

Dual-Space NeRF: Learning Animatable Avatars and Scene Lighting in Separate Spaces

Yihao Zhi^{1*} Shenhan Qian^{1*} Xinhao Yan^{1*} Shenghua Gao^{1,2,3†}
 {zhiyh, qianshh, yanxh, gaoshh}@shanghaitech.edu.cn

¹ShanghaiTech University

²Shanghai Engineering Research Center of Intelligent Vision and Imaging

³Shanghai Engineering Research Center of Energy Efficient and Custom AI IC

Abstract

Modeling the human body in a canonical space is a common practice for capturing and animation. But when involving the neural radiance field (NeRF), learning a static NeRF in the canonical space is not enough because the lighting of the body changes when the person moves even though the scene lighting is constant. Previous methods alleviate the inconsistency of lighting by learning a per-frame embedding, but this operation does not generalize to unseen poses. Given that the lighting condition is static in the world space while the human body is consistent in the canonical space, we propose a dual-space NeRF that models the scene lighting and the human body with two MLPs in two separate spaces. To bridge these two spaces, previous methods mostly rely on the linear blend skinning (LBS) algorithm. However, the blending weights for LBS of a dynamic neural field are intractable and thus are usually memorized with another MLP, which does not generalize to novel poses. Although it is possible to borrow the blending weights of a parametric mesh such as SMPL, the interpolation operation introduces more artifacts. In this paper, we propose to use the barycentric mapping, which can directly generalize to unseen poses and surprisingly achieves superior results than LBS with neural blending weights. Quantitative and qualitative results on the Human3.6M and the ZJU-MoCap datasets show the effectiveness of our method. Our code is available at: <https://github.com/zyhbili/Dual-Space-NeRF>.

1. Introduction

Human body reconstruction and rendering have long been an active research topic. Multi-view videos are especially suitable for this task because they record not only



Figure 1: Our method learns a human body with a neural radiance field [14] and animates it in arbitrary poses with no need for fine-tuning or additional input. Since we model the subject in the canonical space and the lighting in the world space, we name the method “Dual-Space NeRF”.

the appearance but also the movements and deformation of a person. Classic reconstruction and rendering techniques show limited image realism due to the complexity of decoupling the geometry, material and lighting from images. However, the recently proposed neural radiance field (NeRF) [14] proves it possible to represent a static scene with an MLP without explicitly modeling the above factors. Several recent works [20, 9, 19, 15, 24, 26, 27, 8] have adapted NeRF onto human body reconstruction and animation, but challenges remain in the following aspects:

Recent methods that model the human body with NeRF [14] mostly learn the human body in a canonical space, but the lighting inconsistency in the canonical space lacks exploration [20, 9, 19, 15, 24]. When learning a NeRF in the canonical space, we assume the appearance of a person is consistent across varied body poses. But the fact is when a person moves, a point in the canonical space comes to a different place in the world space, resulting in a change

*Equal contribution.

†Corresponding author.

of the lighting condition. This means that merely modeling the appearance of a person in the canonical space is not enough. Therefore, we propose to model the scene lighting with another MLP in the world space, where the scene lighting is assumed to be static. The lighting MLP takes in a point position, a normal vector, and a view direction in the world space, and outputs a lightness coefficient that adjusts the brightness of the point.

Different from NeRF [14] that only depends on the point position and the viewing direction, our lighting MLP also takes in a normal vector, still due to the complexity of dynamic scenes. For a static scene, the surface normal can be uniquely determined by the point position, while for our setting, the surface normal may change when the subject moves. Unlike IDR [28] that models the lighting condition and the appearance of an object with one appearance MLP, our method learns the two factors with separate MLPs in different spaces to realize correct lighting under unseen poses. Also, our lighting MLP only predicts a scalar coefficient to rescale the color from the body MLP instead of directly predicting a color to prevent overfitting.

To bridge the world space and the canonical space, the key is building pixel-level correspondences across views and frames. A stream of methods bind points on 3D skeletons [15, 24] based on the assumption of local rigidity. Another stream of methods incorporates geometric priors characterized by anchoring points onto SMPL [11], a parametric human body model. Neural Body [20] binds features onto SMPL’s vertices and diffuses them into the space before volumetric rendering. It produces realistic novel views of the training sequence but degrades on novel poses. Animatable NeRF [19] resolves novel-pose synthesis by mapping observed points into a canonical space with inverse linear blend skinning (LBS). Since the LBS weight of a spatial point varies for different poses, Animatable NeRF [19] learns a neural blending weight network conditioned on the pose, which requires additional training for novel poses.

To avoid learning the volatile LBS weights, we seek a pose-independent local position representation that generalizes to novel poses easily. Specifically, we propose a barycentric mapping (BM) as follows. For a point in the space, we first project it onto its closest face on the fitted SMPL mesh. Then we describe this point by the barycentric coordinates of its projected point and its signed height from the face. Finally, its corresponding point in the canonical space is uniquely determined. Note that NPMs [16] leverages a similar barycentric mapping to obtain pseudo ground truths to train a unidirectional deformation field. While in our method, we extend BM to support vector transformation and use it bidirectionally, transforming a point position from the world space to the canonical space and warping a surface normal from the canonical space to the world space. By this barycentric mapping, the body MLP and the lighting

MLP are bridged. BM is parameter-free, enabling pose generalization without additional input or network fine-tuning. Though the parameter-free method seems to have inferior expressiveness, our experiments show its comparable ability on two datasets and clear advantages under challenging poses. It is flexible since it anchors a point on the edge vectors of a face, allowing local deformation along with the face. It also avoids the artifacts of inverse LBS caused by blending weighting interpolation (Fig. 7).

Another challenge of this task is the existence of random variations such as clothing wrinkles. These variations are neither fully determined by the body pose nor consistent across frames, making it harder to learn a stable canonical radiance field. Neural Actor [9] uses the ground-truth texture map to relieve the ambiguity in training images but requires a separate image generation network to infer texture maps from normal maps for testing. Motivated by SCANimate [22], to model the pose-dependent deformation, our model is conditioned on the pose parameters. Besides, per-frame latent embeddings are used to capture the random variations. According to our experiments, the union of these components is sufficient to represent vivid deformations along with pose changes with no need for an extra deformation network.

Our contributions can be summarized as follow:

- To ensure the lighting correctness under unseen poses, we propose dual-space NeRF, which models the static scene lighting in the world space and the human body in the canonical space.
- We propose to use the barycentric mapping to build correspondences between two spaces and validate its comparable expressiveness and superior generalization ability under extreme poses.
- We show the effectiveness and interpretability of our method with quantitative and qualitative results on the Human3.6M [7] and the ZJU-MoCap [20] datasets.

