2.68** | 0.36 |
| Average | 42.46 | 17.02 | 33.28 | **3.32** | 0.99 | 47.90 | 16.29 | 34.79 | **2.80** | 0.32 |

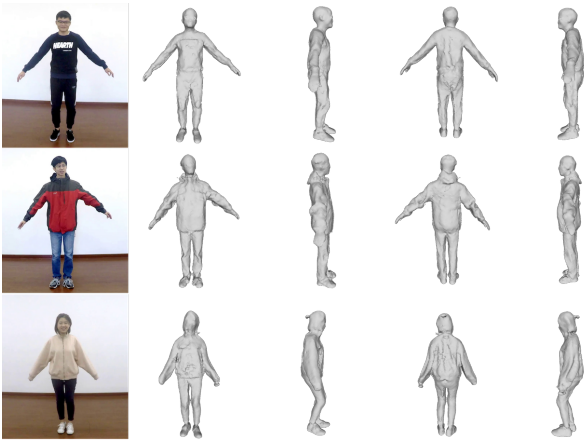Table 2. Quantitative comparison of 3D Reconstruction on Mutil-Garment.



Figure 6. 3D reconstruction results on iPER

of obtaining accurate poses. We also provide qualitative comparisons in Figure 2 with the person examples drawn from the iPER dataset, and Multi-Garment dataset, and we can see that the proposed approach produce more reliable results than implicit control of deformation (NeRF+L). Figure 3 shows the realistic rendering results of more views of the proposed method on more persons with different dresses and hairstyles, indicating the applicability and robustness of the proposed method in real scenarios.

### 4.4. 3D human reconstruction

An quantitative comparsion of different strategies in 3D human reconstruction is shown in Table 2, from which we can see that the proposed approach achieves a much lower P2S and Chamfer distance, demonstrating the superiority of the proposed approach in reconstructing accurate 3D geometry. Figure 5 compares the qualitative results of 3D reconstruction and we can see that NeRF fail to learn reasonable 3D geometry of the human subject with movements, and NeRF+U also produce messy results. Compared with
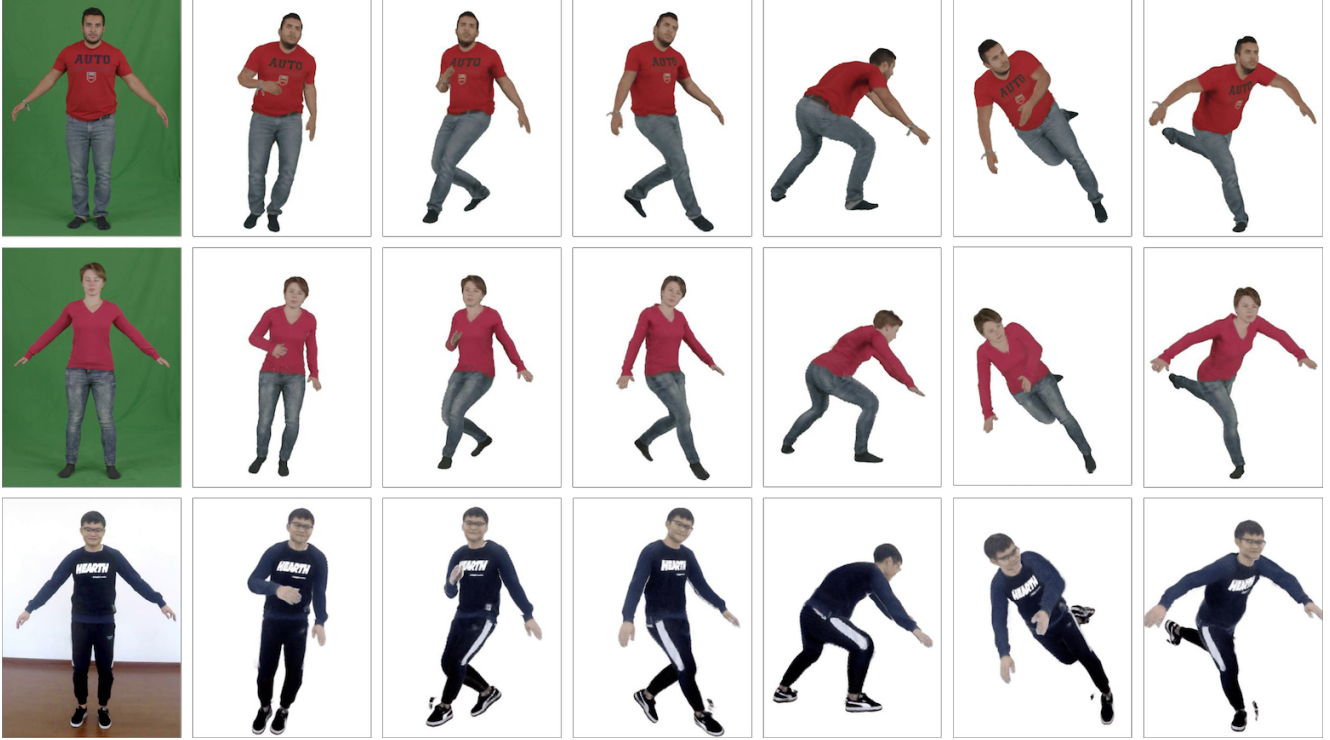
Figure 7. Novel pose synthesis on People-Snapshot and iPER

