# ARAH: Animatable Volume Rendering of Articulated Human SDFs

Shaofei Wang[1], Katja Schwarz[2,3], Andreas Geiger[2,3], and Siyu Tang[1]

[1] ETH Zürich
[2] Max Planck Institute for Intelligent Systems, Tübingen
[3] University of Tübingen

**Abstract.** Combining human body models with differentiable rendering has recently enabled animatable avatars of clothed humans from sparse sets of multi-view RGB videos. While state-of-the-art approaches achieve a realistic appearance with neural radiance fields (NeRF), the inferred geometry often lacks detail due to missing geometric constraints. Further, animating avatars in out-of-distribution poses is not yet possible because the mapping from observation space to canonical space does not generalize faithfully to unseen poses. In this work, we address these shortcomings and propose a model to create animatable clothed human avatars with detailed geometry that generalize well to out-of-distribution poses. To achieve detailed geometry, we combine an articulated implicit surface representation with volume rendering. For generalization, we propose a novel joint root-finding algorithm for simultaneous ray-surface intersection search and correspondence search. Our algorithm enables efficient point sampling and accurate point canonicalization while generalizing well to unseen poses. We demonstrate that our proposed pipeline can generate clothed avatars with high-quality pose-dependent geometry and appearance from a sparse set of multi-view RGB videos. Our method achieves state-of-the-art performance on geometry and appearance reconstruction while creating animatable avatars that generalize well to out-of-distribution poses beyond the small number of training poses.

**Keywords:** 3D Computer Vision, Clothed Human Modeling, Cloth Modeling, Neural Rendering, Neural Implicit Functions

## 1 Introduction

Reconstruction and animation of clothed human avatars is a rising topic in computer vision research. It is of particular interest for various applications in AR/VR and the future metaverse. Various sensors can be used to create clothed human avatars, ranging from 4D scanners over depth sensors to simple RGB cameras. Among these data sources, RGB videos are by far the most accessible and user-friendly choice. However, they also provide the least supervision, making this setup the most challenging for the reconstruction and animation of clothed humans.
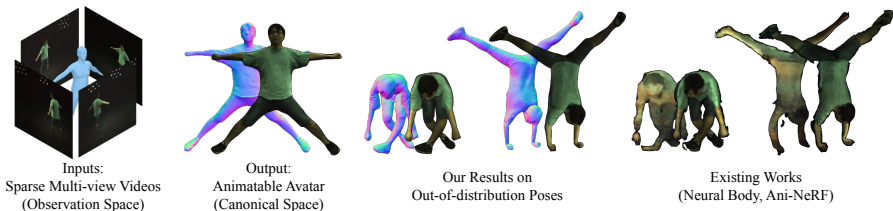
Fig. 1: **Detailed Geometry and Generalization to Extreme Poses.** Given sparse multi-view videos with SMPL fittings and foreground masks, our approach synthesizes animatable clothed avatars with realistic pose-dependent geometry and appearance. While existing works, *e.g.* Neural Body [60] and Ani-NeRF [58], struggle with generalizing to unseen poses, our approach enables avatars that can be animated in extreme out-of-distribution poses.

Traditional works in clothed human modeling use explicit mesh [1, 2, 6, 7, 18, 19, 31, 35, 56, 69, 75, 85, 90] or truncated signed distance fields (TSDFs) of fixed grid resolution [36, 37, 73, 83, 88] to represent the geometry of humans. Textures are often represented by vertex colors or UV-maps. With the recent success of neural implicit representations, significant progress has been made towards modeling articulated clothed humans. PIFu [65] and PIFuHD [66] are among the first works that propose to model clothed humans as continuous neural implicit functions. ARCH [25] extends this idea and develops animatable clothed human avatars from monocular images. However, this line of works does not handle dynamic pose-dependent cloth deformations. Further, they require ground-truth geometry for training. Such ground-truth data is expensive to acquire, limiting the generalization of these methods.

Another line of works removes the need for ground-truth geometry by utilizing differentiable neural rendering. These methods aim to reconstruct humans from a sparse set of multi-view videos with only image supervision. Many of them use NeRF [49] as the underlying representation and achieve impressive visual fidelity on novel view synthesis tasks. However, there are two fundamental drawbacks of these existing approaches: (1) the NeRF-based representation lacks proper geometric regularization, leading to inaccurate geometry. This is particularly detrimental in a sparse multi-view setup and often results in artifacts in the form of erroneous color blobs under novel views or poses. (2) Existing approaches condition their NeRF networks [60] or canonicalization networks [58] on inputs in observation space. Thus, they cannot generalize to unseen out-of-distribution poses.

