# Animatable Neural Radiance Fields from Monocular RGB Videos

Jianchuan Chen, Ying Zhang, Di Kang, Xuefei Zhe, Linchao Bao, Xu Jia, Huchuan Lu, *Senior Member, IEEE*

*Abstract*—We present animatable neural radiance fields (*animatable NeRF*) 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 novel views, and 3) animation of the human with novel poses.

*Index Terms*—Novel view synthesis, 3D human reconstruction, Image generation, Animation, Neural Rendering, Reconstruction algorithms.

## I. INTRODUCTION

**T**HE 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 [1], and RGB-D sensor [2], limiting their applications in practical scenarios. On the other hand, for various 3D human reconstruction approaches [3], [4], [5], [6], [7], 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, which are one of the most accessible video forms in daily life. Different from existing approaches based on pre-scanned human models [8] or parametric body models [5], [7], [6], our approach implicitly reconstructs the human geometry and appearance via generalizing Neural Radiance Fields (NeRF) [9], 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 [9] 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 [10] and non-rigid deformations [11], [12]. 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 [13], 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 view 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 novel 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.

## II. 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 [3], [4], [14], [15], [16], multi-view images [3], [4], [17], [18], RGB videos [6], [5], [19], [7], [20], or RGB-D videos [2], [21]. One stream of these approaches [5], [6], [7] utilize a parametric body model such as SMPL [13] 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 [3] and PIFuHD [4] 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. To handle more complex pose inputs, these methods [16], [18], [22], [23] combine implicit representations and parametric models to obtain more robust results and are animatable.

**Neural Representations**. Representing a scene with neural networks has achieved stunning success in recent years. SRN [24] proposes an implicit neural representation that assigns feature vectors to 3D positions, and uses a differentiable ray marching algorithm for image generation. NeRF [9] establishes a static scene that maps 3D coordinates and viewing direction to density and color using a neural network. These methods [24], [9], [25] can render very realistic images, but they are all limited to static scenes. Dynamic NeRFs [11], [26], [25], [12], [27], [28] extend NeRF to dynamic scenes by introducing the latent deformation field or scene flow fields. These methods where the deformations are learned by networks allow to handle more general deformation, and synthesize novel poses by using interpolation in the latent space. However, it is difficult to implicitly control the complex non-rigid deformation of human body motion. These works [29], [30], [31], [32], [33], [34] combine scene representation network with parametric models [35], [13] 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 and reconstruct the human body in the canonical pose. At the same time, this explicit method allows us to fine-tune the parameters of the SMPL which is very practical in real scenarios. Similar ideas with us have been used in recent works[30], [36], [37], but these methods usually require multi-view images or accurate registered SMPL. It is more challenging for monocular videos because of the difficulty of SMPL estimation.

**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 [38], [39], [40], [41] use image-to-image translation [42], [43] 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. To better preserve the appearance of the source subject, these methods [44], [45] use UV map to transform features from screen space to UV texture space to obtain the neural texture, then render the feature maps in the target pose by neural rendering network. In addition to these general approaches between arbitrary subjects, there are other person-specific methods [1], [46], [47]. Textured Neural Avatar [1] learns a uniform neural texture from different views and poses. SMPLpix[47] and NHR[46] project the point clouds to 2D images and then feed them into an image-to-image translation network. However, these neural rendering methods fail to generate photorealistic results for novel poses that were not seen during training.

## III. 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 Fig. 1. Given a $n$-frame video sequence $\{I_t\}_{t=1}^n$ of a single human subject turning around before the camera and holding an A-pose or T-pose, we estimate the SMPL [13], [48] parameters $M(\theta_t, \beta_t)$ and camera intrinsics $K_t$ of each frame using existing human body shape and pose estimation models [49]. In order to avoid the influence of background changes caused by camera movement, we first use a segmentation network [50] to obtain the foreground human mask and set the background color to white uniformly. Our animatable NeRF (Sec. III-A) can be decomposed into pose-guided deformation (Sec. III-B) and a neural radiance field (NeRF) defined in the canonical space. We can use the volume rendering (Sec. III-C) 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 (Sec. III-D). We also introduce background regularization and pose regularization to improve the robustness of optimization (Sec. III-E).

### A. Animatable Neural Radiance Fields

To model human appearance and geometry with complex non-rigid deformation, we introduce the parameterized human model SMPL [13] into the neural radiance field and present the animatable neural radiance fields (animatable NeRF) $F$ 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 $\sigma$:

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

where $D(\mathbf{x}, \theta_t, \beta_t)$ 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 the skin and clothes of humans are mainly diffuse reflective materials, so we remove the viewing direction from the input. We will discuss the impact of viewing direction in Sec. V-D.

### B. Pose-guided Deformation

In contrast to [12], [11] that implicitly control of the deformation of spatial points, we use the parametric body model - SMPL, to explicitly guide the deformation of spatial points. Here we define the observed image as the observation space and attempt to learn a template human in the canonical space. The articulated SMPL model enables the explicit transformation of the spatial points (i.e. from observation space to canonical space), which facilitates the learning of a specified meaningful canonical space, and reduces the reliance of diverse input poses to generalize to unseen poses so that we can learn the NeRF space from dynamic scenes (containing moving person) and animate this person after training. The template pose in the canonical space is defined as X-pose $\theta^0$ (as shown in Fig. 1), due to its good visibility and separability of each body part. By using the inverse transformation of the
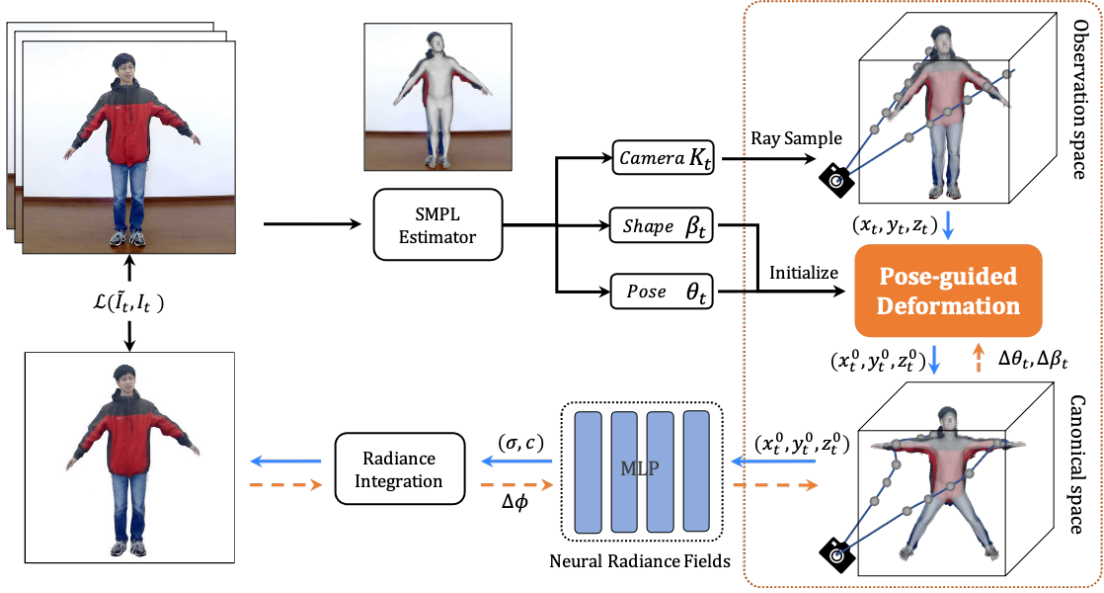
Fig. 1. Overview of the proposed Animatable Neural Radiance Fields. Given a video sequence, we estimate the camera $K_t$ and SMPL parameters $M(\theta_t, \beta_t)$ of the human subject for initialization. We use volume rendering to sample points $(x_t, y_t, z_t)$ along the camera ray in observation space, 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 $\theta_t, \beta_t$ by minimizing the error $\mathcal{L}(\tilde{I}_t, I_t)$ between the rendered image $\tilde{I}_t$ and the ground truth image $I_t$ with the mask.

