

NeuMan: Neural Human Radiance Field from a Single Video

Wei Jiang^{1,2}, Kwang Moo Yi^{1,2}, Golnoosh Samei¹,
Oncel Tuzel¹, and Anurag Ranjan¹

¹ Apple

² The University of British Columbia

{wjiang7, golnoosh, otuzel, anuragr}@apple.com, kmyi@cs.ubc.ca

Abstract. Photorealistic rendering and reposing of humans is important for enabling augmented reality experiences. We propose a novel framework to reconstruct the human and the scene that can be rendered with novel human poses and views from just a single in-the-wild video. Given a video captured by a moving camera, we train two NeRF models: a human NeRF model and a scene NeRF model. To train these models, we rely on existing methods to estimate the rough geometry of the human and the scene. Those rough geometry estimates allow us to create a warping field from the observation space to the canonical pose-independent space, where we train the human model in. Our method is able to learn subject specific details, including cloth wrinkles and accessories, from just a 10 seconds video clip, and to provide high quality renderings of the human under novel poses, from novel views, together with the background.

1 Introduction

The quality of novel view synthesis has been dramatically improved since the introduction of Neural Radiance Fields (NeRF) [27]. While originally proposed to reconstruct a static scene with a set of posed images, it has since been quickly extended to dynamic scenes [29,15,30] and uncalibrated scenes [45,17]. Recent efforts also focus on animation of these radiance field models [19,32,31,12,40] of human, with the aid of large controlled datasets, further extending the application domain of radiance-field-based modeling to enable augmented reality experiences.

In this work, we are interested in the scenario where only one single video is provided, and our goal is to reconstruct the human model and the static scene model, and enable novel pose rendering of the human, without any expensive multi-cameras setups or manual annotations. However, even with the recent advancements in NeRF methods, this is far from being trivial. Existing methods [33,20] require multi-cameras setup, consistent lighting and exposure, clean backgrounds, and accurate human geometry to train the NeRF models. As shown in Table 1, HyperNeRF[30] models a dynamic scene based on a single video, but cannot be driven by human poses. ST-NeRF [12] reconstructs each individual

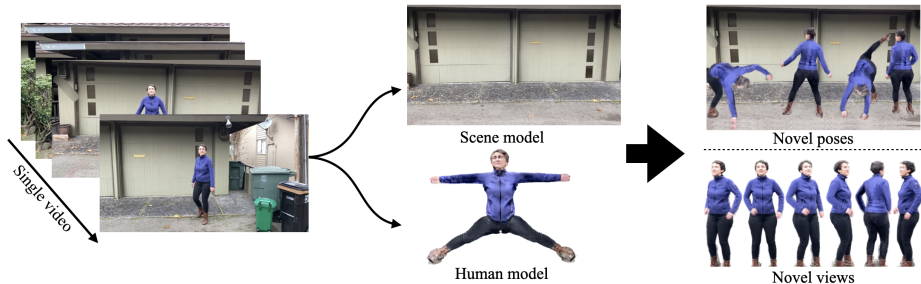


Fig. 1: **Teaser** – We train neural human and scene radiance fields from a single in-the-wild video to enable rendering with novel human poses and novel views.

	Scene Background	Novel Poses	Single Video	Compositionality
HyperNeRF [30]	✓	✗	✓	✗
ST-NeRF [12]	✓	✗	✗	✓
Neural Actor [19]	✗	✓	✗	✗
HumanNeRF [47]	✗	✗	✓	✗
Vid2Actor [46]	✗	✓	✓	✗
Ours	✓	✓	✓	✓

Table 1: **Method capacity comparison** – We illustrate what existing methods are capable of. Compared to other methods, ours is the only one with the ability to render both the scene and the reposed human from a single video.

with a time dependent NeRF model from multiple cameras, but the editing is limited to the transformation of the bounding box. Neural Actor [19] can generate novel poses of a human but requires multiple videos. HumanNeRF [47] builds a human model based on a single video with manually annotated masks, but doesn’t show generalization to novel poses. Vid2Actor [46] generates novel poses of a human with a model trained on a single video but cannot model the background. We address these problems by introducing *NeuMan*, that reconstructs both the human and the scene with the ability to render novel human poses and novel views, from a single in-the-wild video.

NeuMan is a novel framework for training NeRF models for both the human and the scene, which allows high-quality pose-driven rendering as shown in Figure 1. Given a video captured by a moving camera, we first estimate the human pose, human shape, human masks, as well as the camera poses, sparse scene model, and depth maps using conventional off-the-shelf methods [38,42,10,3,8,25].

We then train two NeRF models, one for the human and one for the scene guided by the segmentation masks estimated from Mask-RCNN [10]. Additionally, we regularize the scene NeRF model by fusing together depth estimates from both multi-view reconstruction [38] and monocular depth regression [25].

We train the human NeRF model in a pose independent canonical volume guided by a statistical human shape and pose model, SMPL [23] following Liu et al. [20]. We refine the SMPL estimates from ROMP [42] to better serve the training. However, these refined estimates are still not perfect. Therefore, we jointly optimize the SMPL estimates together with the human NeRF model in an end-to-end fashion. Furthermore, since our static canonical human NeRF cannot represent the dynamics that is not captured by the SMPL model, we introduce an error-correction network to counter it. The SMPL estimates and the error-correction network are jointly optimized during the training.

In summary,

- we propose a framework for neural rendering of a human and a scene from a single video without any extra devices or annotations;
- we show that our method allows high quality rendering of human under novel poses, from novel views, together with the scene;
- we introduce an end-to-end SMPL optimization and an error-correction network to enable training with erroneous estimates of the human geometry;
- our approach allows for the composition of the human and the scene NeRF models enabling applications such as telegathering.

2 Related Work

As our work is mainly based on neural radiance fields, we first review works on NeRF with a focus on works that aim to control and condition the radiance fields—a necessity for rendering a human in the scene in the context of creating visual and immersive experiences [33,19]. We also briefly review works that aim to reanimate and perform novel view synthesis of provided scenes.