2. Related Work

3D human reconstruction. As a popular 3D articulated human model, SMPL [11] learns a template mesh from 3D scans and deforms the mesh with the linear blend skinning (LBS) algorithm. The mesh-based 3D human model is easy to manipulate but limited to a fixed topology. Therefore, neural implicit functions are adopted to model a 3D avatar that can be animated by SMPL [22, 3, 4, 13]. LEAP [13], SCANimate *et al.* [22], and SNARF [3] learn the human body in a canonical space and articulate the body with neural blending weights. NASA [4] is a part-based method that binds an implicit function on each bone. These methods use 3D data as the input and only account the shape of the human body.

Dynamic neural radiance field. Neural radiance fields

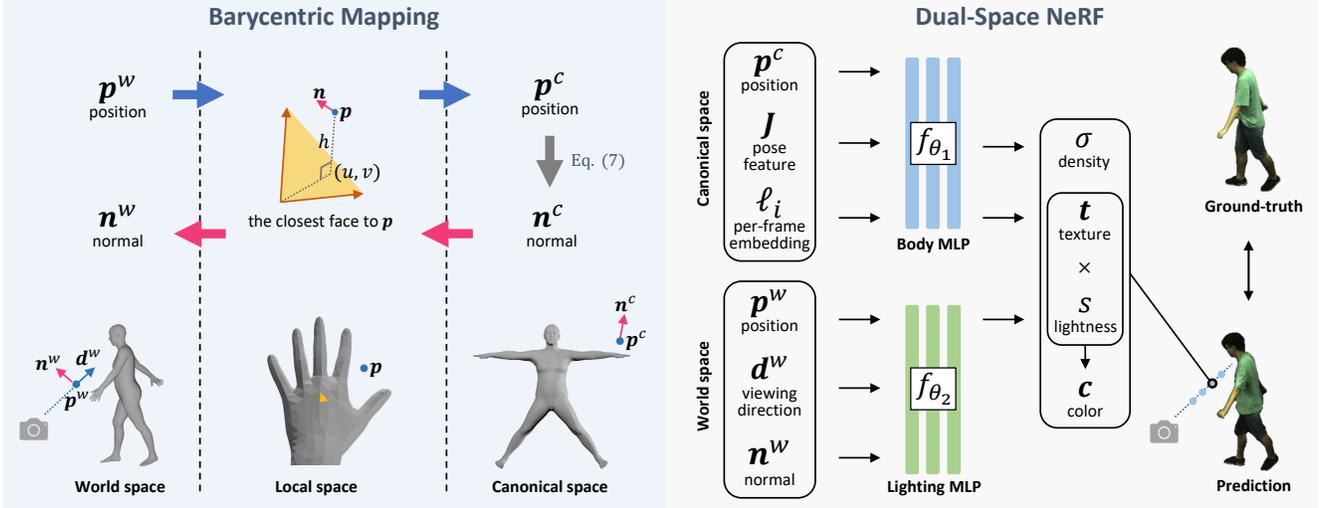


Figure 2: The pipeline of our method. Given a point p^w in the world space, we use the barycentric mapping to warp a point p^w from the world space into the canonical space to query the body properties. In the canonical space, we compute the surface normal n^c and warp it back to obtain the normal vector n^w in the world space. We learn a Body MLP to model a human body in the canonical space and a Lighting MLP to capture the lighting condition in the world space. Finally, we render an image with volumetric rendering.

(NeRF) [14] can synthesize photorealistic images from arbitrary views with no need for 3D data. However, the vanilla NeRF is designed for only static scenes. The challenge for NeRF to model dynamic scenes lies in building correspondences across the timeline. Recent methods define a transformation field that warps an observed point into a canonical space [21, 25, 17, 18, 6]. HumanNeRF [26] records the variation of a scene by learning a motion field from monocular video. Kwon *et al.* [8] resort to a temporal transformer to aggregate skeletal, temporal, and spatial features. These methods exhibit strong ability of replaying events in novel views but are incapable of generating new contents such as animating a human body under unseen poses.

Human body animation with NeRF. To animate a NeRF of a human body, a straightforward solution is incorporating 3D human priors. Peng *et al.* [20] use a set of latent code to encode the local geometry and appearance of the human body and bind them onto SMPL [11] vertices. Liu *et al.* [9] introduce Neural Actor, animating NeRF with blending weights sampled from the nearest vertex of SMPL. Additionally, an image translation network is used to infer texture maps to provide residual deformations and appearance details for novel poses. AniNeRF [19] learns a neural blending weight field to learn the LBS weights for each particular pose. Since the blending weight field varies with poses, AniNeRF relies on a per-frame latent vector as a condition for training poses and requires fine-tuning for novel poses. Xu *et al.* [27] learn NeRF upon imGHUM [1], a statistical human body model represented by neural implicit functions.

3. A Revisit of NeRF

NeRF [14] represents a scene by density σ and color c at each spatial point p . To render an image in an arbitrary view, σ and c are accumulated along viewing rays. Formally, we denote a viewing ray emitted from the optical center of a camera through a given pixel on the image plane by $r(m) = o + md$, then an approximation of the pixel color is

$$\hat{C}(r) = \mathcal{R}(r, c, \sigma) = \sum_{k=1}^K T(m_k) \alpha(\sigma(m_k) \delta_k) c(m_k), \quad (1)$$

where $\mathcal{R}(r, c, \sigma)$ is the volumetric rendering of the color c with the density σ ; $\{m_k\}_{k=1}^K$ is a set of discretely sampled points between the near and the far plane of the camera; $\delta_k = m_{k+1} - m_k$ is the distance between the current sampling point and the next one; $T(m_k) = \exp\left(-\sum_{k'=1}^{k-1} \sigma(m_{k'}) \delta_{k'}\right)$, and $\alpha(x) = 1 - \exp(-x)$. NeRF learns a radiance field with an MLP in the form of

$$[\sigma(m), c(m)] = f_{\theta}(\gamma_p(r(m)), \gamma_d(d)), \quad (2)$$

where θ is the model parameter, $\gamma_p(\cdot)$ and $\gamma_d(\cdot)$ are fixed positional encoding functions for positions and directions.

The network parameters are optimized by the loss

$$\mathcal{L} = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M \left\| \mathbf{C}(r_{ij}) - \hat{\mathbf{C}}(r_{ij}) \right\|_2^2, \quad (3)$$

where N is the number of images, M is the number of rays in each image, and r_{ij} is the j^{th} ray in the i^{th} image.

4. Method

Our method learns to reconstruct a person from synchronized multi-view video frames and animates the subject with novel poses. This is achieved by learning a canonical neural radiance field [14] of a human body in X-pose. This canonical radiance field is anchored to SMPL [11] so that we can animate the radiance field by manipulating SMPL. In Fig. 2, we show the pipeline of our method with barycentric mapping (Sec. 4.1) and dual-space NeRF (Sec. 4.2). The dual-space NeRF includes two networks: a Body MLP (Sec. 4.2.1) to model a human body in the canonical space and a Lighting MLP (Sec. 4.2.2) to capture the location-dependent lighting in the world space. And the barycentric mapping bridges the canonical space and the world space.