linear skinning of SMPL, the pose $\theta_t$ in observation space can be transformed into the X-pose $\theta^0$ in canonical space. Considering 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 [18] we define the transformation of a point $\mathbf{x}$ from observation space to canonical space as

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

(2)

where $\mathcal{N}(\mathbf{x})$ 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 $\omega_i$ that the vertex $v_i$ affects the point $\mathbf{x}$ is defined as

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

(3)

where $\mathbf{b}_i$ is the blend weight of $v_i$ and $\hat{\mathbf{b}}$ is the blend 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 blend skinning weight which characterizes the movement patterns of a vertex along with the SMPL joints [13], to strengthen

the movement impact of the nearest neighbor. $\omega$ is used for weight normalization.

Following SMPL[13], the transformation matrix $\mathbf{M}_i(\beta, \theta)$ of mesh vertex $v_i$ from rest pose to $\theta$-pose is computed by

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

(4)

where $\mathbf{G}_j \in \mathbb{R}^{4 \times 4}$ is the world transformation of joint $j$, $b_{i,j}$ is the blend skinning weight representing how much the rotation of part $j$ affects 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.

### C. Volume Rendering

We use the same volume rendering techniques as in NeRF [9] 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 3D *mask* (detailed in the following) used to provide geometric prior guidance and deal with ambiguity during pose-guided deformation.

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;

$$\eta_t(\mathbf{x}_k) = d(\mathbf{x}_k) \le \delta$$
$$d(\mathbf{x}_k) = \sum_{v_i \in \mathcal{N}(\mathbf{x}_k)} \frac{\omega_j}{\omega} \|\mathbf{x}_k - v_i\| \qquad (6)$$

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. In experiments we follow NeRF [9] 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 networks, respectively.

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

Our proposed method learns an animatable NeRF 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 [49] 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 blurry results. To address this problem, we propose to fine-tune the SMPL parameters during training. Specifically, we use VIBE [49] to estimate SMPL parameters $M(\theta_t, \beta_t)$ for each frame $I_t$ as initialization of variables, 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 SMPL can better fit the input image, and helps to obtain clearer and sharper results as shown in Fig. 2 and Fig. 4.

### E. Objective Functions

Given a monocular video sequence, we learn the animatable NeRF by optimizing the following objective function

$$\mathcal{L} = \mathcal{L}_c + \mathcal{L}_p + \lambda_d * \mathcal{L}_d \qquad (7)$$

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

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

where $\mathbf{r}$ is a camera ray passing through the image $I_t$. $C_t(\mathbf{r})$ is the ground truth color of the pixel intersected by ray $\mathbf{r}$ on the observed image $I_t$. And $\tilde{C}_t^c(\mathbf{r})$, $\tilde{C}_t^f(\mathbf{r})$ are the corresponding rendered colors from the coarse and fine networks respectively (see Sec. III-C).

**Pose Regularization**. To obtain stable and smooth pose parameters, we add the following pose regularization term to encourage the optimized pose parameter to stay close to the initial pose, and the pose parameters between adjacent frames to be similar.

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

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

**Background Regularization**. We only focus on reconstructing the human no matter what the background is, which means, ideally, density exists only inside the human. To better achieve this goal, we first set the background color to white with the help of an off-the-shelf segmentation network. We minimize the difference between the rendered integral density and the mask obtained by segmentation. Since the foreground (i.e. human) region is 1 and the background region is 0 in the mask, we are encouraging the human region's integral density to be 1 and encouraging the background region's integral density to be 0, resulting in a much cleaner empty space estimation and more solid and clearer person estimation in our canonical NeRF space. Mathematically, our background regularization term is defined as follows,

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

where $\tilde{D}_t^c$ and $\tilde{D}_t^f$ is the rendered integral 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 and $D_t(\mathbf{r}) = 1$ in the foreground region and $D_t(\mathbf{r}) = 0$ in the background region.

### F. Applications

The proposed approach learns an animatable NeRF, 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. For the original NeRF [9], novel view images (Sec. IV-C) can be rendered through volume rendering, and the surface geometry (Sec. IV-D) of the scene can be extracted with the Marching Cubes algorithm [51]. Since we have explicitly incorporated the SMPL model into the NeRF training process, we can deform the neural radiance field to desired poses for rendering by our pose-guided deformation. This makes our NeRF *animatable*, and thus a new application that can demonstrate new poses or animating the reconstructed people (Sec. IV-E) is enabled as shown in Fig. 8 and Table III.

## IV. EXPERIMENTS

### A. Implementation Details

Following NeRF [9], we use coarse and fine networks to represent the human body, and use 64 coarse and 64 + 32 fine rays samples for all experiments. Focusing on the foreground