Neural Radiance Fields (NeRF). Since its first introduction [27], NeRF has become a popular way [5] to model scenes and render them from novel views thanks to its high quality rendering. Representing a scene as a radiance field has the advantage that *by-construction* you will be able to render the scene from any supported views through volume rendering. Efforts have been made to adapt NeRF to dynamic scenes [29,30,15], and to even edit and compose scenes with various NeRF models [50,9], widening their potential application. While these methods have shown interesting and exciting results, they often require separate training of editable instances [9] or careful curation of training data [33]. In this work, we are interested in an *in-the-wild* setup.

Particularly related to our task of interest, various efforts have been made towards NeRF models conditioned by explicit human models, such as SMPL [23] or 3D skeleton [19,32,31,12,40]. Neural Body [32] associates a latent code to each SMPL vertex, and use sparse convolution to diffuse the latent code into the volume in observation space. Neural Actor [19] learns the human in the canonical space by a volume warping based on the SMPL [23] mesh transformation, it also utilize a texture map to improve the final rendering quality. Animatable NeRF [31] learns a blending weight field in both observation space and canonical

space, and optimize for a new blending weight field for novel poses. ST-NeRF [12] separates the human into each 3D bounding box, and learns the dynamic human within each bounding box. It doesn’t require to estimate the precise human geometry, but it cannot extrapolate to unseen poses since it is dependent on time(frame). However, all these methods require an expensive multi-cameras setup to obtain the ground truth bounding boxes, 3D poses or SMPL estimates. In other words, they cannot be used for our purpose of reconstructing and neural rendering of human from a *single* video, *without* extra devices or annotations, and *with* potential pose estimation errors.

HumanNeRF [47], a concurrent work, aims to create free-viewpoint rendering of human from a single video. While similar to our work, there are two main differences between HumanNeRF [47] and ours. First, HumanNeRF [47] relies on manual mask annotation to separate the human from the background, while our method learns the decomposition of the human and the scene with the help of modern detectors. Second, HumanNeRF [47] represents motion as a combination of the skeletal and the non-rigid transformations, causing ambiguous or unknown transformations under novel poses, while ours mitigates the ambiguity by using explicit human mesh. Another similar work is Vid2Actor [46], which builds animatable human from a single video by learning a voxelized canonical volume and skinning weights jointly. Although with similar goals, our method is able to reconstruct sharp human geometry with less than 40 images, comparing to thousands frames are required for Vid2Actor, our method is data-efficient.

Neural Rendering of Humans. Majority of the literature [1,24,28,22,36] only consider the problem of reposing a human from a source image to a target image without changing the viewing angle which is essential for enabling new immersive experiences. Grigorev et al [7] tackled a similar problem to ours, namely resynthesizing a human image with a novel pose view given a single input image. They divide the problem into estimating the full texture map of the human body surface from partial texture observations and synthesizing a novel view given the estimated texture map from the first step. They employ two convolutional neural networks (CNN) for each of these steps. With the second one responsible for generating the novel pose consuming the output of the first CNN. This method does not explicitly model the source and target pose, so it is unclear how well it would perform for a novel pose. Sarkar et al [37] additionally consider re-rendering in a novel view as well as novel pose from a single input image. They use a parametric 3D human mesh to recover body pose and shape and a high dimensional UV feature map to encode appearance.

2.1 Preliminary: Neural Radiance Fields (NeRF)

For completeness, we quickly review the standard NeRF model [27]. We denote the NeRF network as \mathcal{F}_{Θ} with parameters Θ , that estimates the RGB color \mathbf{c} and density σ of a given a 3D location \mathbf{x} and a viewing direction \mathbf{d} as

$$\mathbf{c}, \sigma = \mathcal{F}_{\Theta}(\mathbf{x}, \mathbf{d}) . \quad (1)$$

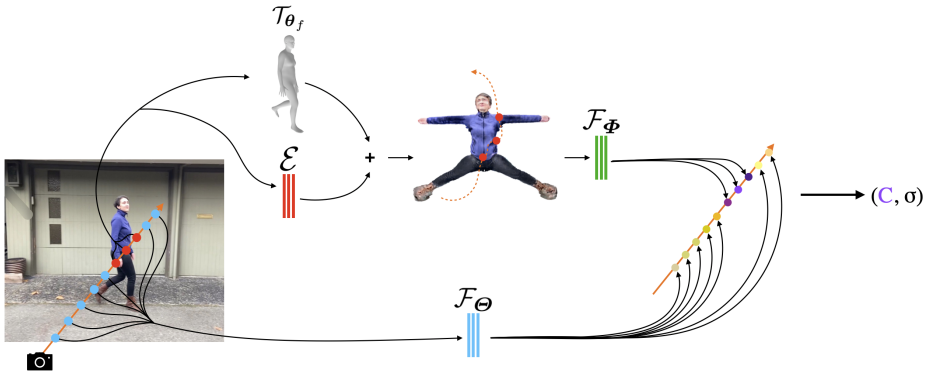


Fig. 2: **Overview** – The blue samples are for the scene branch, and the red samples for the human branch. The human samples are warped to canonical space based on the estimated SMPL mesh and the error-correction network. After the RGB and opacity are evaluated, the two sets of samples are merged for the final integral to obtain the pixel color. See 3.2 for more details.

The radiance field function \mathcal{F}_Θ is often implemented with multi-layer perceptrons (MLPs) with positional encodings [43,27] and periodic activations [2]. The pixel color is then obtained by integrating a discretized ray \mathbf{r} that consists of N samples from the camera in the view direction \mathbf{d} through the volume given by

$$\mathbf{C}(\mathbf{r}) = \sum_{i=1}^N w_i \mathbf{c}_i; \quad \text{where, } w_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right) (1 - \exp(-\sigma_i \delta_i)), \quad (2)$$

and δ_i is the distance between two adjacent samples. Here, the accumulated alpha value of a pixel, which represents transparency, can be obtained by $\alpha(\mathbf{r}) = \sum_{i=1}^N w_i$.