4.1. Barycentric Mapping

Considering the sparsity of views and the ambiguity caused by smooth regions, it is tough to learn robust correspondences across frames purely with images. Therefore, we adopt SMPL [11] as a geometric prior of the human body by anchoring spatial points on the faces of SMPL.

4.1.1 Position mapping

For a point \mathbf{p}^w in the world space (Fig. 3), we first determine its closest face F_i^w by measuring its distances to the mean of vertex positions of each face. Then, we set up a local description of the point \mathbf{p}^w by (u, v, h) , where (u, v) is the barycentric coordinate of the projection of \mathbf{p}^w on the face F_i^w , and h is the signed distance from F_i^w . Based on the corresponding face of F_i^w in the canonical space, *i.e.*, F_i^c , we can compute the corresponding point of \mathbf{p}^w as:

$$\mathbf{p}^c = \mathbf{o}^c + uu^c + vv^c + h \frac{\mathbf{u}^c \times \mathbf{v}^c}{\|\mathbf{u}^c \times \mathbf{v}^c\|}, \quad (4)$$

where \mathbf{o}^c is the first vertex of the face F_i^c , \mathbf{u}^c and \mathbf{v}^c are two edge vectors of F_i^c starting from the vertex \mathbf{o}^c . Note that the mapping can be conducted in an inverse direction.

4.1.2 Direction mapping

Based on the position mapping, we can bridge direction vectors between spaces, also in a differentiable manner. For the example in Fig. 3, we first represent the direction vector \mathbf{n}^c by its starting point \mathbf{p}^c and its ending point $\mathbf{p}_e^c = \mathbf{p}^c + \mathbf{n}^c$. Then we apply the position mapping described above to get the corresponding positions in the world space, *i.e.*, \mathbf{p}^w and \mathbf{p}_e^w . Finally, the warped direction vector in the world space can be obtained by:

$$\mathbf{n}^w = \frac{\mathbf{p}_e^w - \mathbf{p}^w}{\|\mathbf{p}_e^w - \mathbf{p}^w\|}. \quad (5)$$

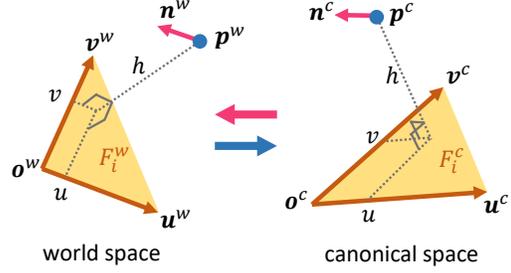


Figure 3: The barycentric mapping of positions and directions. For a point \mathbf{p}^w in the world space, we first find its closest face F_i^w , whose corresponding face in the canonical space is F_i^c . Then, we represent \mathbf{p}^w as the barycentric coordinate (u, v) of its projection on F_i^w and its signed distance h from F_i^w . Finally, we compute the counterpart of \mathbf{p}^w in the canonical space, *i.e.*, \mathbf{p}^c , with the same local representation upon F_i^c . The barycentric mapping can also be used to transform direction vectors.

4.2. Dual-Space NeRF

NeRF [14] astonishes the community for its high rendering realism, especially for the view-dependent visual effects. Most importantly, NeRF is formulated as a function of merely a point position \mathbf{p} and a viewing direction \mathbf{d} . Physically, the shading of a point also depends on the lighting condition and the surface normal, but NeRF omits them because they can be fully determined by the point position for a static scene.

However, for animatable human reconstruction and animation, the lighting condition is static only in the world space while the body shape and the surface normal is consistent only in the canonical space. Therefore, we have to learn a Body MLP and a Lighting MLP in separate spaces.

4.2.1 Body MLP

Given a point position \mathbf{p}^w in the world space (also called the observed space), we use the barycentric mapping (Sec. 4.1.1) to obtain its corresponding point in the canonical space, *i.e.*, \mathbf{p}^c , which is the main input of the Body MLP. Motivated by SCANimate [22], we encode the quaternion matrix of the joints of SMPL with a tiny MLP and get the pose feature \mathbf{J} . To prevent artifacts and blur in the results, we also learn a latent embedding $\ell_i \in \mathbb{R}^8$ for each video frame i to model the random variations that cannot be fully determined by the body pose.

Formally, we feed the canonical point position \mathbf{p}^c , pose features \mathbf{J} , and latent embedding ℓ_i into the Body MLP to predict density σ and texture $\mathbf{t} \in \mathbb{R}^3$. Here, Eq. (2) is reformulated as:

$$[\sigma, \mathbf{t}] = f_{\theta_1}(\gamma_{\mathbf{p}}(\mathbf{p}^c), \mathbf{J}, \ell_i). \quad (6)$$

For the Body MLP, density σ models the static shape and texture t models the true color of a human body in X-pose, both independent of the viewing direction. Moreover, the surface normal at the position p^c can be obtained by the normalized gradient of density σ with respect to p^c [23, 2]:

$$\mathbf{n}^c = -\frac{\nabla\sigma}{\|\nabla\sigma\|}. \quad (7)$$

4.2.2 Lighting MLP

To illustrate the necessity of a separate Lighting MLP, we consider a point p^c on a hand of a subject in the canonical space. When the subject waves the hand, the corresponding location of p^c in the world space moves, resulting in a change of the lighting condition. Therefore, from the perspective of the point p^c , the lighting condition varies with the body pose. Although the per-frame (or per-pose) lighting embeddings used by previous methods [19, 17] do help relieve the inconsistency of lighting in the training frames, they cannot generalize to novel poses. In contrast, we learn a Lighting MLP in the world space, making the lighting condition independent of the body pose.

Concretely, we use the Lighting MLP to predict a lightness coefficient

$$s = f_{\theta_2}(\mathbf{p}^w, \mathbf{d}^w, \mathbf{n}^w), \quad (8)$$

where \mathbf{p}^w is the point position, \mathbf{d}^w is the viewing direction, and \mathbf{n}^w is the surface normal, all in the world space. During ray casting and point sampling, \mathbf{p}^w and \mathbf{d}^w are directly available while the surface normal \mathbf{n}^w is not. To obtain \mathbf{n}^w , we first use the barycentric mapping (Sec. 4.1.1) to find the corresponding point of \mathbf{p}^w in the canonical space, *i.e.*, p^c , then compute the surface normal of p^c , *i.e.*, \mathbf{n}^c , finally map \mathbf{n}^c back to the world space also with the barycentric mapping (Sec. 4.1.2). The lightness coefficient s is meant to scale the lightness of the texture t for shading with the color

$$\mathbf{c} = s\mathbf{t}, \quad (9)$$

and the final result is obtained by volumetric rendering with the density σ and the color \mathbf{c} . Here, we model the lighting condition with simply a lightness coefficient instead of a color vector or more complex models because the task is highly under-constrained, and suppressing the expressiveness of the Lighting MLP helps prevent overfitting (Fig. 8).