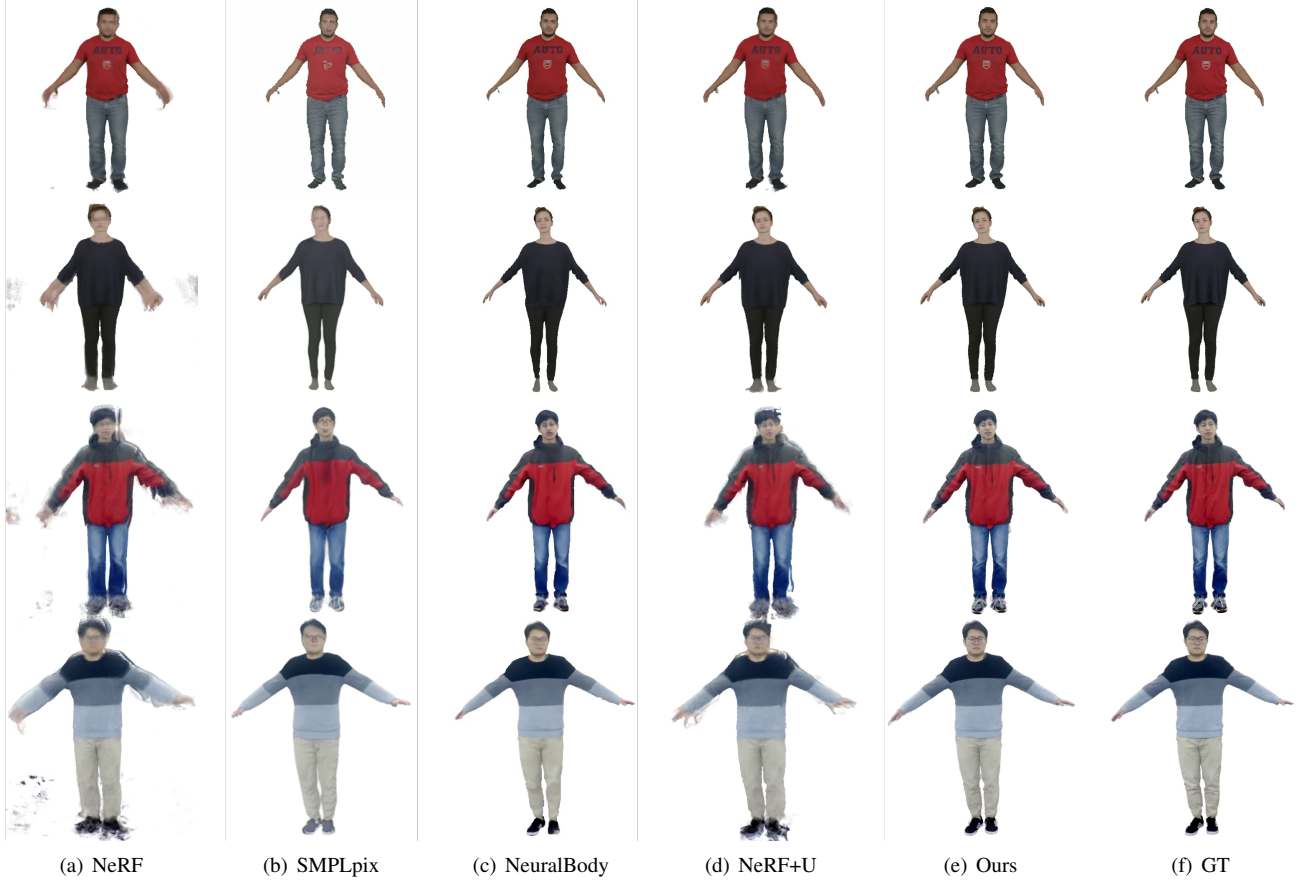| (a) NeRF | (b) SMPLpix | (c) NeuralBody | (d) NeRF+U | (e) Ours | (f) GT |

Fig. 2. Visual comparison of different methods about novel view synthesis on People-snapshot[6](1-2 rows) and iPER[40](3-4 rows). NeRF[9] is struggling to handle dynamic scenes because the movement of the subject violates the multi-view consistency requirement. With the help of our proposed pose-guide deformation, NeRF+U (NeRF + Unpose) achieves much better results (row 1&2) if the estimated SMPL poses are accurate but still produces blurry results (row 3&4) if they are not. Further adding pose refinement (ours) greatly improves the robustness as long as the estimated SMPL pose is reasonably good. Compared with NeuralBody[30] and SMPLpix[47], our approach can produce realistic images with well preserved identity and cloth details.

subject, we set 90% of the rays to be sampled from the foreground, and the remaining 10% 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 512×512 image in all experiments. For training the model, we adopt the Adam optimizer [52], and it spends about 13 hours on 2 Nvidia GeForce RTX 3090 24GB GPUs.

*B. Datasets and Evaluation*

**Datasets**. To evaluate the effectiveness of the proposed method, we conduct experiments on 3 different datasets, including People-Snapshot [6], iPER [40], and Multi-Garment [53]. People-Snapshot[6] and iPER[40] datasets both contain different monocular RGB videos captured in real-world scenes, where the subjects hold an A-pose and turn around before a fixed camera. In addition, iPER dataset also contains videos of the same person with random motion sequences. Multi-Garment [53] dataset contains 3D scanned human body models and textures and the corresponding registered SMPLD models that can be used for animation. We selected 4 human body models to synthesize the videos, according to motion sequences which the subjects rotate while holding an A-pose in People-Snapshot dataset. People-

Snapshot and iPER datasets are mainly used for the evaluation of novel view synthesis and novel pose synthesis experiments. And the synthetic data from Multi-Garment dataset are used to evaluate the quality of the 3D reconstructions.

**Evaluation**. In our experiments, we use A-pose frames (2 circles) for training, and the remaining A-pose frames (1 circle) for testing novel view synthesis and random pose frames for testing novel pose synthesis. Since there are depth and scale ambiguities in optimizing the SMPL parameters for monocular videos, we also optimize the SMPL parameters of the test frames for quantitative evaluation. Note that the parameters of the neural radiance field network remain fixed. For quantitative evaluation, we evaluate our method for novel view synthesis and novel pose synthesis using the following metrics: peak signal-to-noise ratio (PSNR), structural similarity index (SSIM [54]), and learned perceptual image patch similarity (LPIPS [55]). For 3D reconstruction, we use point-to-surface Euclidean distance (P2S) and Chamfer distance [56] (in cm) between the reconstructed and the ground truth surfaces. We register our meshes to ground truth geometry for comparison in consideration of scale and depth ambiguities. The datasets from the real scenes don't have the corresponding ground truth geometry, and we only provide qualitative results.
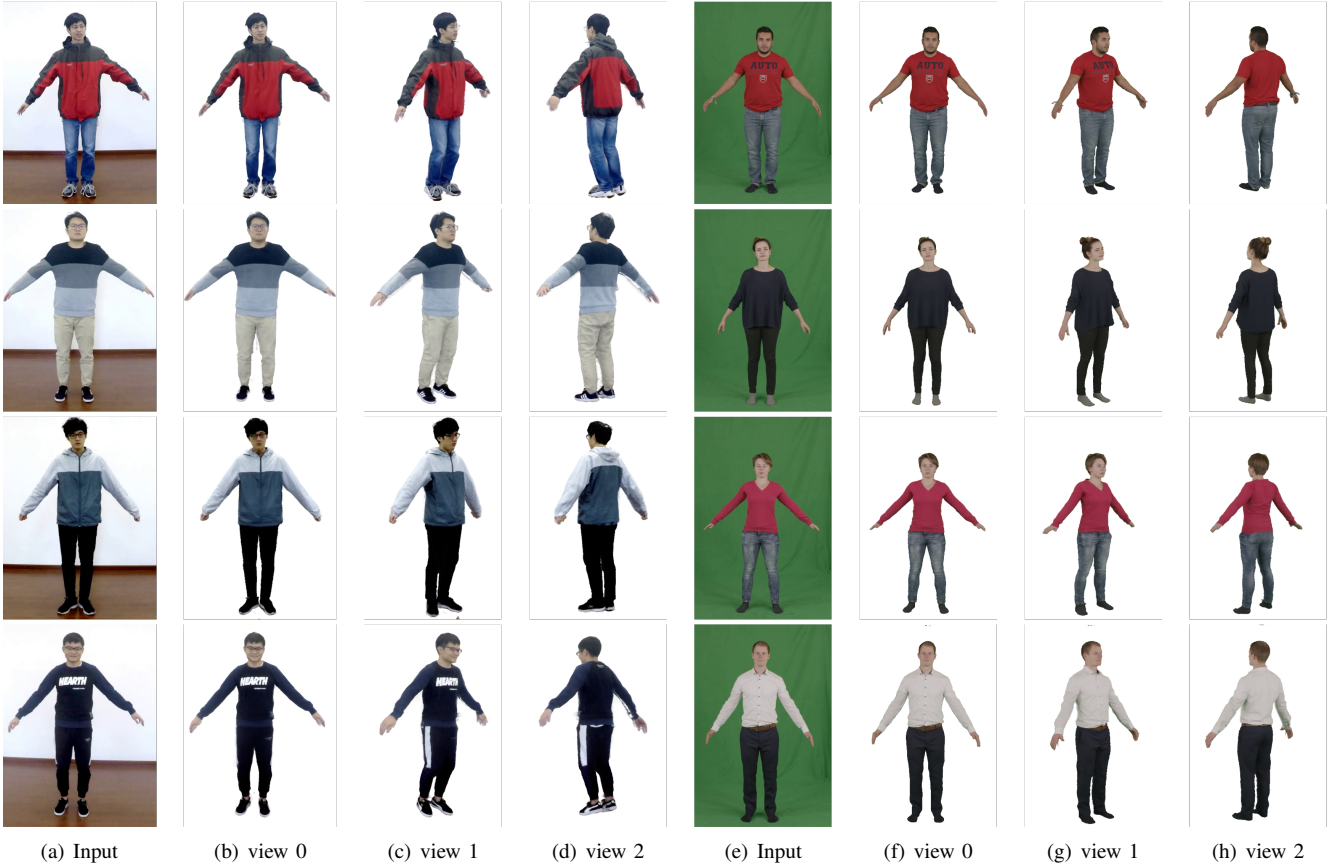
Fig. 3. Results of Novel View Synthesis on iPER (a-d) and People-Snapshot (e-h). Our method can synthesize realistic and multi-view consistent results from different camera views while maintaining the subject pose fixed.