3 Method

An overview of our framework is shown in Figure 2. Our framework is composed mainly of two NeRF networks³: the *human* NeRF that encodes the appearance and geometry of the human in the scene, conditioned on the human pose; and the *scene* NeRF that encodes how the background looks like. We train the scene NeRF first, then train the human NeRF conditioned on the trained scene NeRF.

3.1 The Scene NeRF Model

The scene NeRF model is analogous to the background model in traditional motion detection work [41,6,16], except it’s a NeRF. For the scene NeRF model,

³ In our work, we assume a single human being in the scene, but this can be trivially extended.

we construct a NeRF model and train it with only the pixels that are deemed to be from the background.

For a ray \mathbf{r} , given the human segmentation mask as $\mathcal{M}(\mathbf{r})$ where $\mathcal{M}(\mathbf{r}) = 1$ if the ray corresponds to the human and $\mathcal{M}(\mathbf{r}) = 0$ corresponding the background, we formulate the reconstruction loss for the scene NeRF model as

$$\mathcal{L}_{s,rgb}(\mathbf{r}) = (1 - \mathcal{M}(\mathbf{r})) \left\| \mathbf{C}_s(\mathbf{r}) - \hat{\mathbf{C}}(\mathbf{r}) \right\|, \quad (3)$$

where $\hat{\mathbf{C}}(\mathbf{r})$ corresponds to the ground-truth RGB color value and $\mathbf{C}_s(\mathbf{r})$ corresponds to rendered color value from the scene NeRF model.

As in Video-NeRF [49], simply minimizing Eq. 3 leads to ‘hazy’ objects floating in the scene. Therefore, following Video-NeRF [49], we resolve this by adding a regularizer on the estimated density, and forcing it to be zero for space that should be empty—the space between the camera and the scene. For each ray \mathbf{r} , we sample the terminating depth value $\hat{z}_{\mathbf{r}} = \mathbf{D}_{fuse}(\mathbf{r})$ and minimize

$$\mathcal{L}_{s,empty}(\mathbf{r}) = \int_{t_n}^{\alpha \hat{z}_{\mathbf{r}}} \sigma_s(\mathbf{r}(t)) dt, \quad (4)$$

where $\alpha = 0.8$ is a slack margin to avoid strong regularization when the depth estimates are inaccurate. The final loss that we use to train our scene NeRF is

$$\mathcal{L}_s = \mathcal{L}_{s,rgb}(\mathbf{r}) + \lambda_{empty} \mathcal{L}_{s,empty}(\mathbf{r}), \quad (5)$$

where $\lambda_{empty} = 0.1$ is a hyper parameter controlling the emptiness regularizer in all our experiments.

Preprocessing. Given a video sequence, we use COLMAP [39,38] to obtain the camera poses, sparse scene model, and multi-view-stereo (MVS) depth maps. Typically, MVS depth maps \mathbf{D}_{mvs} contain holes, which we fill with the help of dense monocular depth maps \mathbf{D}_{mono} using Miangoleh et al. [25]. We fuse \mathbf{D}_{mvs} and \mathbf{D}_{mono} together to obtain a fused depth map \mathbf{D}_{fuse} with consistent scale. In more detail, we find a linear mapping between the two depth maps using the pixels that have both estimates. We then transform the values of \mathbf{D}_{mono} with this mapping to match the depth scale in \mathbf{D}_{mvs} to obtain a fused depth map \mathbf{D}_{fuse} by filling in the holes. For retrieving human segmentation maps we apply Mask-RCNN [10]. We further dilate the human masks by 4% to ensure the human is completely masked out. With the estimated camera poses and the background masks, we train the scene NeRF model only over the background.

3.2 The Human NeRF Model

To build a human model that can be pose-driven, we require the model to be pose independent. Therefore, we define a canonical space based on the 6-pose (Da-pose) SMPL [23] mesh, similar to [31,19]. In comparison to the traditional T-pose, Da-pose avoids volume collision when warping from observation space to canonical space for the legs.

To render a pixel of a human in the observation space with this model, we transform the points along that ray into the canonical space. The difficulty in doing so is how one expands the transformation of SMPL meshes into the entire observation space to allow this ray tracing in canonical space. Similar to Liu et al. [19], we use a simple strategy to extend the mesh skinning into a volume warping field.

In each frame f , given a 3D point $\mathbf{x}_f = \mathbf{r}_f(t)$ in observation space and the corresponding estimated SMPL mesh $\boldsymbol{\theta}_f$ obtained from preprocessing 3.2, we transform it into the canonical space by following the rigid transformation of its closest point on the mesh; see Figure 2. We denote this mesh-based transform as \mathcal{T} such that $\mathbf{x}'_f = \mathcal{T}_{\boldsymbol{\theta}_f}(\mathbf{x}_f)$. This transformation, however, relies completely on the accuracy of $\boldsymbol{\theta}_f$, which is not reliable even with the recent state of the art. To mitigate the misalignment between the SMPL estimates and the underlying human, we propose to jointly optimize $\boldsymbol{\theta}_f$ together with the neural radiance field while training. In Figure 3, we show the effect of online optimization on correcting the estimates. Furthermore, to account for the details that can not be expressed by the SMPL model, we introduce the error-correction network \mathcal{E} , an MLP that corrects for the errors in the warping field. Finally, the mapping between the points in the observation space to the corrected points in the canonical space $\mathbf{x}_f \rightarrow \tilde{\mathbf{x}}'_f$ is obtained as

$$\tilde{\mathbf{x}}'_f = \mathcal{T}_{\boldsymbol{\theta}_f}(\mathbf{x}_f) + \mathcal{E}(\mathbf{x}_f, f) . \quad (6)$$

The error-correction net is only used during training, and is discarded for rendering with validation and novel poses. Since a single canonical space is used to explain all poses, the error-correction network naturally overfits to each frame and makes the canonical volume more generalized.