4.3. Implementation Details

We apply the neutral SMPL [11] model for body mesh fitting. Personalized shape parameters are used for each subject. The X-pose is chosen as the canonical pose. Our network consists of 2 parts: the Body MLP is an 8-layers MLP with a shortcut connected to the fifth layer, and the Lighting MLP is a 4-layers MLP. We use a 3-layer MLP

to learn the pose feature \mathbf{J} . Hidden layers are activated by ReLU. The per-frame embeddings are initialized with Gaussian distribution ($\mathcal{N}(0, 1)$). We use the Adam optimizer to train our network for 200 epochs. We set the learning rate to 0.0005 and exponentially downscale it until the last epoch to 10 times lower. Weight decay is not used. We use a batch size of 1 with 5000 rays per batch. We sample 64 points on each ray. To accelerate training and inference, we abandon the coarse-to-fine strategy [14, 20, 19]. Instead, we adopt the geometry-guided ray marching [9], which produces a tighter bound. Since we have instance-level human-parsing masks, we ensure that 5 percent of rays are sampled around the face in each iteration, and the other rays are randomly sampled in the 2D bounding box. All experiments are conducted on a GeForce RTX 2080 Ti GPU and take around two days to converge.

5. Experiments

5.1. Settings

Datasets. ZJU-MoCap [20] is a multi-view dataset containing 9 performers captured by 21 synchronized cameras. It provides estimated SMPL [11] parameters and instance-level human-parsing masks generated by an established method [5]. We follow the experimental settings of Neural Body [20] and AniNeRF [19]. Images corresponding to four uniformly distributed cameras are used for training and the rest for evaluation. We conduct experiments on 8 performers. Human3.6M [7] contains four-view videos with human poses collected by a marker-based motion capture system. Images corresponding to three views are used for training and one for evaluation. We use the same protocol as Neural Body [20] to generate SMPL [11] parameters and masks.

Metrics. We adopt three metrics to evaluate the rendering quality, including PSNR, SSIM, and LPIPS [29]. PSNR is a pixel-wise metric based on the mean squared error, which is sensitive to noises and random variations. SSIM measures the structural similarity based on luminance, contrast, and structure comparisons. LPIPS [29] measures the perceptual distance between an image pair in a deep feature space. Since most pixels in the datasets belong to the background, we calculate PSNR and SSIM only within the 2D foreground mask, which is obtained by projecting the 3D bounding box on to the image plane.

5.2. Image Synthesis in Novel Poses

Baselines. Since we focus on the generalization ability of the model under novel poses, we compare our method with two state-of-the-art methods [20, 19] on novel pose synthesis. For results on novel view synthesis, please refer to our supplementary material. Note that we cannot compare with some recent works [9, 27] since the official code

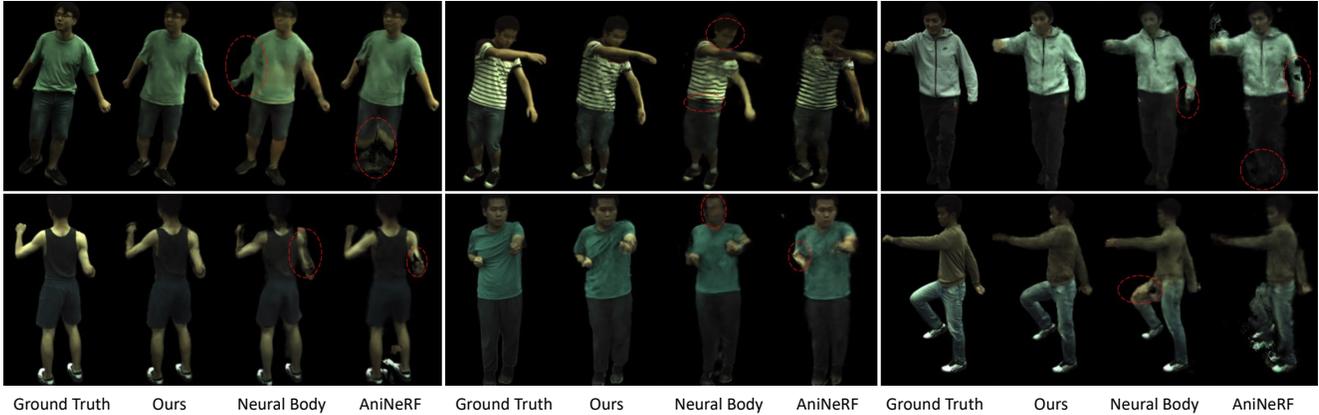


Figure 4: Results of novel pose synthesis on the ZJU-MoCap [20] dataset. Our results have fewer artifacts and are more visually pleasing. Note how our results better preserve the details on the faces of the subjects.

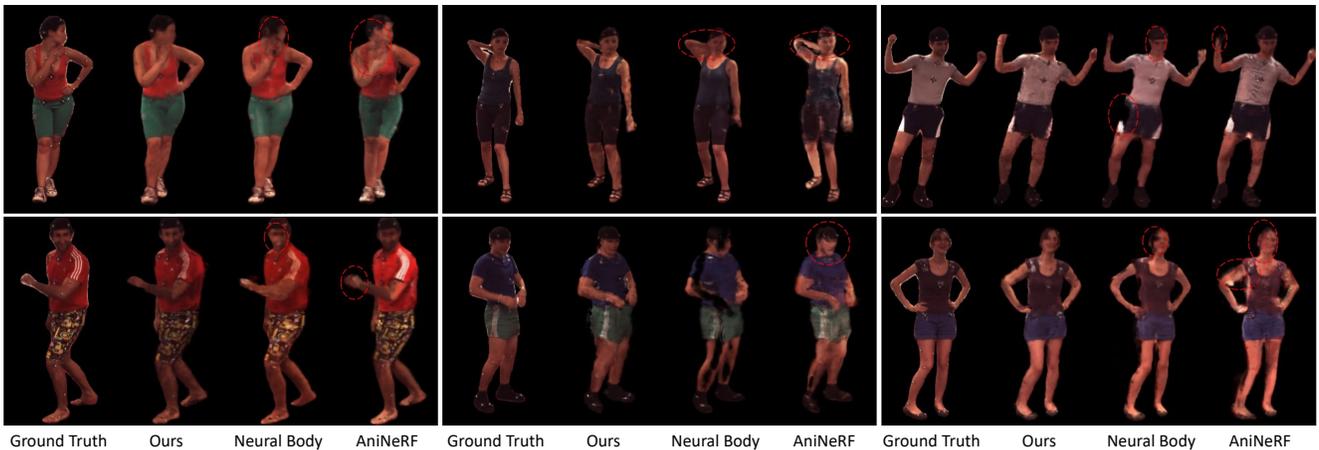


Figure 5: Results of novel pose synthesis on the Human3.6M [7] dataset. Our results show higher visual quality with fewer artifacts. Also, note the realistic lighting and shading effects in our results.

is not released so far. Neural Body [20] provides result images on both ZJU-MoCap and Human3.6M datasets. So we directly evaluate their results with our metrics. AniNeRF [19] releases results only on Human3.6M. So we run the official code of AniNeRF on ZJU-MoCap and conduct the same evaluation as above.