TABLE I
QUANTITATIVE COMPARISON OF NOVEL VIEW SYNTHESIS ON PEOPLE-SNAPSHOT[6] AND iPER[40].

| Subject ID | PSNR↑ | | | | | SSIM↑ | | | | | LIPIS↓ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | NeRF | SMPLpix | NB | NeRF+U | OURS | NeRF | SMPLpix | NB | NeRF+U | OURS | NeRF | SMPLpix | NB | NeRF+U | OURS |
| male-3-casual | 20.64 | 23.74 | 24.94 | 23.88 | **29.37** | .8993 | .9229 | .9428 | .9329 | **.9703** | .1008 | .0222 | .0326 | .0438 | **.0168** |
| male-4-casual | 20.29 | 22.43 | 24.71 | 23.13 | **28.37** | .8803 | .9095 | .9469 | .9276 | **.9605** | .1445 | .0305 | .0423 | .0554 | **.0268** |
| female-3-casual | 17.43 | 22.33 | 23.87 | 22.45 | **28.91** | .8605 | .9288 | .9504 | .9413 | **.9743** | .1696 | .0270 | .0346 | .0498 | **.0215** |
| female-4-casual | 17.63 | 23.35 | 24.37 | 23.13 | **28.90** | .8578 | .9258 | .9451 | .9276 | **.9678** | .1827 | .0239 | .0382 | .0556 | **.0174** |
| iper-009-4-1 | 19.54 | 20.25 | 25.46 | 21.56 | **30.23** | .7870 | .9018 | .9378 | .8667 | **.9466** | .2641 | **.0293** | .0558 | .1197 | .0335 |
| iper-023-1-1 | 17.41 | 19.48 | 25.44 | 20.25 | **27.26** | .7623 | .8945 | .9330 | .8656 | **.9457** | .2769 | .0442 | .0493 | .1109 | **.0285** |
| iper-002-1-1 | 16.01 | 19.64 | 23.06 | 18.75 | **26.99** | .7500 | .8886 | .9394 | .8708 | **.9502** | .3363 | .0392 | .0476 | .1205 | **.0285** |
| iper-026-1-1 | 17.09 | 19.03 | 23.77 | 18.48 | **26.85** | .7580 | .8574 | .9351 | .8623 | **.9542** | .2928 | .0494 | .0550 | .1282 | **.0315** |

## C. Novel view synthesis

Like the original NeRF [9], our animatable NeRF can be rendered from arbitrary views (of the same pose). Since the monocular video does not have corresponding novel view images, we use the first part of the A-pose video frames to train our model and the remaining frames to test the rendered novel view images[1]. To compare against the original NeRF, we first transform the global rotation and translations of the SMPL estimation from every video frame to camera pose as if we are handling a collection of multi-view images of an almost static scene since the subject is holding an A-pose. In order to verify the effectiveness of our pose fine-tuning strategy, we first use VIBE [49] to estimate the parameters of SMPL, and conduct comparative experiments between "Ours" (NeRF + Unpose + Pose Refinement) and "NeRF+U" (NeRF + Unpose) as shown in Fig. 2. We also compare the proposed method with several state-of-the-art (SOTA) methods, including NeuralBody[30](NB) and SMPLpix[47]. NeuralBody, which also combines SMPL and NeRF, is able to reconstruct dynamic human bodies from monocular video. SMPLpix takes SMPL pose as input to generate images via a neural rendering network. Table I quantitatively compares the results of different approaches about novel view synthesis

[1]Technically, this is not a "novel" view apart from slightly different human pose

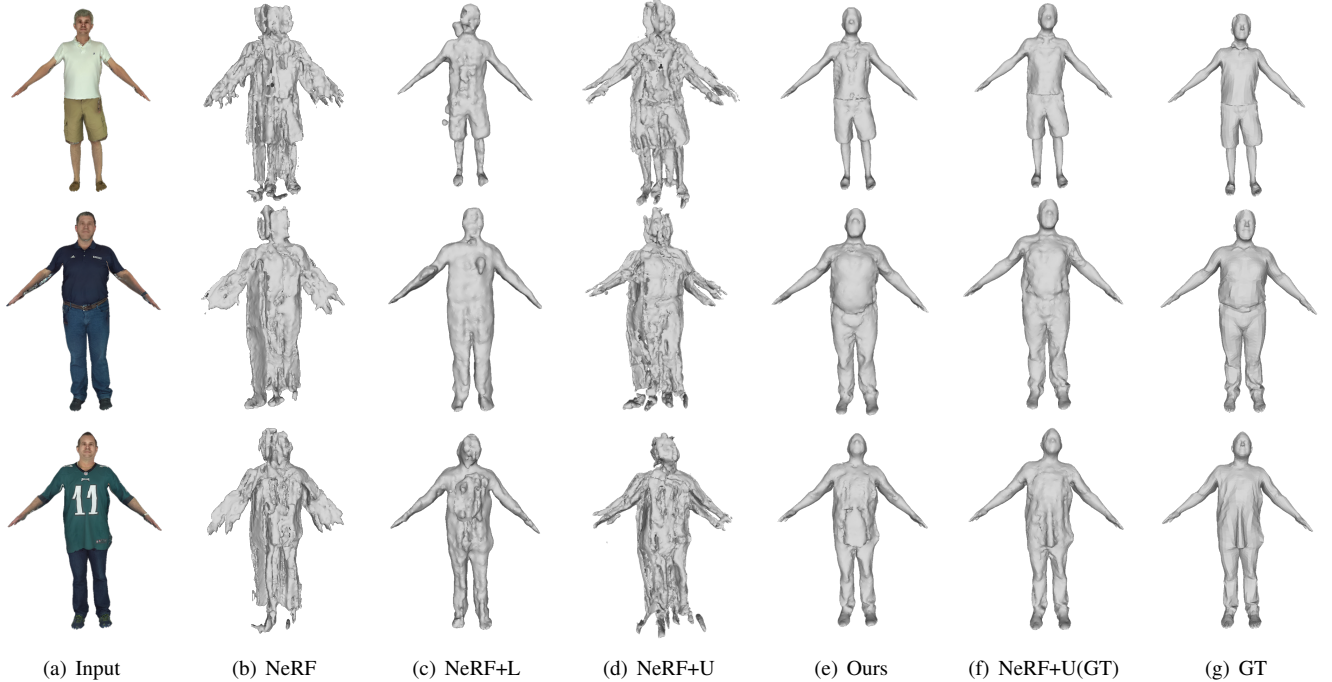|     (a) Input |     (b) NeRF |     (c) NeRF+L |     (d) NeRF+U |     (e) Ours |     (f) NeRF+U(GT) |     (g) GT |