NeRF+L which produces over-smooth or under-smooth results, the proposed approaches better capture the geometric details in cloth wrinkles, faces and hairs. In Figure 4 we compare the reconstruction results with video avatar [3], which deforms vertices of the SMPL model to fit the 2D human silhouettes over the video sequence. We can see that the implicit learning of subject geometry with animatable neural radiance fields generates better quality of details, including cloth wrinkles, hair, and accessories. In Figure 6 we show the reconstruction results of persons with varied clothes and hairstyles from the iPER dataset, and it is interesting to see that the proposed method is capable of capturing the high-quality 3D geometry details, such as the hood (second line) and pigtail (third line).

| Subject ID | PSNR↑ | | SSIM↑ | | LPIPS↓ | |
|---|---|---|---|---|---|---|
| | NB | OURS | NB | OURS | NB | OURS |
| female3c | 23.92 | 26.97 | 0.96 | 0.97 | 0.110 | 0.076 |
| female4c | 24.61 | 24.92 | 0.96 | 0.96 | 0.079 | 0.058 |
| male3c | 24.99 | 27.21 | 0.95 | 0.96 | 0.089 | 0.058 |
| male4c | 24.88 | 27.45 | 0.96 | 0.97 | 0.093 | 0.064 |

Table 3. Quantitative comparison with NeuralBody(NB)[27] on the People-snapshot dataset.



Figure 8. Visualization results of novel pose synthesis between NeuralBody(top) and Ours(bottom)

## 4.5. Novel pose synthesis

Due to our explicit control of deformation via SMPL, our method can synthesize images under unseen poses even with only simple A-pose sequences as input. As shown in the Figure 7, we provide a qualitative visualization of novel pose synthesis on People-snapshot and iPER. We feed the trained animatable neural radiance fields with different

pose parameters to synthesize images of the performers in the corresponding poses. Despite the significant differences between the input novel poses and the training poses, the results show that our method can still produce realistic images while preserving the identity consistency of the subjects. NeuralBody[27] is the most similar work to ours since it also combines NeRF with SMPL. Compared to ours, it handles complex cloth geometry (which is not modeled by SMPL) better due to its use of latent code. However, every vertex's latent code affects a much larger region after their sparse convolution layers, resulting in unpredictable artifacts for novel pose synthesis (see Figure 8). Table 3 shows that our method achieves much better results than NeuralBody on novel pose synthesis.



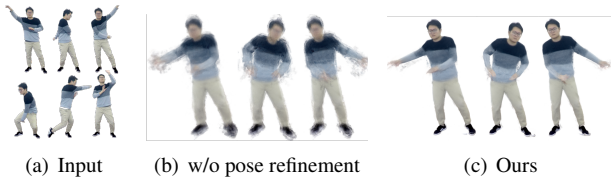    (a) Input    (b) w/o pose refinement    (c) Ours

Figure 9. Visualization result of novel pose synthesis on a complex pose video 009-4-2 of the iPER dataset.