An implementation detail of NeuralBody and AniNeRF is that they only cast rays within the ground-truth human mask, which leaves them an advantage in comparisons. We argue that this operation is unreasonable because ground-truth masks are not always available for novel poses. Therefore, our method does not leverage the ground-truth human masks despite being at a disadvantage. Since the Lighting MLP is undefined beyond the movement range of a subject in the training frames, we place the avatar at the mean position of the training sequence when querying the Lighting MLP for novel poses. Likewise, the latent embedding is also not defined for novel poses, so we set it to zeros as done by previous work [31].

Comparisons on novel pose synthesis. As shown in Tab. 1 and Tab. 2, our method achieves the best PSNR and SSIM scores compared to two strong baselines. Since previous works [30, 10, 9] show that higher PSNR and SSIM scores do not guarantee better visual quality of images, we report LPIPS as a perceptual measurement, on which our method also shows advantages. According to Fig. 4 and Fig. 5, our method produces fewer artifacts than AniNeRF [19], indicating a better correspondences across frames. Meanwhile, the lighting conditions in our results are closer to the ground truths and the details are easier to recognize thanks to the reasonable decoupling of body properties and the environmental lighting. As shown in Fig. 4 and Fig. 5, Neural Body [20] tends to produce wrong body structures on Human3.6M when an unseen pose is far from the seen ones in the training set. While, our method produces visually pleasing results under novel poses, demonstrating the robustness of the barycentric mapping and the correctness of our lighting model.

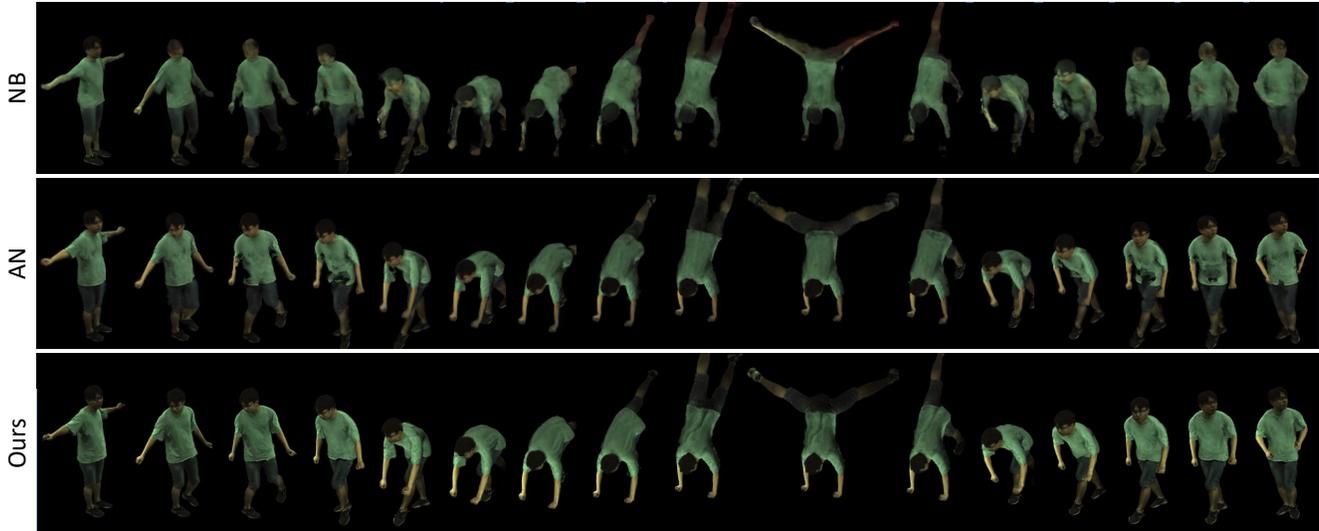


Figure 6: Results of extreme pose synthesis on the “hand stand” sequence from AMASS [12]. “NB” means Neural Body[20], and “AN” means AniNeRF [19]. Neural Body produces corrupted limbs and faces. AniNeRF makes artifacts and blurs. Our method renders sharp images with clear details and realistic lighting.

	PSNR \uparrow			SSIM \uparrow			LPIPS \downarrow		
	NB	AN	Ours	NB	AN	Ours	NB	AN	Ours
Twirl	23.853	22.800	24.300	0.902	0.863	0.910	0.056	0.078	0.045
Taichi	19.606	18.470	19.533	0.853	0.795	0.862	0.054	0.092	0.047
Warmup	23.907	23.280	24.675	0.909	0.901	0.919	0.036	0.056	0.031
Punch1	25.671	25.550	26.042	0.881	0.872	0.889	0.044	0.053	0.035
Punch2	21.595	21.916	22.395	0.870	0.838	0.881	0.058	0.089	0.050
Swing1	25.736	18.438	25.776	0.908	0.670	0.914	0.049	0.212	0.045
Swing2	23.802	21.870	24.360	0.878	0.836	0.888	0.055	0.090	0.049
Swing3	23.248	17.694	23.247	0.893	0.792	0.894	0.055	0.206	0.053
Average	23.427	21.252	23.791	0.887	0.821	0.894	0.051	0.110	0.044

Table 1: Comparison with baselines on novel pose synthesis on ZJU-MoCap [20], “NB” means Neural Body [20], and “AN” means AniNeRF [19]. Our method outperforms both baselines with a clear margin.

	PSNR \uparrow			SSIM \uparrow			LPIPS \downarrow		
	NB	AN	Ours	NB	AN	Ours	NB	AN	Ours
S1	21.932	19.955	23.206	0.873	0.855	0.886	0.026	0.029	0.027
S5	23.332	20.022	23.025	0.893	0.840	0.886	0.022	0.025	0.021
S6	23.263	23.637	24.059	0.888	0.882	0.893	0.041	0.046	0.038
S7	22.398	21.762	22.913	0.888	0.869	0.885	0.029	0.033	0.027
S8	20.779	21.631	22.659	0.872	0.877	0.889	0.035	0.032	0.031
S9	22.868	21.948	24.143	0.880	0.871	0.887	0.029	0.034	0.028
S11	23.538	22.547	24.842	0.879	0.875	0.894	0.032	0.030	0.029
Average	22.587	21.643	23.550	0.882	0.867	0.889	0.030	0.033	0.029

Table 2: Comparison with baselines on novel pose synthesis on the Human3.6M [7] dataset. “NB” means Neural Body [20], and “AN” means AniNeRF [19]. Our method outperforms both baselines in most cases.

Comparison on extreme pose synthesis. Since the difference between the training and the test poses in a dataset may not be large enough, we compare the methods on a challenging pose sequence from the AMASS [12] database. As shown in Fig. 6, Neural Body produces corrupted limbs and faces, which show the limitation of the convolution-based solution. AniNeRF makes artifacts and blurs due to the instability of the spatially interpolated LBS weights and the poor generalization ability of the neural blending weights. Our method renders sharp images with clear details and realistic lighting thanks to the stable correspondences and reasonable decoupling of body properties and the environmental lighting. Please refer to our supplementary video for animated results.