Fig. 4. Visualization of 3D reconstruction on Multi-Garment. NeRF[9] and NeRF+U (NeRF + Unpose) fail to reconstruct 3D geometry due to the movement of the subject and the inaccurate SMPL. Compared with NeRF+L (NeRF + Latent) which produces over-smooth or under-smooth results, our results are more reasonable. As a reference, NeRF+U(GT) uses GT SMPL and learns geometry with very high precision, demonstrating the effectiveness of our pose-guided deformation and showing the importance of obtaining accurate SMPL for 3D reconstruction tasks.

TABLE II
QUANTITATIVE COMPARISON OF 3D RECONSTRUCTION ON MULTI-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 |

on the People-Snapshot and iPER datasets. As described in the table, our proposed approach achieves higher PSNR and SSIM scores compared to other approaches. We also provide qualitative comparisons in Fig. 2 with the person examples drawn from the People-Snapshot the iPER datasets. We can see that the proposed approach produces more realistic and reliable results. NeRF fails to handle such dynamic scenes since the movement of the subject violates the multi-view consistency requirement. Experiments in Fig. 2(d) and Table I (see NeRF+U) also show that inaccurate SMPL parameters cause a very negative impact. In contrast, after taking pose refinement into consideration, the quality of novel view synthesis has been significantly improved. As shown in Fig. 3(b)(c), SMPLpix and NeuralBody seem to overfit the training frames, while our results better preserve the details such as faces and hands. Fig. 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.

### D. 3D human reconstruction

On this task, we compare against the original NeRF [9] and NeRF+L baselines. NeRF+L extends NeRF to condition it on a (per-frame) learnable latent deformation code to handle dynamic scenes as shown in Fig 4(c) and Table II. For synthetic data, we also show results using ground truth SMPL parameters (NeRF+U(GT)) for pose-guided deformation as the upper bound as shown in Fig 4(f) and Table II. Quantitative compassion of different strategies in 3D human reconstruction is shown in Table II. We can see that the proposed approach achieves much lower P2S and Chamfer distances, demonstrating the superiority of the proposed approach in reconstructing accurate 3D geometry. Fig. 4 compares the qualitative results of 3D reconstruction. We can see that NeRF fails to learn reasonable 3D geometry of the human subject with *small* movements. NeRF+U(see Sec. IV-C) also produces messy results. Compared with NeRF+L, which produces over-smooth or under-smooth results, the proposed approaches better capture the geometric details such as cloth wrinkles,

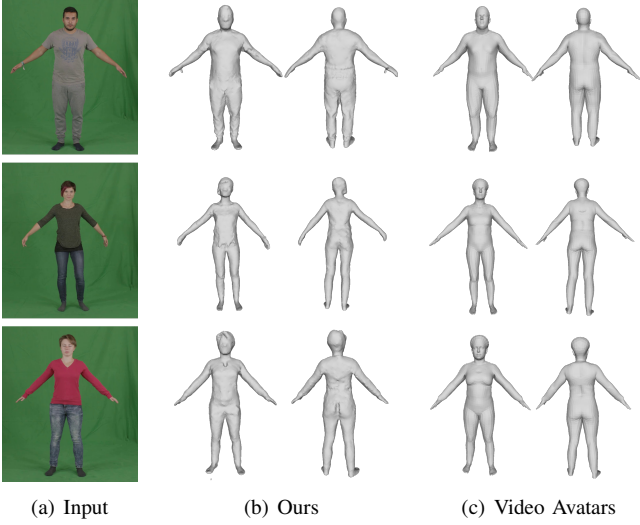(a) Input       (b) Ours       (c) Video Avatars

Fig. 5. Comparisons of 3D reconstruction results on People-Snapshot with video avatars [6]. Compared with Video Avatars[6], our approach can generate more details such as hairs and clothes wrinkles.
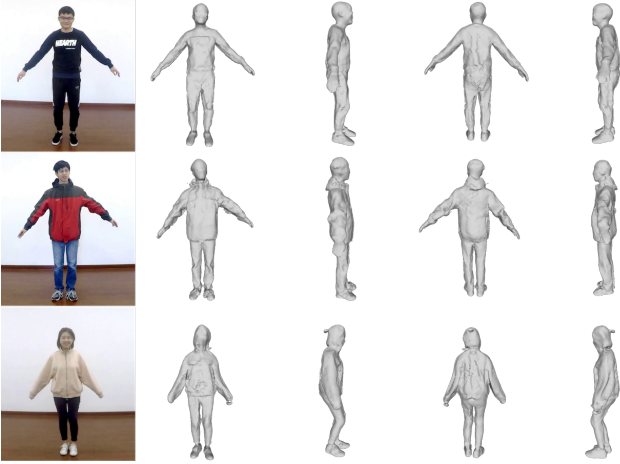


Fig. 6. Visualization of our reconstructed geometry on iPER from different views.



Fig. 7. Comparisons between NeuralBody[30] (first row) and Ours (second row) on novel pose synthesis task. In contrast to NeuralBody, which fails to synthesize novel poses, our approach generalizes better on novel poses that are very different from the training poses.

### E. 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. To quantitatively evaluate the capability of our method on novel pose synthesis, we train the model using A-pose videos and test it using random pose videos of the same person. NeuralBody[30] is the most similar work to ours in the sense that 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, each vertex's latent code would affect a much larger region after several sparse convolution layers, resulting in unpredictable artifacts for novel pose synthesis (see Fig. 7). Table III shows that our method achieves much better results than NeuralBody on novel pose synthesis. Qualitative visualizations of novel pose synthesis on People-snapshot and iPER are provided in Fig. 8. Specifically, different poses are fed into the trained animatable NeRF to obtain the aforementioned renderings. Despite the significant differences between the test novel poses and the training poses, the results show that our method can still produce realistic images with well preserved identity and cloth details of the subjects.

faces, and hairs. With ground truth SMPL parameters, the P2S and Chamfer distance are much lower than all the approaches, which demonstrates the necessity of obtaining accurate poses and the effectiveness of our approximated pose-guided deformation. In Fig. 5, we compare the reconstruction results with video avatar [6], 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 NeRF generates a better quality of details, including cloth wrinkles, hair, and accessories. In Fig. 6 we show the reconstruction results of persons with varied clothes and hairstyles from iPER dataset. Although our pose-guided deformation does not take the deformation of clothes into account, the proposed method is capable of capturing the high-quality 3D geometry details, such as the hood (second line) and pigtail (third line), as long as the clothes do not have violent deformation.