In this work, we address these two major drawbacks of existing approaches. (1) We improve geometry by building an articulated signed-distance-field (SDF) representation for clothed human bodies to better capture the geometry of clothed humans and improve the rendering quality. (2) In order to render the SDF, we develop an efficient joint root-finding algorithm for the conversion from observation space to canonical space. Specifically, we represent clothed human

avatars as a combination of a forward linear blend skinning (LBS) network, an implicit SDF network, and a color network, all defined in canonical space and do not condition on inputs in observation space. Given these networks and camera rays in observation space, we apply our novel joint root-finding algorithm that can efficiently find the iso-surface points in observation space and their correspondences in canonical space. This enables us to perform efficient sampling on camera rays around the iso-surface. All network modules can be trained with a photometric loss in image space and regularization losses in canonical space.

We validate our approach on the ZJU-MoCap [60] and the H36M [26] dataset. Our approach generalizes well to unseen poses, enabling robust animation of clothed avatars even under out-of-distribution poses where existing works fail, as shown in Fig. 1. We achieve significant improvements over state-of-the-arts for novel pose synthesis and geometry reconstruction, while also outperforming state-of-the-arts in the novel view synthesis task on training poses. Code and data are available at https://neuralbodies.github.io/arah/.

## 2   Related Works

**Clothed Human Modeling with Explicit Representations:** Many explicit mesh-based approaches represent cloth deformations as deformation layers [1, 2, 6–8] added to minimally clothed parametric human body models [5, 21, 28, 39, 54, 57, 82]. Such approaches enjoy compatibility with parametric human body models but have difficulties in modeling large garment deformations. Other mesh-based approaches model garments as separate meshes [18, 19, 31, 35, 56, 69, 75, 85, 90] in order to represent more detailed and physically plausible cloth deformations. However, such methods often require accurate 3D-surface registration, synthetic 3D data or dense multi-view images for training and the garment meshes need to be pre-defined for each cloth type. More recently, point-cloud-based explicit methods [40, 42, 89] also showed promising results in modeling clothed humans. However, they still require explicit 3D or depth supervision for training, while our goal is to train using sparse multi-view RGB supervision alone.

**Clothed Humans as Implicit Functions:** Neural implicit functions [13, 44, 45, 55, 61] have been used to model clothed humans from various sensor inputs including monocular images [22, 23, 25, 33, 64–66, 72, 80, 93], multi-view videos [30, 38, 52, 58, 60, 81], sparse point clouds [6, 14, 16, 77, 78, 94], or 3D meshes [11, 12, 15, 47, 48, 67, 74]. Among the image-based methods, [4, 23, 25] obtain animatable reconstructions of clothed humans from a single image. However, they do not model pose-dependent cloth deformations and require ground-truth geometry for training. [30] learns generalizable NeRF models for human performance capture and only requires multi-view images as supervision. But it needs images as inputs for synthesizing novel poses. [38, 52, 58, 60, 81] take multi-view videos as inputs and do not need ground-truth geometry during training. These methods generate

personalized per-subject avatars and only need 2D supervision. Our approach follows this line of work and also learns a personalized avatar for each subject.