5.3. Ablation Studies

To verify the effectiveness of our main components, we conduct ablation studies on the “Twirl” sequence of ZJU-MoCap [20] in terms of novel pose and novel view synthesis. All models are trained for the same number of epochs (100) for a fair comparison.

Barycentric mapping. To validate the virtue of the barycentric mapping, we replace it with the inverse LBS algorithm. For the blending weights, we follow the strategy of AniNeRF [19], which interpolates the blending weights from nearby SMPL vertices. The second row of Tab. 3 shows the clear superiority of the barycentric mapping, especially on novel poses. In the visual comparison (Fig. 7), inverse LBS produces artifacts around movement-frequent places like armpits and feet, while the barycentric mapping still performs well.

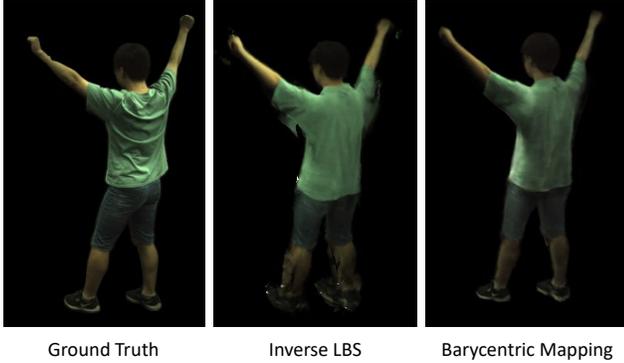


Figure 7: Visual comparison between the barycentric mapping and inverse LBS. Inverse LBS with interpolated blending weights produces artifacts near movement-frequent places like armpits and feet while the barycentric mapping renders clear results.

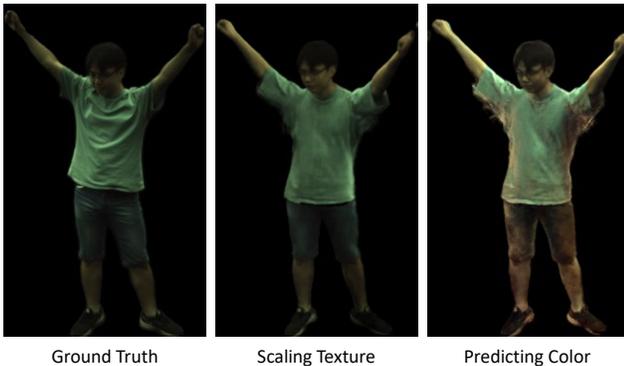


Figure 8: Ablation study of the Lighting MLP. We test with an alternative design of Lighting MLP that directly predicts RGB values instead of predicting the lightness coefficient for texture scaling. However, the Lighting MLP that directly predicts colors tends to overfit the training frames and produces distorted colors on the novel pose.

	PSNR \uparrow		SSIM \uparrow		LPIPS \downarrow	
	View	Pose	View	Pose	View	Pose
Full model	<u>31.090</u>	24.216	<u>0.970</u>	0.911	<u>0.023</u>	0.044
Replace BM with inverse LBS	30.758	23.301	0.968	0.895	0.026	0.055
w/o Lighting MLP	30.696	23.465	0.967	<u>0.906</u>	0.027	0.049
Lighting MLP (predicting color)	31.270	<u>23.570</u>	0.971	0.905	0.019	<u>0.047</u>

Table 3: Ablation studies. “View” refers to novel view synthesis, and “Pose” refers to novel pose synthesis. “BM” means the barycentric mapping. Bold values are the best scores, and underlined values are the second best.

w/o Lighting MLP. The Lighting MLP plays a vital role in solving the lighting inconsistency. It models the correct location-dependent lighting in the world space and benefits high-fidelity rendering. To validate its usefulness, we disable the Lighting MLP when training and rendering. Then,

our model degenerates to a vanilla NeRF defined in the canonical space. The third row in Tab. 3 shows significant degradation in all metrics.

Lighting MLP predicts color directly. We also try an alternative design of the Lighting MLP that takes in the body texture and outputs the color instead of predicting the lightness coefficient for texture scaling. As shown in the fourth row in Tab. 3, this alternative Lighting MLP perform better on novel view synthesis due to higher expressiveness but perform worse on novel poses. Thus, the alternative Lighting MLP is just overfitting the colors in the training frames instead of learning the environmental lighting. Similar conclusion can be drawn from Fig. 8, where the alternative Lighting MLP predicts distorted colors on the novel pose. We show animated results in our supplementary video.

6. Conclusion

In this paper, we focus on the generalization problem of human body reconstruction and animation. We propose to model the human body and the lighting condition in separate spaces. To bridge the canonical space and the world space, we propose the barycentric mapping, which helps us to transform point positions and surface normals of a human body between the two spaces, enabling rendering in the world space with body properties from the canonical space. Most importantly, the barycentric mapping can directly generalize to novel poses without additional input or network training. Thanks to the reasonable decoupling of body properties and lighting conditions, we obtain clear improvements upon two strong baselines.

7. Limitations and Potential Impacts

Our method uses SMPL [11] as a proxy to build connections between the world space and the canonical space. Therefore, it strongly relies on an accurate SMPL fitting. In scenarios where SMPL parameters cannot be precisely obtained, our method is likely to fail. Also, our approach does not model long-range dependencies and thus is unable to deal with a performer in a long dress. Our work reconstructs the appearances of subjects and animates them with public video datasets. Currently, the rendering realism is far from fooling people, but attention should be paid to future versions of related technologies for potential misusing.

Acknowledgments: The work is supported by National Key R&D Program of China (2018AAA0100704), NSFC #61932020, #62172279, Science and Technology Commission of Shanghai Municipality (Grant No. 20ZR1436000), Program of Shanghai Academic Research Leader, and “Shuguang Program” supported by Shanghai Education Development Foundation and Shanghai Municipal Education Commission.