## 5. Limitations

Although our method can reconstruct detailed 3D human body models and render realistic images from monocular videos while the subject rotates while holding an A-pose or T-pose, our method does not take into account the complex poses and the detailed deformation of clothes when the human body is moving. Our input video requires the performer to do simple poses and keep the clothes as unchanged as possible. We also train on videos of complex poses as shown in Figure 9, and our method yields reasonably good results. It is more challenging for videos with complex poses due to the inaccurate SMPL estimation. Our pose refinement strategy significantly improves the negative impact. Our method can't reconstruct invisible parts, such as the underarms and the inner thighs, so the input video needs to cover the whole body of the human body as much as possible. Our explicit pose-guided deformation depends on SMPL, so the shape of two clothes belonging to different body parts should be close to the surface of SMPL, otherwise it will cause deformation ambiguity. We will discuss these limitations in detail in the supplementary material.

## 6. Conclusion

In this paper, we propose to learn an animatable neural radiance field from monocular videos, which allows us to perform visually realistic novel-view synthesis results, reconstruct 3D geometry of the human subject with high-quality details, and animate the human with arbitrary poses.

We extend the neural radiance field to dynamic scenes with human movements via introducing explicit pose-guided deformation and analysis-by-synthesis pose refinement strategy. The pose guided deformation attempts to deform the 3d position along with neighboring SMPL vertices to learn a meaningful human template in the canonical space, as well as to learn accurate 3d geometry for the observation space. The pose refinement strategy compensates for the impact of inaccurate pose estimation from existing approaches, and provides better guidance for learning better geometry and appearance. Experiments on both synthetic data and real data demonstrate the superiority of the proposed approach. Works will be continued to address the above mentioned limitations in the future.

# References

[1] Thiemo Alldieck, Marcus A. Magnor, Bharat Lal Bhatnagar, Christian Theobalt, and Gerard Pons-Moll. Learning to reconstruct people in clothing from a single RGB camera. In *CVPR*, pages 1175–1186, 2019. 2

[2] Thiemo Alldieck, Marcus A. Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. Detailed human avatars from monocular video. In *3DV*, 2018. 1, 2

[3] Thiemo Alldieck, Marcus A. Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. Video based reconstruction of 3d people models. In *CVPR*, pages 8387–8397, 2018. 1, 2, 5, 6, 8

[4] Thiemo Alldieck, Gerard Pons-Moll, Christian Theobalt, and Marcus A. Magnor. Tex2shape: Detailed full human body geometry from a single image. In *ICCV*, pages 2293–2303, 2019. 2

[5] Guha Balakrishnan, Amy Zhao, Adrian V. Dalca, Frédo Durand, and John V. Guttag. Synthesizing images of humans in unseen poses. In *CVPR*, pages 8340–8348, 2018. 2

[6] Bharat Lal Bhatnagar, Garvita Tiwari, Christian Theobalt, and Gerard Pons-Moll. Multi-garment net: Learning to dress 3d people from images. In *ICCV*, pages 5419–5429, 2019. 5

[7] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *SIGGRAPH*, pages 187–194, 1999. 2

[8] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A. Efros. Everybody dance now. In *ICCV*, pages 5932–5941, 2019. 2

[9] Haoqiang Fan, Hao Su, and Leonidas J. Guibas. A point set generation network for 3d object reconstruction from a single image. In *CVPR*, pages 2463–2471, 2017. 6

[10] Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *CVPR*, pages 8649–8658, 2021. 2

[11] Zeng Huang, Tianye Li, Weikai Chen, Yajie Zhao, Jun Xing, Chloe LeGendre, Linjie Luo, Chongyang Ma, and Hao Li. Deep volumetric video from very sparse multi-view performance capture. In *ECCV*, volume 11220, pages 351–369, 2018. 2

[12] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. ARCH: animatable reconstruction of clothed humans. In *CVPR*, pages 3090–3099, 2020. 2

[13] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, pages 5967–5976, 2017. 2

[14] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 5

[15] Muhammed Kocabas, Nikos Athanasiou, and Michael J. Black. VIBE: video inference for human body pose and shape estimation. In *CVPR*, pages 5252–5262, 2020. 2, 4, 6

[16] Verica Lazova, Eldar Insafutdinov, and Gerard Pons-Moll. 360-degree textures of people in clothing from a single image. In *3DV*, pages 643–653, 2019. 2