**Neural Rendering of Animatable Clothed Humans:** Differentiable neural rendering has been extended to model animatable human bodies by a number of recent works [52, 58, 60, 63, 72, 81]. Neural Body [60] proposes to diffuse latent per-vertex codes associated with SMPL meshes in observation space and condition NeRF [49] on such latent codes. However, the conditional inputs of Neural Body are in the observation space. Therefore, it does not generalize well to out-of-distribution poses. Several recent works [52, 58, 72] propose to model the radiance field in canonical space and use a pre-defined or learned backward mapping to map query points from observation space to this canonical space. A-NeRF [72] uses a deterministic backward mapping defined by piecewise rigid bone transformations. This mapping is very coarse and the model has to use a complicated bone-relative embedding to compensate for that. Ani-NeRF [58] trains a backward LBS network that does not generalize well to out-of-distribution poses, even when fine-tuned with a cycle consistency loss for its backward LBS network for each test pose. Further, all aforementioned methods utilize a volumetric radiance representation and hence suffer from noisy geometry [53, 76, 86, 87]. In contrast to these works, we improve geometry by combining an implicit surface representation with volume rendering and improve pose generalization via iterative root-finding. H-NeRF [81] achieves large improvements in geometric reconstruction by co-training SDF and NeRF networks. However, code and models of H-NeRF are not publicly available. Furthermore, H-NeRF's canonicalization process relies on imGHUM [3] to predict an accurate signed distance in *observation space*. Therefore, imGHUM needs to be trained on a large corpus of posed human scans and it is unclear whether the learned signed distance fields generalize to out-of-distribution poses beyond the training set. In contrast, our approach does not need to be trained on any posed scans and it can generalize to extreme out-of-distribution poses.

**Concurrent Works:** Several concurrent works extend NeRF-based articulated models to improve novel view synthesis, geometry reconstruction, or animation quality [10, 24, 27, 32, 46, 59, 71, 79, 84, 92]. [92] proposes to jointly learn forward blending weights, a canonical occupancy network, and a canonical color network using differentiable surface rendering for head-avatars. In contrast to human heads, human bodies show much more articulation. Abrupt changes in depth also occur more frequently when rendering human bodies, which is difficult to capture with surface rendering [76]. Furthermore, [92] uses the secant method to find surface points. For each secant step, this needs to solve a root-finding problem from scratch. Instead, we use volume rendering of SDFs and formulate the surface-finding task of articulated SDFs as a joint root-finding problem that only needs to be solved once per ray. We remark that [27] proposes to formulate surface-finding and correspondence search as a joint root-finding problem to tackle geometry reconstruction from photometric and mask losses. However, they use pre-defined skinning fields and surface rendering. They also require esti-
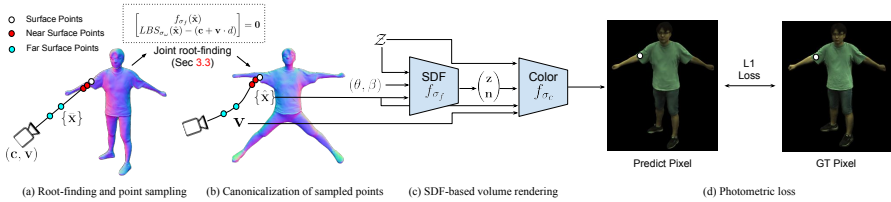
Fig. 2: **Overview of Our Pipeline.** (a) Given a ray $(\mathbf{c}, \mathbf{v})$ with camera center $\mathbf{c}$ and ray direction $\mathbf{v}$ in observation space, we jointly search for its intersection with the SDF iso-surface and the correspondence of the intersection point via a novel joint root-finding algorithm (Section 3.3). We then sample near/far surface points $\{\bar{\mathbf{x}}\}$. (b) The sampled points are mapped into canonical space as $\{\hat{\mathbf{x}}\}$ via root-finding. (c) In canonical space, we run an SDF-based volume rendering with canonicalized points $\{\hat{\mathbf{x}}\}$, local body poses and shape $(\theta, \beta)$, an SDF network feature $\mathbf{z}$, surface normals $\mathbf{n}$, and a per-frame latent code $\mathcal{Z}$ to predict the corresponding pixel value of the input ray (Section 3.4). (d) All network modules, including the forward LBS network $LBS_{\sigma_\omega}$, the canonical SDF network $f_{\sigma_f}$, and the canonical color network $f_{\sigma_c}$, are trained end-to-end with a photometric loss in image space and regularization losses in canonical space (Section 3.5).

mated normals from PIFuHD [66] while our approach achieves detailed geometry reconstructions without such supervision.