TABLE III
QUANTITATIVE COMPARISON ABOUT NOVEL POSE SYNTHESIS WITH
NEURALBODY(NB)[30] ON THE iPER DATASET.

| Subject ID | PSNR↑ | | SSIM↑ | | LPIPS↓ | |
|---|---|---|---|---|---|---|
| | NB | OURS | NB | OURS | NB | OURS |
| iper-009-4-2 | 20.95 | **24.11** | **.9035** | .8927 | .0980 | **.0782** |
| iper-023-1-2 | 20.28 | **21.98** | **.9009** | .8940 | .0870 | **.0644** |
| iper-026-1-2 | 17.42 | **19.27** | **.8795** | .8713 | .1192 | **.0990** |
| iper-002-1-2 | 19.07 | **23.47** | .8957 | **.9165** | .0749 | **.0483** |

Fig. 8. Novel pose synthesis on People-Snapshot[6] and iPER[40]. We can feed novel SMPL pose parameters to the trained animatable NeRF to synthesize novel pose images. Although trained only on A-pose images, our animatable NeRF has the capability to stably render new images containing complex poses.

## V. DISCUSSION

In the following, we will discuss the proposed approach in details from the following aspects: analysis of pose refinement (Sec. V-A) and canonical poses (Sec. V-B), Impact of background regularization (Sec. V-C) and view direction (Sec. V-D).

### A. Analysis of Pose Refinement

Here we discuss the impact of pose refinement on our approach. Our method relies on SMPL parameters for explicit deformation, so inaccurate SMPL estimation may lead to catastrophic results as shown in Fig. 2(d). We initialize the SMPL parameters with estimations from VIBE[49], which is the state-of-the-art pose and shape estimation method. However, pose estimation from monocular videos usually suffers from depth ambiguity. As shown in Fig. 9(a), although our input video is a simple A-pose, the SMPL estimated by VIBE is usually misaligned at the foot joints. After pose refinement, the SMPL model is better fitting to the input image as shown in Fig. 9(b).

### B. Analysis of Canonical Poses

In this section, we discuss the effect of using different canonical poses in canonical space on the reconstruction results. Since the pose-guided deformation is explicitly based on SMPL[13], we will get different canonical spaces with different canonical poses. Therefore, the choice of canonical poses has a crucial impact on the reconstruction and novel pose synthesis. Here we will discuss the reconstruction results of three different canonical poses: A-pose, T-pose, and X-pose.
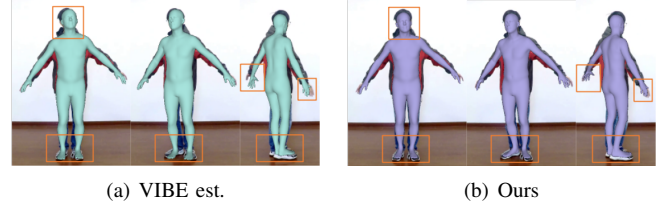


(a) VIBE est.          (b) Ours

Fig. 9. Visual comparison before and after pose refinement on iPER[40]. After pose refinement, the SMPL model is better aligned with the input image (e.g. foot joints).

A-pose is the average pose of the body poses of the training frames, which is the closest to the poses in the training frames. T-pose is the SMPL model's rest pose, where the arms are far away from the body, but the legs are closer to each other. In comparison, our customized X-pose offers more spread body parts (see Fig. 10(d)).

As shown in Fig. 10, using A-pose as the canonical pose offers the best quality for canonical space NeRF learning, while using T-pose and X-pose as canonical poses result in some artifacts under the axilla and the thighs. If a point is close to two different SMPL body parts (e.g. body and arm, two legs), it is hard to decide which part the point belongs to since SMPL models unclothed human body only without taking the offset of the clothes into consideration.

For reconstruction, A-pose is the best choice, but X-pose is more suitable for new pose synthesis. As shown in Fig. 11. When using A-pose and T-pose as canonical poses for synthesizing a pose different from the poses in the training frames, there exist some unacceptable artifacts (e.g. multiple
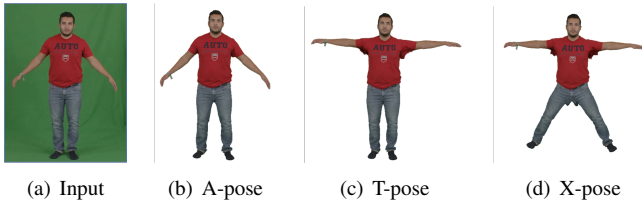
(a) Input  (b) A-pose  (c) T-pose  (d) X-pose

Fig. 10. Visualization of different canonical NeRF spaces with different canonical poses during training on People-Snapshot[6].

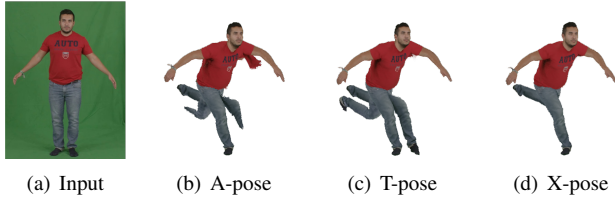

(a) Input  (b) A-pose  (c) T-pose  (d) X-pose

Fig. 11. Novel pose synthesis with different canonical poses on People-Snapshot[6]. X-pose is a better choice for novel pose synthesis compared to A-pose and T-pose, which produce unacceptable artifacts (e.g. multiple legs).

legs). This is because the different body parts are too close to each other in the canonical space of A-pose or T-pose, so that one body part (e.g. left leg) will be deformed by the transformation of another body part (e.g. right leg), resulting in multiple legs in Fig. 11(b)(c).

### C. Impact of Background Regularization

In this section, we discuss the benefits of background regularization for our approach. Our approach focuses on the animatable NeRF of the human from a monocular video. To avoid possible negative influence from the background, we set the background to white uniformly (with the help of an off-the-shelf segmentation network). However, it is common to appear some noisy density regions in the empty space (i.e. non-human space) after training. As shown in Fig. 12(b), although the background of the image is white, we notice some noisy non-zero density regions from the depth map. This is because there is an ambiguity between the background (white in our case) and the cloth which happens to be the same color as the background. To deal with this problem, we introduce background regularization to encourage the density of the background region to be zero. With background regularization, the artifacts in the empty space are significantly reduced as shown in Fig. 12(c).



(a) Input  (b) w/o background reg.  (c) w/ background reg.

Fig. 12. Impact of background regularization on iPER[40]. The background regularization can effectively reduce the artifacts in the background region.