References

- [1] Thiemo Alldieck, Hongyi Xu, and Cristian Sminchisescu. imghum: Implicit generative models of 3d human shape and articulated pose. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5461–5470, 2021. 3
- [2] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T. Barron, Ce Liu, and Hendrik P.A. Lensch. Nerf: Neural reflectance decomposition from image collections. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 5
- [3] Xu Chen, Yufeng Zheng, Michael J Black, Otmar Hilliges, and Andreas Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. In *International Conference on Computer Vision (ICCV)*, 2021. 2
- [4] Boyang Deng, John P Lewis, Timothy Jeruzalski, Gerard Pons-Moll, Geoffrey Hinton, Mohammad Norouzi, and Andrea Tagliasacchi. Nasa neural articulated shape approximation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VII 16*, pages 612–628. Springer, 2020. 2
- [5] Ke Gong, Xiaodan Liang, Yicheng Li, Yimin Chen, Ming Yang, and Liang Lin. Instance-level human parsing via part grouping network. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 770–785, 2018. 5
- [6] Yudong Guo, Keyu Chen, Sen Liang, Yongjin Liu, Hujun Bao, and Juyong Zhang. Ad-nerf: Audio driven neural radiance fields for talking head synthesis. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 3
- [7] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. 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, 2013. 2, 5, 6, 7, 11, 13
- [8] Youngjoong Kwon, Dahun Kim, Duygu Ceylan, and Henry Fuchs. Neural human performer: Learning generalizable radiance fields for human performance rendering. *Advances in Neural Information Processing Systems*, 34, 2021. 1, 3
- [9] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. *ACM Trans. Graph.(ACM SIGGRAPH Asia)*, 2021. 1, 2, 3, 5, 6
- [10] Wen Liu, Zhixin Piao, Jie Min, Wenhan Luo, Lin Ma, and Shenghua Gao. Liquid warping gan: A unified framework for human motion imitation, appearance transfer and novel view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5904–5913, 2019. 6
- [11] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, Oct. 2015. 2, 3, 4, 5, 8, 11
- [12] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. AMASS: Archive of motion capture as surface shapes. In *International Conference on Computer Vision*, pages 5442–5451, Oct. 2019. 7, 11
- [13] Marko Mihajlovic, Yan Zhang, Michael J Black, and Siyu Tang. LEAP: Learning articulated occupancy of people. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2021. 2
- [14] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 1, 2, 3, 4, 5
- [15] Atsuhiko Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. Neural articulated radiance field. In *International Conference on Computer Vision*, 2021. 1, 2
- [16] Pablo Palafox, Aljaž Božič, Justus Thies, Matthias Nießner, and Angela Dai. Npms: Neural parametric models for 3d deformable shapes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12695–12705, 2021. 2
- [17] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021. 3, 5
- [18] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228*, 2021. 3
- [19] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *ICCV*, 2021. 1, 2, 3, 5, 6, 7, 11, 12
- [20] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *CVPR*, 2021. 1, 2, 3, 5, 6, 7, 11, 12
- [21] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021. 3
- [22] Shunsuke Saito, Jinlong Yang, Qianli Ma, and Michael J. Black. SCANimate: Weakly supervised learning of skinned clothed avatar networks. In *Proceedings IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2021. 2, 4
- [23] Pratul P. Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T. Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *CVPR*, 2021. 5
- [24] Shih-Yang Su, Frank Yu, Michael Zollhöfer, and Helge Rhodin. A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose. In *Advances in Neural Information Processing Systems*, 2021. 1, 2
- [25] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhofer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view

- synthesis of a dynamic scene from monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12959–12970, 2021. 3
- [26] Chung-Yi Weng, Brian Curless, Pratul P. Srinivasan, Jonathan T. Barron, and Ira Kemelmacher-Shlizerman. HumanNeRF: Free-viewpoint rendering of moving people from monocular video. *CVPR*, 2022. 1, 3
- [27] Hongyi Xu, Thiemo Alldieck, and Cristian Sminchisescu. H-nerf: Neural radiance fields for rendering and temporal reconstruction of humans in motion. *Advances in Neural Information Processing Systems*, 34, 2021. 1, 3, 5
- [28] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33, 2020. 2
- [29] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 5
- [30] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 6
- [31] Zerong Zheng, Han Huang, Tao Yu, Hongwen Zhang, Yandong Guo, and Yebin Liu. Structured local radiance fields for human avatar modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15893–15903, 2022. 6

8. Implementation Details

Empirically, the points that are too far away from the SMPL [11] mesh should not follow the movement of the human body. To make the learned neural radiance field less noisy, we exclude these outlier points during volumetric rendering by setting their density values to zeros, and this operation brings a slight improvement. To determine whether a given point is an outlier, we apply the following algorithm:

Algorithm 1 Outlier detection

Input: a sampled point p ; SMPL faces \mathcal{F} ; α ; β ; γ

Output: outlier mask m

- 1: $f \leftarrow \text{Find_NN_Mesh}(p, \mathcal{F})$
 - 2: $u, v, h \leftarrow \text{Compute_UV_SignedDistance}(p, f)$
 - 3: return $(u, v < \alpha)$ OR $(u, v > \beta)$ OR $(|h| > \gamma)$
-

We set $\alpha = -4$, $\beta = 5$, and $\gamma = 0.1$ experimentally.

9. Novel View Synthesis

Although our method targets novel pose synthesis, we still report results on novel view synthesis. We compare the methods on the ZJU-Mocap [20] dataset in Tab. 4, where Neural Body [20] exhibits superiority on PSNR and SSIM, and achieves comparable LPIPS to our method. Neural Body anchors features on the vertices of SMPL and diffuses them into a feature grid before volumetric rendering, avoiding establishing correspondences across frames. It is favorable for novel view synthesis but degrades on novel poses. Our method exceeds AniNeRF [19], which explicitly builds correspondences across frames like our method, by a large margin in all metrics. And the results of quantitative comparisons are consistent with the rendering results shown in Fig. 9. Our barycentric mapping establishes more robust correspondences across poses, producing fewer structural artifacts such as the extra feet in Fig. 9. On the Human3.6M [7] dataset, our method shows outstanding performance compared to both baselines as shown in Tab. 5. Corresponding visual comparisons are shown in Fig. 10, where our method produces realistic textures, lights, and shades. Note that Human3.6M [7] is noisier with higher errors in the fitted SMPL parameters and unclear boundaries in foreground masks compared to ZJU-MoCap [20]. This explains the obvious degradation of Neural Body on this dataset and indicates higher robustness towards imperfect SMPL fits and noisy data of our method.

10. Novel Pose Synthesis

To further verify the generalization ability of our method on novel poses, we show animated results in our supplementary video. We animate the subjects of each dataset with a

	PSNR \uparrow			SSIM \uparrow			LPIPS \downarrow		
	NB	AN	ours	NB	AN	ours	NB	AN	ours
Twirl	30.455	28.050	31.623	0.966	0.940	0.972	0.028	0.038	0.024
Taichi	27.199	19.660	25.829	0.960	0.849	0.950	0.022	0.065	0.026
Warmup	27.962	24.970	<u>27.347</u>	0.952	0.920	<u>0.950</u>	0.026	0.045	0.026
Punch1	28.659	25.760	28.487	0.928	0.870	0.925	<u>0.027</u>	0.054	0.026
Punch2	25.866	22.551	<u>25.208</u>	0.927	0.862	0.918	0.045	0.083	0.044
Swing1	29.618	23.717	<u>29.226</u>	0.946	0.869	<u>0.942</u>	0.030	0.074	<u>0.031</u>
Swing2	28.632	23.793	<u>28.494</u>	0.939	0.877	<u>0.933</u>	0.034	0.069	<u>0.035</u>
Swing3	27.583	17.351	<u>27.199</u>	0.936	0.760	0.930	0.034	0.205	<u>0.035</u>
Average	28.247	23.232	<u>27.927</u>	0.944	0.868	<u>0.940</u>	0.031	0.079	0.031

Table 4: Comparison with baselines on novel view synthesis on the ZJU-MoCap [20] dataset. “NB” means Neural Body, and “AN” means AniNeRF [19]. Bold values are the best scores, and underlined values are the second best. Our method outperforms AniNeRF [19] and is comparable to Neural Body [20] on the perceptual metric.