## 3  Method

Our pipeline is illustrated in Fig. 2. Our model consists of a forward linear blend skinning (LBS) network (Section 3.1), a canonical SDF network, and a canonical color network (Section 3.2). When rendering a specific pixel of the image in observation space, we first find the intersection of the corresponding camera ray and the observation-space SDF iso-surface. Since we model a canonical SDF and a forward LBS, we propose a novel joint root-finding algorithm that can simultaneously search for the ray-surface intersection and the canonical correspondence of the intersection point (Section 3.3). Such a formulation does not condition the networks on observations in observation space. Consequently, it can generalize to unseen poses. Once the ray-surface intersection is found, we sample near/far surface points on the camera ray and find their canonical correspondences via forward LBS root-finding. The canonicalized points are used for volume rendering to compose the final RGB value at the pixel (Section 3.4). The predicted pixel color is then compared to the observation using a photometric loss (Section 3.5). The model is trained end-to-end using the photometric loss and regularization losses. The learned networks represent a personalized animatable avatar that can robustly synthesize new geometries and appearances under novel poses (Section 4.1).

### 3.1   Neural Linear Blend Skinning

Traditional parametric human body models [5, 21, 39, 54, 57, 82] often use linear blend skinning (LBS) to deform a template model according to rigid bone transformations and skinning weights. We follow the notations of [78] to describe LBS. Given a set of $N$ points in canonical space, $\hat{\mathbf{X}} = \{\hat{\mathbf{x}}^{(i)}\}_{i=1}^{N}$, LBS takes a set of rigid bone transformations $\{\mathbf{B}_b\}_{b=1}^{24}$ as inputs, each $\mathbf{B}_b$ being a $4 \times 4$ rotation-translation matrix. We use 23 local transformations and one global transformation with an underlying SMPL [39] model. For a 3D point $\hat{\mathbf{x}}^{(i)} \in \hat{\mathbf{X}}$ [4], a skinning weight vector is defined as $\mathbf{w}^{(i)} \in [0,1]^{24}, \text{s.t.} \sum_{b=1}^{24} \mathbf{w}_b^{(i)} = 1$. This vector indicates the affinity of the point $\hat{\mathbf{x}}^{(i)}$ to each of the bone transformations $\{\mathbf{B}_b\}_{b=1}^{24}$. Following recent works [12, 48, 67, 78], we use a neural network $f_{\sigma_\omega}(\cdot) : \mathbb{R}^3 \mapsto [0,1]^{24}$ with parameters $\sigma_\omega$ to predict the skinning weights of any point in space. The set of transformed points $\bar{\mathbf{X}} = \{\bar{\mathbf{x}}^{(i)}\}_{i=1}^{N}$ is related to $\hat{\mathbf{X}}$ via:

$$\bar{\mathbf{x}}^{(i)} = LBS_{\sigma_f}\left(\hat{\mathbf{x}}^{(i)}, \{\mathbf{B}_b\}\right), \quad \forall i = 1, \ldots, N$$

$$\Longleftrightarrow \bar{\mathbf{x}}^{(i)} = \left(\sum_{b=1}^{24} f_{\sigma_\omega}(\hat{\mathbf{x}}^{(i)})_b \mathbf{B}_b\right)\hat{\mathbf{x}}^{(i)}, \quad \forall i = 1, \ldots, N \tag{1}$$

where Eq. (1) is referred to as the forward LBS function. The process of applying Eq. (1) to all points in $\hat{\mathbf{X}}$ is often referred to as *forward skinning*. For brevity, for the remainder of the paper, we drop $\{\mathbf{B}_b\}$ from the LBS function and write $LBS_{\sigma_\omega}(\hat{\mathbf{x}}^{(i)}, \{\mathbf{B}_b\})$ as $LBS_{\sigma_\omega}(\hat{\mathbf{x}}^{(i)})$.

### 3.2   Canonical SDF and Color Networks

We model an articulated human as a neural SDF $f_{\sigma_f}(\hat{\mathbf{x}}, \theta, \beta, \mathcal{Z})$ with parameters $\sigma_f$ in canonical space, where $\hat{\mathbf{x}}$ denotes the canonical query point, $\theta$ and $\beta$ denote local poses and body shape of the human which capture pose-dependent cloth deformations, and $\mathcal{Z}$ denotes a per-frame optimizable latent code which compensates for time-dependent dynamic cloth deformations. For brevity, we write this neural SDF as $f_{\sigma_f}(\hat{\mathbf{x}})$ in the remainder of the paper.