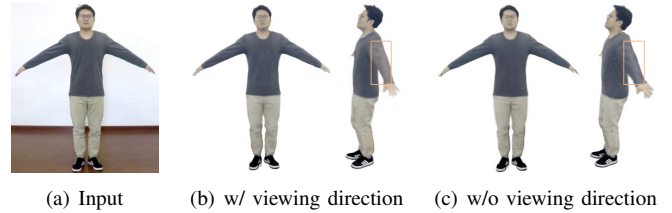(a) Input  (b) w/ viewing direction  (c) w/o viewing direction

Fig. 13. Impact of viewing direction on novel view synthesis on iPER. After removing viewing direction from the input, our model produces more consistent result across different views.
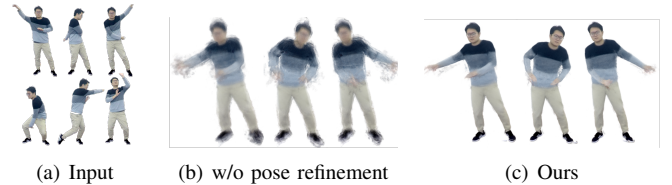


(a) Input  (b) w/o pose refinement  (c) Ours

Fig. 14. Visualization result of novel pose synthesis on a complex pose video 009-4-2 of the iPER[40].

### D. Impact of viewing direction

Unlike NeRF[9], which maps the 3D position and viewing direction to color and density, our approach excludes viewing direction from our input for robust dynamic reconstruction. NeRF's viewing direction is mainly used to handle specular reflections for materials such as glass and metal. For dynamic scenes, it is very difficult to deal with changes of illumination, so we assume that the appearance of the subject is not view-dependent in our experiment. Also, human skin and clothes are mainly diffuse reflective materials, and there exist very few specular reflections. As shown in the Fig. 13, the novel view synthesis results generated with viewing direction as a condition during training show unpredictable artifacts.

## VI. LIMITATIONS

Our method reconstructs a detailed 3D human body model and renders realistic images from a monocular video. Typically, the training videos capture subjects turning around before the camera and holding an A-pose or T-pose. When trained on a video containing complex poses, our method still obtains reasonable results (Fig. 14). But noticeable losses of details are observed compared to previous simple training videos. The main reason is that it is more challenging to obtain accurate enough SMPL estimations for videos containing complex poses. Another limitation is that it is difficult for our method to handle extremely loose clothes or complex non-rigid deformations of the garments, because our explicit pose-guided deformation associates spatial points to SMPL mesh, without explicit modeling of the garments. Thus, the novel view synthesis results inevitably lose some details on the clothes. So, to the get best results, the performer should slowly turn around and hold a simple pose so that their clothes almost remain still relative to their body for high-quality rendering. Like all NeRF based methods trained for only one scene, Our method cannot 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.

## VII. Conclusion

In this paper, we propose to learn an animatable neural radiance field from a monocular video, which allows us to perform photo-realistic novel-view synthesis, reconstruct the 3D geometry of the person with high-quality details, and animate the person with novel poses. To achieve these goals, we extend the neural radiance field to dynamic scenes with human movements via introducing an explicit pose-guided deformation module and an analysis-by-synthesis pose refinement strategy. Specifically, the pose-guided deformation attempts to deform the 3d position according to the neighboring SMPL vertices to learn a good and controllable human template in the canonical space, as well as to learn accurate 3d geometry. The pose refinement strategy compensates for the negative impact of inaccurate pose estimation from existing approaches and provides more consistent guidance for learning better geometry (i.e. density) and appearance (i.e. RGB). Experiments on both synthetic data and real data demonstrate the effectiveness of the proposed approach.

## References

[1] A. Shysheya, E. Zakharov, K. Aliev, R. Bashirov, E. Burkov, K. Iskakov, A. Ivakhnenko, Y. Malkov, I. Pasechnik, D. Ulyanov, A. Vakhitov, and V. S. Lempitsky, "Textured neural avatars," in *CVPR*, 2019, pp. 2387–2397. 1, 2

[2] Z. Li, T. Yu, C. Pan, Z. Zheng, and Y. Liu, "Robust 3d self-portraits in seconds," in *CVPR*, 2020, pp. 1341–1350. 1

[3] S. Saito, Z. Huang, R. Natsume, S. Morishima, H. Li, and A. Kanazawa, "Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization," in *ICCV*, 2019. 1

[4] S. Saito, T. Simon, J. M. Saragih, and H. Joo, "Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization," in *CVPR*, 2020. 1

[5] T. Alldieck, M. A. Magnor, W. Xu, C. Theobalt, and G. Pons-Moll, "Detailed human avatars from monocular video," in *3DV*, 2018. 1

[6] T. Alldieck, M. A. Magnor, W. Xu, C. Theobalt, and G. Pons-Moll, "Video based reconstruction of 3d people models," in *CVPR*, 2018, pp. 8387–8397. 1, 5, 6, 8, 9, 10

[7] D. Xiang, F. Prada, C. Wu, and J. K. Hodgins, "Monoclothcap: Towards temporally coherent clothing capture from monocular RGB video," in *3DV*, 2020. 1

[8] D. Vlasic, I. Baran, W. Matusik, and J. Popovic, "Articulated mesh animation from multi-view silhouettes," *ACM Trans. Graph.*, 2008. 1

[9] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," in *ECCV*, vol. 12346, 2020, pp. 405–421. 1, 2, 3, 4, 5, 6, 7, 10

[10] R. Martin-Brualla, N. Radwan, M. S. M. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth, "Nerf in the wild: Neural radiance fields for unconstrained photo collections," in *CVPR*, 2021. 1

[11] K. Park, U. Sinha, J. T. Barron, S. Bouaziz, D. B. Goldman, S. M. Seitz, and R. Martin-Brualla, "Deformable neural radiance fields," *ICCV*, 2021. 1, 2

[12] E. Tretschk, A. Tewari, V. Golyanik, M. Zollhöfer, C. Lassner, and C. Theobalt, "Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a deforming scene from monocular video," in *ICCV*, 2021. 1, 2

[13] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black, "SMPL: a skinned multi-person linear model," *ACM Trans. Graph.*, pp. 248:1–248:16, 2015. 1, 2, 3, 9

[14] V. Lazova, E. Insafutdinov, and G. Pons-Moll, "360-degree textures of people in clothing from a single image," in *3DV*, 2019, pp. 643–653. 1

[15] T. Alldieck, G. Pons-Moll, C. Theobalt, and M. A. Magnor, "Tex2shape: Detailed full human body geometry from a single image," in *ICCV*, 2019, pp. 2293–2303. 1

[16] Z. Huang, Y. Xu, C. Lassner, H. Li, and T. Tung, "ARCH: animatable reconstruction of clothed humans," in *CVPR*, 2020, pp. 3090–3099. 1, 2

[17] Z. Huang, T. Li, W. Chen, Y. Zhao, J. Xing, C. LeGendre, L. Luo, C. Ma, and H. Li, "Deep volumetric video from very sparse multi-view performance capture," in *ECCV*, vol. 11220, 2018, pp. 351–369. 1

[18] Z. Zheng, T. Yu, Y. Liu, and Q. Dai, "Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction," *TPAMI*, 2021. 1, 2, 3

[19] T. Alldieck, M. A. Magnor, B. L. Bhatnagar, C. Theobalt, and G. Pons-Moll, "Learning to reconstruct people in clothing from a single RGB camera," in *CVPR*, 2019, pp. 1175–1186. 1

[20] T. Yu, J. Zhao, Z. Zheng, K. Guo, Q. Dai, H. Li, G. Pons-Moll, and Y. Liu, "Doublefusion: Real-time capture of human performances with inner body shapes from a single depth sensor," *TPMAI*, 2020. 1