[17] Zhe Li, Tao Yu, Chuanyu Pan, Zerong Zheng, and Yebin Liu. Robust 3d self-portraits in seconds. In *CVPR*, pages 1341–1350, 2020. 1, 2

[18] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In *NeurIPS*, 2020. 2

[19] Wen Liu, Zhixin Piao, Jie Min, Wenhan Luo, Lin Ma, and Shenghua Gao. Liquid warping GAN: A unified framework for human motion imitation, appearance transfer and novel view synthesis. In *ICCV*, pages 5903–5912, 2019. 2, 5

[20] Wen Liu, Zhixin Piao, Zhi Tu, Wenhan Luo, Lin Ma, and Shenghua Gao. Liquid warping GAN with attention: A unified framework for human image synthesis. *arXiv: 2011.09055*, 2020. 2

[21] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: a skinned multi-person linear model. *ACM Trans. Graph.*, pages 248:1–248:16, 2015. 1, 2, 3

[22] William E. Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *SIGGRAPH*, pages 163–169, 1987. 5

[23] 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. *arXiv: 2008.02268*, 2020. 1

[24] 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*, volume 12346, pages 405–421, 2020. 1, 2, 4, 5, 6

[25] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B. Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Deformable neural radiance fields. *arXiv: 2011.12948*, 2020. 1, 2

[26] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3d hands, face, and body from a single image. In *CVPR*, pages 10975–10985, 2019. 2

[27] 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. 2, 8, 9

[28] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural Radiance Fields for Dynamic Scenes. In *CVPR*, 2020. 2

[29] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Hao Li, and Angjoo Kanazawa. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *ICCV*, 2019. 1, 2

[30] Shunsuke Saito, Tomas Simon, Jason M. Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *CVPR*, 2020. 1, 2

[31] Aliaksandra Shysheya, Egor Zakharov, Kara-Ali Aliev, Renat Bashirov, Egor Burkov, Karim Iskakov, Aleksei

Ivakhnenko, Yury Malkov, Igor Pasechnik, Dmitry Ulyanov, Alexander Vakhitov, and Victor S. Lempitsky. Textured neural avatars. In *CVPR*, pages 2387–2397, 2019. 1

[32] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *NeurIPS*, 2019. 2

[33] Shih-Yang Su, Frank Yu, Michael Zollhöfer, and Helge Rhodin. A-nerf: Surface-free human 3d pose refinement via neural rendering. *arXiv: 2102.06199*, 2021. 2

[34] Shao-Hua Sun, Minyoung Huh, Yuan-Hong Liao, Ning Zhang, and Joseph J. Lim. Multi-view to novel view: Synthesizing novel views with self-learned confidence. In *ECCV*, volume 11207, pages 162–178, 2018. 6

[35] 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 deforming scene from monocular video. *arXiv: 2012.12247*, 2020. 1, 2

[36] Daniel Vlasic, Ilya Baran, Wojciech Matusik, and Jovan Popovic. Articulated mesh animation from multi-view silhouettes. *ACM Trans. Graph.*, 2008. 1

[37] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *CVPR*, 2018. 2

[38] Donglai Xiang, Fabian Prada, Chenglei Wu, and Jessica K. Hodgins. Monoclothcap: Towards temporally coherent clothing capture from monocular RGB video. In *3DV*, 2020. 1, 2

[39] Lu Yang, Qing Song, Zhihui Wang, Mengjie Hu, Chun Liu, Xueshi Xin, Wenhe Jia, and Songcen Xu. Renovating parsing R-CNN for accurate multiple human parsing. In *ECCV*, volume 12357, pages 421–437, 2020. 2

[40] Tao Yu, Jianhui Zhao, Zerong Zheng, Kaiwen Guo, Qionghai Dai, Hao Li, Gerard Pons-Moll, and Yebin Liu. Doublefusion: Real-time capture of human performances with inner body shapes from a single depth sensor. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020. 2

[41] Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai. Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. *arXiv:2007.03858*, abs/2007.03858, 2020. 2, 3

[42] Tiancheng Zhi, Christoph Lassner, Tony Tung, Carsten Stoll, Srinivasa G. Narasimhan, and Minh Vo. Texmesh: Reconstructing detailed human texture and geometry from RGB-D video. In *ECCV*, volume 12355, pages 492–509, 2020. 2