Due to the nature of the warping field, a straight line in the observation space is curved in the canonical space after warping. Therefore, we recompute the viewing angles by taking into account how the light rays *actually* travel in the canonical space by looking at where the previous sample is,

$$\mathbf{d}(t_i)'_f = \hat{\mathbf{x}}'_f(t_i) - \hat{\mathbf{x}}'_f(t_{i-1}), \quad (7)$$

where $\hat{\mathbf{x}}'_f(t_i)$ and $\mathbf{d}(t_i)'_f$ are the coordinate and viewing angel of the i -th sample along the curved ray in canonical space. Finally, with the canonical space coordinates $\tilde{\mathbf{x}}'_f$ and the corrected viewing angles \mathbf{d}'_f the radiance field values for the human model is obtained by

$$\mathbf{c}_h, \sigma_h = \mathcal{F}_{\Phi}(\tilde{\mathbf{x}}'_f, \mathbf{d}'_f) , \quad (8)$$

where Φ are the parameters of the human NeRF \mathcal{F}_{Φ} .

To render a pixel, we shoot two rays, one for the human NeRF, and the other for the scene NeRF. We evaluate the colors and densities for the two sets of samples along the rays. We then sort the colors and the densities in the ascending order based on their depth values, similar to ST-NeRF [12]. Finally, we integrate over these values to obtain the pixel using Eq. (2).

Training. To train the human radiance field, we sample rays on the regions covered by the human mask and minimize

$$\mathcal{L}_{h,rgb}(\mathbf{r}) = \mathcal{M}(\mathbf{r}) \left\| \mathbf{C}_h(\mathbf{r}) - \hat{\mathbf{C}}(\mathbf{r}) \right\|, \quad (9)$$

where $\mathbf{C}_h(\mathbf{r})$ is the rendered color from the human NeRF model. Similar to HumanNeRF [47], we also use LPIPS [51] as an additional loss term \mathcal{L}_{lpiips} by sampling a 32×32 patch. We use \mathcal{L}_{mask} to enforce the accumulated alpha map from the human NeRF to be similar to the detected human mask.

$$\mathcal{L}_{mask}(\mathbf{r}) = \mathcal{M}(\mathbf{r}) \|1 - \alpha_h(\mathbf{r})\|, \quad (10)$$

where α_h corresponds to accumulated density over the ray as defined in Sec. 2.1.

To avoid blobs in the canonical space and semi-transparent canonical human, we enforce the volume inside the canonical SMPL mesh to be solid, while enforcing the volume outside the canonical SMPL mesh to be empty, given by

$$\mathcal{L}_{smpl}(\hat{\mathbf{x}}'_f, \sigma_h) = \begin{cases} \|1 - \sigma_h\|, & \text{if } \hat{\mathbf{x}}'_f \text{ inside SMPL mesh} \\ |\sigma_h|, & \text{otherwise} \end{cases}, \quad (11)$$

Moreover, we utilize hard surface loss \mathcal{L}_{hard} [35] to mitigate the halo around the canonical human. To be specific, we encourage the weight of each sample to be either 1 or 0 given by,

$$\mathcal{L}_{hard} = -\log(e^{-|w|} + e^{-|1-w|}) \quad (12)$$

where w refers to the transparency where the ray terminates as defined in Sec. 2.1. However, this penalty alone is not enough to obtain a sharp canonical shape, we also add a canonical edge loss, \mathcal{L}_{edge} . By rendering a random straight ray in the canonical volume, we encourage the accumulated alpha values to be either 1 or 0. This is given by,

$$\mathcal{L}_{edge} = -\log(e^{-|\alpha_c|} + e^{-|1-\alpha_c|}) \quad (13)$$

where α_c is the accumulated alpha value obtained from a random straight ray in canonical space. Thus, the final loss is given by,

$$\begin{aligned} \mathcal{L} = & \mathcal{L}_{h,rgb} + \lambda_{lpiips} \mathcal{L}_{lpiips} + \lambda_{mask} \mathcal{L}_{mask} \\ & + \lambda_{smpl} \mathcal{L}_{smpl} + \lambda_{hard} \mathcal{L}_{hard} + \lambda_{edge} \mathcal{L}_{edge}. \end{aligned} \quad (14)$$

To train, we jointly optimize θ_f , \mathcal{E} , and \mathcal{F}_{Φ} by minimizing this loss. We set $\lambda_{lpiips} = 0.01$, $\lambda_{mask} = 0.01$, $\lambda_{smpl} = 1.0$, $\lambda_{hard} = 0.1$, and $\lambda_{edge} = 0.1$. Since the detected masks are inaccurate, we linearly decay λ_{mask} to 0 through the training.

Preprocessing. We utilize ROMP [42] to estimate the SMPL [23] parameters of the human in the videos. However, the estimated SMPL parameters are not accurate. Therefore, we refine the SMPL estimates by optimizing the

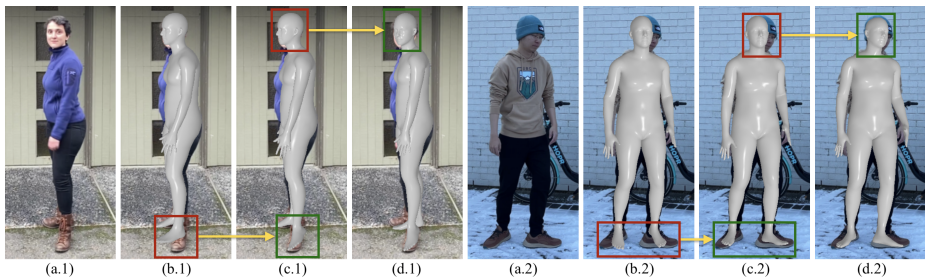


Fig. 3: **Pose optimization examples**– Our proposed preprocessing and end-to-end optimization over SMPL parameters effectively produce better fits to the human, see the bounding boxes. (a) Original image. (b) ROMP estimates. (c) Preprocessed SMPL. (d) End-to-end optimized SMPL.