[21] T. Zhi, C. Lassner, T. Tung, C. Stoll, S. G. Narasimhan, and M. Vo, "Texmesh: Reconstructing detailed human texture and geometry from RGB-D video," in *ECCV*, vol. 12355, 2020, pp. 492–509. 1

[22] Z. Yang, S. Wang, S. Manivasagam, Z. Huang, W.-C. Ma, X. Yan, E. Yumer, and R. Urtasun, "S3: Neural shape, skeleton, and skinning fields for 3d human modeling," in *CVPR*, 2021, pp. 13 284–13 293. 2

[23] B. L. Bhatnagar, C. Sminchisescu, C. Theobalt, and G. Pons-Moll, "Combining implicit function learning and parametric models for 3d human reconstruction," in *ECCV*, vol. 12347, 2020, pp. 311–329. 2

[24] V. Sitzmann, M. Zollhöfer, and G. Wetzstein, "Scene representation networks: Continuous 3d-structure-aware neural scene representations," in *NeurIPS*, 2019. 2

[25] L. Liu, J. Gu, K. Z. Lin, T. Chua, and C. Theobalt, "Neural sparse voxel fields," in *NeurIPS*, 2020. 2

[26] A. Pumarola, E. Corona, G. Pons-Moll, and F. Moreno-Noguer, "D-NeRF: Neural Radiance Fields for Dynamic Scenes," in *CVPR*, 2020. 2

[27] Z. Li, S. Niklaus, N. Snavely, and O. Wang, "Neural scene flow fields for space-time view synthesis of dynamic scenes," in *CVPR*, 2021. 2

[28] K. Park, U. Sinha, P. Hedman, J. T. Barron, S. Bouaziz, D. B. Goldman, R. Martin-Brualla, and S. M. Seitz, "Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields," *CoRR*, 2021. 2

[29] G. Gafni, J. Thies, M. Zollhöfer, and M. Nießner, "Dynamic neural radiance fields for monocular 4d facial avatar reconstruction," in *CVPR*, 2021, pp. 8649–8658. 2

[30] S. Peng, Y. Zhang, Y. Xu, Q. Wang, Q. Shuai, H. Bao, and X. Zhou, "Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans," in *CVPR*, 2021. 2, 5, 6, 8

[31] S. Su, F. Yu, M. Zollhöfer, and H. Rhodin, "A-nerf: Surface-free human 3d pose refinement via neural rendering," *arXiv: 2102.06199*, 2021. 2

[32] S. Saito, J. Yang, Q. Ma, and M. J. Black, "SCANimate: Weakly supervised learning of skinned clothed avatar networks," in *CVPR*, 2021. 2

[33] E. Corona, A. Pumarola, G. Alenyà, G. Pons-Moll, and F. Moreno-Noguer, "Smplicit: Topology-aware generative model for clothed people," in *CVPR*, 2021. 2

[34] B. L. Bhatnagar, C. Sminchisescu, C. Theobalt, and G. Pons-Moll, "Loopreg: Self-supervised learning of implicit surface correspondences, pose and shape for 3d human mesh registration," in *NeurIPS*, 2020. 2

[35] V. Blanz and T. Vetter, "A morphable model for the synthesis of 3d faces," in *SIGGRAPH*, 1999, pp. 187–194. 2

[36] S. Peng, J. Dong, Q. Wang, S. Zhang, Q. Shuai, H. Bao, and X. Zhou, "Animatable neural radiance fields for human body modeling," *CoRR*, vol. abs/2105.02872, 2021. 2

[37] L. Liu, M. Habermann, V. Rudnev, K. Sarkar, J. Gu, and C. Theobalt, "Neural actor: Neural free-view synthesis of human actors with pose control," *CoRR*, vol. abs/2106.02019, 2021. 2

[38] G. Balakrishnan, A. Zhao, A. V. Dalca, F. Durand, and J. V. Guttag, "Synthesizing images of humans in unseen poses," in *CVPR*, 2018, pp. 8340–8348. 2

[39] C. Chan, S. Ginosar, T. Zhou, and A. A. Efros, "Everybody dance now," in *ICCV*, 2019, pp. 5932–5941. 2

[40] W. Liu, Z. Piao, J. Min, W. Luo, L. Ma, and S. Gao, "Liquid warping GAN: A unified framework for human motion imitation, appearance transfer and novel view synthesis," in *ICCV*, 2019, pp. 5903–5912. 2, 5, 6, 9, 10

[41] W. Liu, Z. Piao, Z. Tu, W. Luo, L. Ma, and S. Gao, "Liquid warping GAN with attention: A unified framework for human image synthesis," *arXiv: 2011.09055*, 2020. 2

[42] P. Isola, J. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *CVPR*, 2017, pp. 5967–5976. 2

[43] T. Wang, M. Liu, J. Zhu, A. Tao, J. Kautz, and B. Catanzaro, "High-resolution image synthesis and semantic manipulation with conditional gans," in *CVPR*, 2018. 2

[44] K. Sarkar, D. Mehta, W. Xu, V. Golyanik, and C. Theobalt, "Neural re-rendering of humans from a single image," in *ECCV*, vol. 12356, 2020, pp. 596–613. 2

[45] K. Sarkar, L. Liu, V. Golyanik, and C. Theobalt, "Humangan: A generative model of humans images," *CoRR*, vol. abs/2103.06902, 2021. 2

[46] M. Wu, Y. Wang, Q. Hu, and J. Yu, "Multi-view neural human rendering," in *CVPR*, 2020, pp. 1679–1688. 2

[47] S. Prokudin, M. J. Black, and J. Romero, "Smplpix: Neural avatars from 3d human models," in *WACV*, 2021, pp. 1809–1818. 2, 5, 6

[48] G. Pavlakos, V. Choutas, N. Ghorbani, T. Bolkart, A. A. A. Osman, D. Tzionas, and M. J. Black, "Expressive body capture: 3d hands, face, and body from a single image," in *CVPR*, 2019, pp. 10 975–10 985. 2

[49] M. Kocabas, N. Athanasiou, and M. J. Black, "VIBE: video inference for human body pose and shape estimation," in *CVPR*, 2020, pp. 5252–5262. 2, 4, 6, 9

[50] L. Yang, Q. Song, Z. Wang, M. Hu, C. Liu, X. Xin, W. Jia, and S. Xu, "Renovating parsing R-CNN for accurate multiple human parsing," in *ECCV*, vol. 12357, 2020, pp. 421–437. 2

[51] W. E. Lorensen and H. E. Cline, "Marching cubes: A high resolution 3d surface construction algorithm," in *SIGGRAPH*, 1987, pp. 163–169. 4

[52] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015. 5

[53] B. L. Bhatnagar, G. Tiwari, C. Theobalt, and G. Pons-Moll, "Multi-garment net: Learning to dress 3d people from images," in *ICCV*, 2019, pp. 5419–5429. 5

[54] S. Sun, M. Huh, Y. Liao, N. Zhang, and J. J. Lim, "Multi-view to novel view: Synthesizing novel views with self-learned confidence," in *ECCV*, vol. 11207, 2018, pp. 162–178. 5

[55] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *CVPR*, 2018, pp. 586–595. 5

[56] H. Fan, H. Su, and L. J. Guibas, "A point set generation network for 3d object reconstruction from a single image," in *CVPR*, 2017, pp. 2463–2471. 5