	PSNR \uparrow			SSIM \uparrow			LPIPS \downarrow		
	NB	AN	ours	NB	AN	ours	NB	AN	ours
S1	<u>22.716</u>	22.415	24.494	<u>0.893</u>	0.890	0.913	0.032	0.034	0.032
S5	<u>24.439</u>	23.228	24.819	<u>0.914</u>	0.891	0.915	<u>0.023</u>	0.027	0.022
S6	22.668	<u>22.689</u>	24.294	<u>0.884</u>	0.866	0.894	<u>0.029</u>	0.034	0.028
S7	<u>22.991</u>	21.793	23.933	0.911	0.886	0.909	0.025	0.030	<u>0.026</u>
S8	21.570	<u>22.666</u>	23.234	0.890	<u>0.897</u>	0.911	0.034	<u>0.031</u>	0.027
S9	24.121	<u>24.694</u>	25.691	<u>0.907</u>	<u>0.907</u>	0.914	0.029	0.034	0.029
S11	23.537	<u>24.594</u>	25.622	0.892	<u>0.903</u>	0.914	0.039	<u>0.035</u>	0.032
Average	23.149	<u>23.154</u>	24.584	<u>0.899</u>	0.891	0.910	<u>0.030</u>	0.032	0.028

Table 5: Comparison with baselines on novel view synthesis on the Human3.6M [7] dataset. “NB” means Neural Body [20], and “AN” means AniNeRF [19]. Bold values are the best scores, and underlined values are the second best. Our method achieves the highest performances on PSNR, SSIM, and LPIPS.

pose sequence from the other dataset. In pursuit of diversity and complexity, we select the “Swing3” sequence from ZJU-Mocap and the “S9” sequence from Human3.6M. Besides, we compare our method with the baselines on three more challenging pose sequences from the AMASS [12] database. Neural Body [20] can hardly generalize to extreme poses. AniNeRF [19] lacks high-fidelity details such as wrinkles and lighting. In extreme pose #3, AniNeRF produces artifacts like corrupted faces, while our method gives stable results.

11. Lighting MLP Validation

To further interpret our Lighting MLP, we visualize its effect in our supplementary video by manipulating the querying points of the Lighting MLP. By rotating the querying points along a certain axis, we can observe a change in the lighting on the human body. While, if the subject walks out of the moving boundary in the training frames, the rendering results will be dim because those locations are undefined for the Lighting MLP.



Figure 9: Results of novel view synthesis on the ZJU-MoCap [20] dataset. The results of Neural Body [20] and our method exhibit fewer artifacts compared with AniNeRF [19].

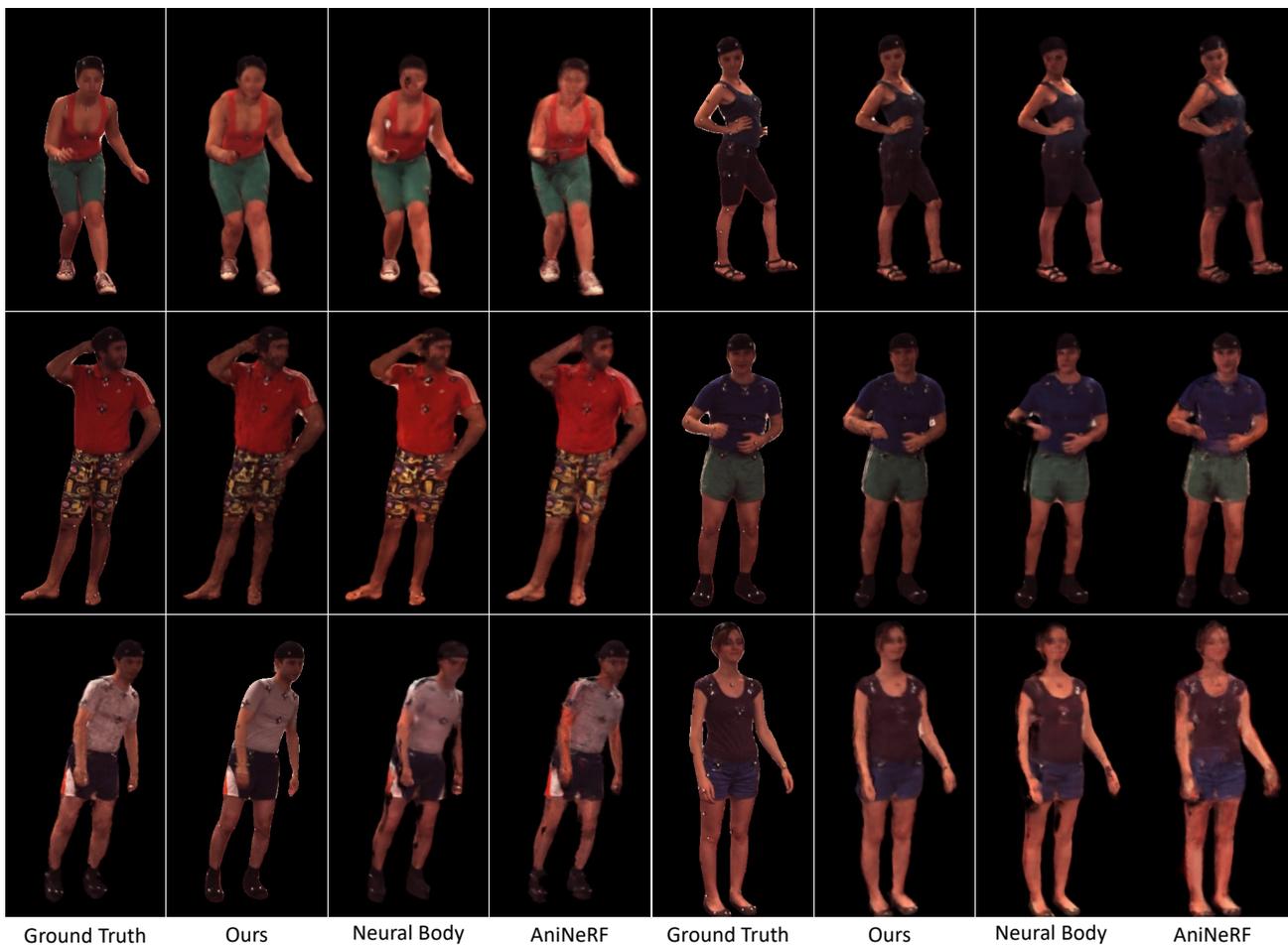


Figure 10: Results of novel view synthesis on the Human3.6M [7] dataset. Our results show higher fidelity with clear textures and realistic lighting and shading.