SMPL parameters using silhouette estimated from [8,48], and 2D joints estimated from [3,4] as detailed in the supplementary material. We then align the SMPL estimates in the scene coordinates.



Fig. 4: **Visualization of SMPL and scene alignment**– We show sampled video frames(first row), and the estimated SMPL meshes overlaying on top of the scene point cloud(second row). The human is in the scene with a correct scale as the foot is touching the ground plane. SMPL meshes are colored based on time.

Scene-SMPL Alignment. To compose a scene with a human in novel view and pose, and to train the two NeRF models, we align the coordinate systems in which the two NeRF models lie. This is, in fact, a non-trivial problem, as human body pose estimators [13,42,18] operate in their own camera systems with often near-orthographic camera models. To deal with this issue, we first solve the Perspective-n-Point (PnP) problem [14] between the estimated 3D joints and the projected 2D joints with the camera intrinsics from COLMAP. This solves the alignment up to an arbitrary scale. We then assume that the human is standing on a ground at least in one frame, and solve for the scale ambiguity by finding the scale that allows the feet meshes of the SMPL model to touch the ground plane. We obtain the ground plane by applying RANSAC. We show the results of the aligned SMPL estimates in the scene in Figure 4.

Once the two NeRF models are properly aligned we can render the pixel by shooting two rays, one for the human NeRF model, and the other for the scene NeRF model, as describe above. For the near and far planes to generate samples in Eq. 2, we use the estimated scene point cloud to determine them for scene NeRF, and use the estimated SMPL mesh to determine them for the human NeRF, following the strategy in Liu et al. [19].

4 Experiments

We introduce our dataset, show qualitative and quantitative results of our method and provide ablation studies. Our method is the first method that can render a human together with a scene with novel human poses and novel views from a single video. We show the importance of our geometry correction and novel loss terms in order to obtain a realistic and sharp human NeRF model.

4.1 Dataset

Existing human motion datasets are not suitable for our experiments. Generally, motion capture data [32,11] is captured with a static multi-cameras system in a controlled environment which defeats the purpose of reconstructing from a single video. Other video datasets [26,34] have multiple humans and crowded scenes that are not suitable for our case. Therefore, we introduce NeuMan dataset, a collection of 6 videos about 10 to 20 seconds long each, where a single person performs a walking sequence captured using a mobile phone. Moreover, the camera reasonably pans through the scene to enable multi-view reconstruction. The sequences are named – Seattle, Citron, Parking, Bike, Jogging and Lab.



Fig. 5: **Scene NeRF examples** – We show the training samples(left), together with the renderings from validation views. Our scene models are able to remove the human in the scene effectively, and to produce high quality novel view renderings of the background even with limited coverage of the scene.

For each video sequence, we first subsample the frames from the video and use them to train our NeRF models. We split frames into 80% training frames, 10% validation frames, and 10% test frames. We provide the details in supplementary material.

4.2 Qualitative Results

Scene NeRF Reconstructions. Figure 5 shows the novel view renderings of our scene NeRF models. By reconstructing the background pixels only, our model learns the consistent geometry of the scene, and effectively removes the dynamic human.

Human NeRF Reconstructions. Our framework learns an animatable human model with realistic details. It captures not only the texture details such as the pattern on the cloth, but also the subject specific geometric details, such as the sleeves, collar, even zipper. Notice that these geometric details are beyond the expressiveness of SMPL model. The learned human models can be reposed to novel driving motions, and produce high quality rendering of the human under novel poses from novel views. Although, the training sequence is as simple as walking, the model can perform stunning cartwheeling motion, which shows the



Fig. 6: **Human NeRF model examples** – We show front and side view of the reconstructed canonical human(left), novel human poses and views renderings with the reconstructed human model(middle), and composition of both the human and scene model with novel human poses(right).

ability to extrapolate to unseen poses. By simple composition, we can render both the human and the scene realistically; see Figure 6.

Telegathering. The ability to render reposed human together with the scene further allows us to, for example, create telegathering of multiple individuals as in Figure 7. The results show that our framework can facilitate combining human NeRF models in the same scene without any additional training.



Fig. 7: **Telegathering**— Our method is able to provide telegathering for multiple individuals: (left) is the SMPL meshes and scene point cloud overlay, others are the rendering results from novel views.

4.3 Novel view synthesis

To compare our method to existing solutions that can do similar tasks, we train NeRF with time(NeRF-T) [15] using the same training data as ours but with-

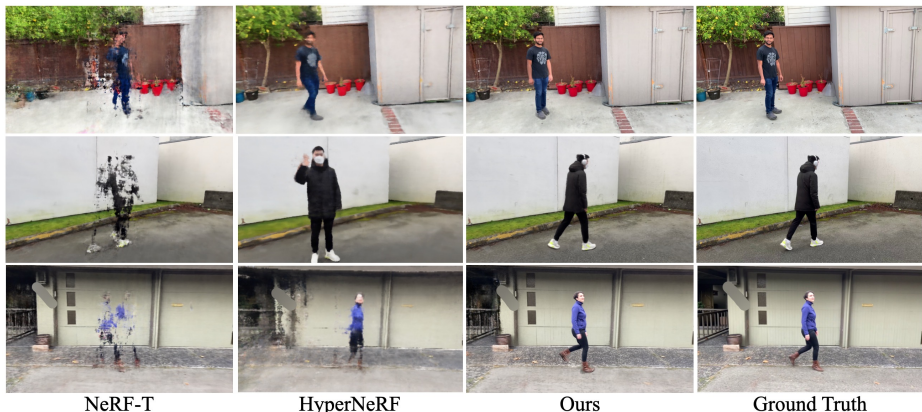


Fig. 8: **Qualitative comparisons**— We show the qualitative comparisons among NeRF-T, HyperNeRF and ours. NeRF-T and HyperNeRF failed to reconstruct the drastically dynamic scenes and to interpolate across space and time, while our method faithfully reconstructed both the human and the scene with high rendering quality.