Similar to the canonical SDF network, we define a canonical color network with parameters $\sigma_c$ as $f_{\sigma_c}(\hat{\mathbf{x}}, \mathbf{n}, \mathbf{v}, \mathbf{z}, \mathcal{Z}) : \mathbb{R}^{9+|\mathbf{z}|+|\mathcal{Z}|} \mapsto \mathbb{R}^3$. Here, $\mathbf{n}$ denotes a normal vector in the observation space. $\mathbf{n}$ is computed by transforming the canonical normal vectors using the rotational part of forward transformations $\sum_{b=1}^{24} f_{\sigma_\omega}(\hat{\mathbf{x}}^{(i)})_b \mathbf{B}_b$ (Eq. (1)). $\mathbf{v}$ denotes viewing direction. Similar to [76, 86, 87], $\mathbf{z}$ denotes an SDF feature which is extracted from the output of the second-last layer of the neural SDF. $\mathcal{Z}$ denotes a per-frame latent code which is shared with the SDF network. It compensates for time-dependent dynamic lighting effects. The outputs of $f_{\sigma_c}$ are RGB color values in the range $[0,1]$.

---

[4] with slight abuse of notation, we also use $\hat{\mathbf{x}}$ to represent points in homogeneous coordinates when necessary.

### 3.3   Joint Root-Finding

While surface rendering [51, 87] could be used to learn the network parameters introduced in Sections 3.1 and 3.2, it cannot handle abrupt changes in depth, as demonstrated in [76]. We also observe severe geometric artifacts when applying surface rendering to our setup, we refer readers to Appendix F for such an ablation. On the other hand, volume rendering can better handle abrupt depth changes in articulated human rendering. However, volume rendering requires multi-step dense sampling on camera rays [76,86], which, when combined naively with the iterative root-finding algorithm [12], requires significantly more memory and becomes prohibitively slow to train and test. We thus employ a hybrid method similar to [53]. We first search the ray-surface intersection and then sample near/far surface points on the ray. In practice, we initialize our SDF network with [78]. Thus, we fix the sampling depth interval around the surface to $[-5cm, +5cm]$.

A naive way of finding the ray-surface intersection is to use sphere tracing [20] and map each point to canonical space via root-finding [12]. In this case, we need to solve the costly root-finding problem during each step of the sphere tracing. This becomes prohibitively expensive when the number of rays is large. Thus, we propose an alternative solution. We leverage the skinning weights of the nearest neighbor on the registered SMPL mesh to the query point $\bar{\mathbf{x}}$ and use the inverse of the linearly combined forward bone transforms to map $\bar{\mathbf{x}}$ to its rough canonical correspondence. Combining this approximate backward mapping with sphere tracing, we obtain rough estimations of intersection points. Then, starting from these rough estimations, we apply a novel joint root-finding algorithm to search the precise intersection points and their correspondences in canonical space. In practice, we found that using a single initialization for our joint root-finding works well already. Adding more initializations incurs drastic memory and runtime overhead while not achieving any noticeable improvements. We hypothesize that this is due to the fact that our initialization is obtained using inverse transformations with SMPL skinning weights rather than rigid bone transformations (as was done in [12]).

Formally, we define a camera ray as $\mathbf{r} = (\mathbf{c}, \mathbf{v})$ where $\mathbf{c}$ is the camera center and $\mathbf{v}$ is a unit vector that defines the direction of this camera ray. Any point on the camera ray can be expressed as $\mathbf{c} + \mathbf{v} \cdot d$ with $d >= 0$. The joint root-finding aims to find canonical point $\hat{\mathbf{x}}$ and depth $d$ on the ray in observation space, such that:

$$f_{\sigma_f}(\hat{\mathbf{x}}) = 0$$
$$LBS_{\sigma_\omega}(\hat{\mathbf{x}}) - (\mathbf{c} + \mathbf{v} \cdot d) = \mathbf{0} \tag{2}$$

in which $\mathbf{c}, \mathbf{v}$ are constants per ray. Denoting the joint vector-valued function as $g_{\sigma_f, \sigma_\omega}(\hat{\mathbf{x}}, d)$ and the joint root-finding problem as:

$$g_{\sigma_f, \sigma_\omega}(\hat{\mathbf{x}}, d) = \begin{bmatrix} f_{\sigma_f}(\hat{\mathbf{x}}) \\ LBS_{\sigma_\omega}(\hat{\mathbf{x}}) - (\mathbf{c} + \mathbf{v} \cdot d) \end{bmatrix} = \mathbf{0} \tag{3}$$