	Seattle			Citron			Parking		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NeRF-T	21.84	0.69	0.37	12.33	0.49	0.65	21.98	0.69	0.46
HyperNeRF	16.43	0.43	0.40	16.81	0.41	0.56	16.04	0.38	0.62
Ours	24.06	0.77	0.27	24.72	0.79	0.27	25.76	0.79	0.32

	Bike			Jogging			Lab		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NeRF-T	21.16	0.71	0.36	20.63	0.53	0.49	20.52	0.75	0.39
HyperNeRF	17.64	0.42	0.43	18.52	0.39	0.52	16.75	0.51	0.23
Ours	25.34	0.82	0.24	22.69	0.67	0.33	24.78	0.85	0.21

Table 2: **Novel view synthesis comparison** – We report the quantitative results on test views. Our method achieves the best rendering quality across all scenes and all metrics by a large margin. Notice that the numbers for our method depends on the estimated human pose.

out the empty space penalty. As HyperNeRF [30] requires a smoothly changing video, we provide it with densely sampled frames from the videos. NeRF-T and HyperNeRF [30] do not perform well at reconstruction of drastically dynamic scenes with humans, while our method is able to do so even with less than 40 training images. We evaluate methods on the test views, and measure PSNR, SSIM [44], and LPIPS [51]. The results are shown in Table 2 and Figure 8. Our method outperforms across all scenes and metrics. Notice that, unlike HyperNeRF, our method depends on the estimated human pose, and we discard the offset network when rendering validation or test views. Therefore, our reported numbers also reflect the errors in the estimated human pose.

4.4 Ablation Studies

Effect of Geometry Correction Our method has 3 components to correct the estimated geometry of the human: the offline SMPL optimization in the preprocessing, the online end-to-end SMPL optimization during the training, and the error-correction network which accounts for the warping errors. We train models(Ours-GC) without the end-to-end SMPL optimization and error-correction network, and starting with raw SMPL estimates from ROMP [42]. Without any offline or online geometry corrections, the canonical volume overfits to each observations causing averaged color and shape instead of sharp details of the human. However, with geometry correction enabled, the human NeRF model can learn both the textural and geometric details of the human. We show the comparison between Ours-Full and Ours-GC in Figure 9.



Fig. 9: **Ablation studies** – We show the canonical(first row) and test view(second row) renderings of our full model, ours without any geometry corrections, and ours without \mathcal{L}_{smpl} . Geometry correction allows us to reconstruct the human with details, and \mathcal{L}_{smpl} encourages a clean canonical volume and suppresses halo in the final renderings.

Effect of \mathcal{L}_{smpl} We train models(Ours- \mathcal{L}_{smpl}) with $\lambda_{smpl} = 0$. When \mathcal{L}_{smpl} is disabled, fog appears in the canonical volume around the human, and causes a halo in the final renderings, see Ours- \mathcal{L}_{smpl} in Figure 9. By encouraging the volume outside the canonical SMPL mesh to be empty, \mathcal{L}_{smpl} effectively removes the unwanted fog in the volume and suppresses the halo in the final renderings.

5 Conclusions

We have proposed a novel framework to reconstruct the human and the scene NeRF models that can be rendered with novel human poses and views from a single in-the-wild video without any extra devices or manual annotations. To do so, we use off-the-shelf methods to extract camera poses and depth of each frame, the scene point cloud, the human body pose, and the human mask. We then utilize these to build two NeRF models, one for the human and the one for the scene. Our scene NeRF model reconstructs only the background, and can faithfully render novel views of the background scene. Our human NeRF model is able to learn texture details such as patterns on cloth, and geometric details such as sleeves, collar even zipper from less than 40 images. By simple composition, we can further create renderings of individuals in the same scene or new different scenes, without any additional training.

Limitations and future work. The most obvious limitation of our method, from the renderings, is that the dynamics beyond SMPL cannot be modeled with our static neural radiance field, those dynamics will degenerate to average shape or color. This is most evident in the hands of the person. When hand gestures are changing over the video, the volume of hands will converge to an average shape. We hope to extend our method with more sophisticated body models, or even ones that can learn them, for a more detailed reconstruction and more expressive rendering in these parts of the human body.

Additionally, our warping function is a simple extension of the SMPL mesh skinning, it could cause volume collision in some extreme cases. A collision-aware volume warping method or a learned one is required to improve the generalization under extreme poses.

Finally, we assume that the human always has at least one contact point with the ground to estimate the scale relative to the scene. Therefore, our method cannot be used for videos where the human has zero contact point with the ground, for example, when the human is doing a backflip, or for videos where the ground is not a simple plane, for example a video of human stepping down stairs. Solving this would require smarter geometric reasoning, which could be a promising direction for future research.

6 Acknowledgement

We thank Ashish Shrivastava, Russ Webb and Miguel Angel Bautista Martin for providing insightful review feedback.

References

1. Balakrishnan, G., Zhao, A., Dalca, A.V., Durand, F., Gutttag, J.: Synthesizing images of humans in unseen poses. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 8340–8348 (2018) [4](#)
2. Chan, E.R., Monteiro, M., Kellnhofer, P., Wu, J., Wetzstein, G.: pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5799–5809 (2021) [5](#)
3. Cheng, B., Xiao, B., Wang, J., Shi, H., Huang, T.S., Zhang, L.: HigherHRNet: Scale-Aware Representation Learning for Bottom-Up Human Pose Estimation. In: CVPR (2020) [2](#), [9](#), [19](#)
4. Contributors, M.: OpenMMLab Pose Estimation Toolbox and Benchmark. <https://github.com/open-mmlab/mmpose> (2020) [9](#)
5. Dellaert, F., Yen-Chen, L.: Neural Volume Rendering: NeRF And Beyond (2021) [3](#)
6. Elgammal, A., Harwood, D., Davis, L.: Non-parametric model for background subtraction. In: European conference on computer vision. pp. 751–767. Springer (2000) [5](#)
7. Grigorev, A., Sevastopolsky, A., Vakhitov, A., Lempitsky, V.: Coordinate-based texture inpainting for pose-guided human image generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12135–12144 (2019) [4](#)
8. Güler, R.A., Neverova, N., Kokkinos, I.: Densepose: Dense human pose estimation in the wild. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7297–7306 (2018) [2](#), [9](#), [19](#)
9. Guo, M., Fathi, A., Wu, J., Funkhouser, T.: Object-Centric Neural Scene Rendering. <https://arxiv.org/abs/2012.08503> (2020) [3](#)
10. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask R-CNN. In: Proceedings of the IEEE international conference on computer vision. pp. 2961–2969 (2017) [2](#), [6](#), [19](#)

11. Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C.: Human3.6M: Large Scale Datasets and Predictive Methods for 3D Human Sensing in Natural Environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **36**(7), 1325–1339 (jul 2014) [10](#)
12. Jiakai, Z., Xinhang, L., Xinyi, Y., Fuqiang, Z., Yanshun, Z., Minye, W., Yingliang, Z., Lan, X., Jingyi, Y.: Editable free-viewpoint video using a layered neural representation. In: *ACM SIGGRAPH (2021)* [1](#), [2](#), [3](#), [4](#), [7](#)
13. Kocabas, M., Athanasiou, N., Black, M.J.: Vibe: Video inference for human body pose and shape estimation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 5253–5263 (2020) [9](#)
14. Lepetit, V., Moreno-Noguer, F., Fua, P.: Epnnp: An accurate $o(n)$ solution to the pnp problem. *International Journal Of Computer Vision* **81**, 155–166 (2009). <https://doi.org/10.1007/s11263-008-0152-6>, <http://infoscience.epfl.ch/record/160138> [9](#)
15. Li, Z., Niklaus, S., Snavely, N., Wang, O.: Neural scene flow fields for space-time view synthesis of dynamic scenes. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 6498–6508 (2021) [1](#), [3](#), [12](#)
16. Lim, L.A., Keles, H.Y.: Foreground segmentation using convolutional neural networks for multiscale feature encoding. *Pattern Recognition Letters* **112**, 256 – 262 (2018). <https://doi.org/https://doi.org/10.1016/j.patrec.2018.08.002>, <http://www.sciencedirect.com/science/article/pii/S0167865518303702> [5](#)
17. Lin, C.H., Ma, W.C., Torralba, A., Lucey, S.: BARF: Bundle-Adjusting Neural Radiance Fields. In: *ICCV (2021)* [1](#)
18. Lin, K., Wang, L., Liu, Z.: End-to-end human pose and mesh reconstruction with transformers. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 1954–1963 (2021) [9](#)
19. Liu, L., Habermann, M., Rudnev, V., Sarkar, K., Gu, J., Theobalt, C.: Neural Actor: Neural Free-view Synthesis of Human Actors with Pose Control. *ACM Trans. Graph.(ACM SIGGRAPH Asia)* (2021) [1](#), [2](#), [3](#), [6](#), [7](#), [10](#)
20. Liu, L., Habermann, M., Rudnev, V., Sarkar, K., Gu, J., Theobalt, C.: Neural actor: Neural free-view synthesis of human actors with pose control. *arXiv preprint arXiv:2106.02019* (2021) [1](#), [3](#)
21. Liu, S., Li, T., Chen, W., Li, H.: Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 7708–7717 (2019) [19](#)
22. Liu, W., Piao, Z., Min, J., Luo, W., Ma, L., Gao, S.: Liquid warping gan: A unified framework for human motion imitation, appearance transfer and novel view synthesis. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 5904–5913 (2019) [4](#)
23. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: Smpl: A skinned multi-person linear model. *ACM transactions on graphics (TOG)* **34**(6), 1–16 (2015) [3](#), [6](#), [8](#)
24. Ma, L., Sun, Q., Georgoulis, S., Van Gool, L., Schiele, B., Fritz, M.: Disentangled person image generation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 99–108 (2018) [4](#)
25. Miangoleh, S.M.H., Dille, S., Mai, L., Paris, S., Aksoy, Y.: Boosting Monocular Depth Estimation Models to High-Resolution via Content-Adaptive Multi-Resolution Merging. In: *CVPR (2021)* [2](#), [6](#)
26. Milan, A., Leal-Taixé, L., Reid, I., Roth, S., Schindler, K.: Mot16: A benchmark for multi-object tracking. *arXiv preprint arXiv:1603.00831* (2016) [10](#)

27. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: European conference on computer vision. pp. 405–421. Springer (2020) **1, 3, 4, 5**
28. Neverova, N., Guler, R.A., Kokkinos, I.: Dense pose transfer. In: Proceedings of the European conference on computer vision (ECCV). pp. 123–138 (2018) **4**
29. Park, K., Sinha, U., Barron, J.T., Bouaziz, S., Goldman, D.B., Seitz, S.M., Martin-Brualla, R.: Nerfies: Deformable Neural Radiance Fields. ICCV (2021) **1, 3**
30. Park, K., Sinha, U., Hedman, P., Barron, J.T., Bouaziz, S., Goldman, D.B., Martin-Brualla, R., Seitz, S.M.: HyperNeRF: A Higher-Dimensional Representation for Topologically Varying Neural Radiance Fields. arXiv preprint arXiv:2106.13228 (2021) **1, 2, 3, 13**
31. Peng, S., Dong, J., Wang, Q., Zhang, S., Shuai, Q., Zhou, X., Bao, H.: Animatable Neural Radiance Fields for Modeling Dynamic Human Bodies. In: ICCV (2021) **1, 3, 6**
32. Peng, S., Zhang, Y., Xu, Y., Wang, Q., Shuai, Q., Bao, H., Zhou, X.: Neural Body: Implicit Neural Representations with Structured Latent Codes for Novel View Synthesis of Dynamic Humans. In: CVPR (2021) **1, 3, 10**
33. Peng, S., Zhang, Y., Xu, Y., Wang, Q., Shuai, Q., Bao, H., Zhou, X.: Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9054–9063 (2021) **1, 3**
34. Ranjan, A., Hoffmann, D.T., Tzionas, D., Tang, S., Romero, J., Black, M.J.: Learning multi-human optical flow. International Journal of Computer Vision **128**(4), 873–890 (2020) **10**
35. Rebain, D., Matthews, M., Yi, K.M., Lagun, D., Tagliasacchi, A.: LOLNeRF: Learn from One Look. arXiv preprint arXiv:2111.09996 (2022) **8**
36. Sanyal, S., Vorobiov, A., Bolkart, T., Loper, M., Mohler, B., Davis, L.S., Romero, J., Black, M.J.: Learning realistic human reposing using cyclic self-supervision with 3d shape, pose, and appearance consistency. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 11138–11147 (2021) **4**
37. Sarkar, K., Mehta, D., Xu, W., Golyanik, V., Theobalt, C.: Neural re-rendering of humans from a single image. In: European Conference on Computer Vision. pp. 596–613. Springer (2020) **4**
38. Schönberger, J.L., Frahm, J.M.: Structure-from-Motion Revisited. In: CVPR (2016) **2, 6**
39. Schönberger, J.L., Zheng, E., Pollefeys, M., Frahm, J.M.: Pixelwise view selection for unstructured multi-view stereo. In: European Conference on Computer Vision (ECCV) (2016) **6**
40. Su, S.Y., Yu, F., Zollhoefer, M., Rhodin, H.: A-NeRF: Surface-free Human 3D Pose Refinement via Neural Rendering. <https://arxiv.org/abs/2102.06199> (2021) **1, 3**
41. Sun, D., Sudderth, E.B., Black, M.J.: Layered segmentation and optical flow estimation over time. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition. pp. 1768–1775. IEEE (2012) **5**
42. Sun, Y., Bao, Q., Liu, W., Fu, Y., Michael J., B., Mei, T.: Monocular, One-stage, Regression of Multiple 3D People. In: ICCV (October 2021) **2, 3, 8, 9, 13, 19**
43. Tancik, M., Srinivasan, P.P., Mildenhall, B., Fridovich-Keil, S., Raghavan, N., Singhal, U., Ramamoorthi, R., Barron, J.T., Ng, R.: Fourier features let networks learn high frequency functions in low dimensional domains. arXiv preprint arXiv:2006.10739 (2020) **5**
44. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image Quality Assessment: From Error Visibility to Structural Similarity. IEEE TIP (2004) **13**

45. Wang, Z., Wu, S., Xie, W., Chen, M., Prisacariu, V.A.: NeRF—: Neural Radiance Fields Without Known Camera Parameters. arXiv preprint arXiv:2102.07064 (2021) **1**
46. Weng, C.Y., Curless, B., Kemelmacher-Shlizerman, I.: Vid2Actor: Free-viewpoint Animatable Person Synthesis from Video in the Wild. arXiv preprint arXiv:2012.12884 (2020) **2, 4**
47. Weng, C.Y., Curless, B., Srinivasan, P.P., Barron, J.T., Kemelmacher-Shlizerman, I.: HumanNeRF: Free-viewpoint Rendering of Moving People from Monocular Video. arXiv preprint arXiv:2201.04127 (2022) **2, 4, 8**
48. Wu, Y., Kirillov, A., Massa, F., Lo, W.Y., Girshick, R.: Detectron2. <https://github.com/facebookresearch/detectron2> (2019) **9**
49. Xian, W., Huang, J.B., Kopf, J., Kim, C.: Space-time Neural Irradiance Fields for Free-Viewpoint Video. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 9421–9431 (2021) **6**
50. Yang, B., Zhang, Y., Xu, Y., Li, Y., Zhou, H., Bao, H., Zhang, G., Cui, Z.: Learning Object-Compositional Neural Radiance Field for Editable Scene Rendering. In: ICCV (October 2021) **3**
51. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In: CVPR (2018) **8, 13**

A.1 Dataset Details

The dataset details are as follows.

Sequence	Total Frames	Train Frames	Validation Frames	Test Frames
Seattle	41	33	4	4
Citron	37	30	4	3
Parking	42	34	4	4
Bike	104	83	11	10
Jogging	102	82	10	10
Lab	103	82	11	10

Table 3: Number of frames in each dataset used for training, validation and test.

A.2 SMPL Refinement

Given an image, we regress the 2D joints j_{2d} and segmentation mask m of the human using HigherHRNet [3] and Densepose [8]. We further estimate the SMPL mesh $M = (V, F)$, a collection of vertices and faces using ROMP [42]. The mesh M is parametrized by SMPL parameters θ such that $M = \text{SMPL}(\theta)$ and includes the 3D joints j_{3d} . The regressed SMPL parameters θ are noisy. Therefore, we use soft-rasterizer [21], Π to refine these estimates. Given a mesh, M and camera θ_c , the rasterizer renders a silhouette $\hat{m} = \Pi(\theta_c, M)$. We also project the 3D joints in the image plane using camera matrix $\hat{j}_{2d} = \mathbf{p}(j_{3d})$ where \mathbf{p} is a projection operator. We obtain the refined SMPL parameters and camera estimates by minimizing

$$\theta^* = \min_{\theta} \| m - \hat{m} \| + \| j_{2d} - \hat{j}_{2d} \|. \quad (15)$$

Notice that we use the estimates from Densepose [8] as the target silhouettes in the preprocessing, while using the estimates from Mask-RCNN [10] as the target masks during the training of human NeRF model. It’s because in the preprocessing phase, we only have a SMPL model and Densepose [8] is trained with such dense SMPL correspondences. However, we wish to learn extra geometry details beyond the SMPL model with the human NeRF model, and Mask-RCNN [10] better estimates those 2D details such as the head and clothes.