we can then solve it via Newton's method

$$\begin{bmatrix} \hat{\mathbf{x}}_{k+1} \\ d_{k+1} \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{x}}_k \\ d_k \end{bmatrix} - \mathbf{J}_k^{-1} \cdot g_{\sigma_f, \sigma_\omega}(\hat{\mathbf{x}}_k, d_k) \tag{4}$$

where:

$$\mathbf{J}_k = \begin{bmatrix} \frac{\partial f_{\sigma_f}}{\partial \hat{\mathbf{x}}}(\hat{\mathbf{x}}_k) & 0 \\ \frac{\partial LBS_{\sigma_\omega}}{\partial \hat{\mathbf{x}}}(\hat{\mathbf{x}}_k) & -\mathbf{v} \end{bmatrix} \tag{5}$$

Following [12], we use Broyden's method to avoid computing $\mathbf{J}_k$ at each iteration.

**Amortized Complexity:** Given the number of sphere-tracing steps as N and the number of root-finding steps as M, the amortized complexity for joint root-finding is $O(M)$ while naive alternation between sphere-tracing and root-finding is $O(MN)$. In practice, this results in about 5× speed up of joint root-finding compared to the naive alternation between sphere-tracing and root-finding. We also note that from a theoretical perspective, our proposed joint root-finding converges quadratically while the secant-method-based root-finding in the concurrent work [92] converges only superlinearly.

We describe how to compute implicit gradients wrt. the canonical SDF and the forward LBS in Appendix C. In the main paper, we use volume rendering which does not need to compute implicit gradients wrt. the canonical SDF.

### 3.4   Differentiable Volume Rendering

We employ a recently proposed SDF-based volume rendering formulation [86]. Specifically, we convert SDF values into density values $\sigma$ using the scaled CDF of the Laplace distribution with the negated SDF values as input

$$\sigma(\hat{\mathbf{x}}) = \frac{1}{b}\left(\frac{1}{2} + \frac{1}{2}\text{sign}(-f_{\sigma_f}(\hat{\mathbf{x}}))\left(1 - \exp(-\frac{|-f_{\sigma_f}(\hat{\mathbf{x}})|}{b})\right)\right) \tag{6}$$

where $b$ is a learnable parameter. Given the surface point found via solving Eq. (3), we sample 16 points around the surface points and another 16 points between the near scene bound and the surface point, and map them to canonical space along with the surface point. For rays that do not intersect with any surface, we uniformly sample 64 points for volume rendering. With $N$ sampled points on a ray $\mathbf{r} = (\mathbf{c}, \mathbf{v})$, we use standard volume rendering [49] to render the pixel color

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^{N} T^{(i)}\left(1 - \exp(-\sigma(\hat{\mathbf{x}}^{(i)})\delta^{(i)})\right) f_{c_\sigma}(\hat{\mathbf{x}}^{(i)}, \mathbf{n}^{(i)}, \mathbf{v}, \mathbf{z}, \mathcal{Z}) \tag{7}$$

$$T^{(i)} = \exp\left(-\sum_{j<i} \sigma(\hat{\mathbf{x}}^{(j)})\delta^{(j)}\right) \tag{8}$$

where $\delta^{(i)} = |d^{(i+1)} - d^{(i)}|$.

### 3.5   Loss Function

Our loss consists of a photometric loss in observation space and multiple regularizers in canonical space

$$\mathcal{L} = \lambda_C \cdot \mathcal{L}_C + \lambda_E \cdot \mathcal{L}_E + \lambda_O \cdot \mathcal{L}_O + \lambda_I \cdot \mathcal{L}_I + \lambda_S \cdot \mathcal{L}_S \qquad (9)$$

$\mathcal{L}_C$ is the L1 loss for color predictions. $\mathcal{L}_E$ is the Eikonal regularization [17]. $\mathcal{L}_O$ is an off-surface point loss, encouraging points far away from the SMPL mesh to have positive SDF values. Similarly, $\mathcal{L}_I$ regularizes points inside the canonical SMPL mesh to have negative SDF values. $\mathcal{L}_S$ encourages the forward LBS network to predict similar skinning weights to the canonical SMPL mesh. Different from [27, 81, 87], we do not use an explicit silhouette loss. Instead, we utilize foreground masks and set all background pixel values to zero. In practice, this encourages the SDF network to predict positive SDF values for points on rays that do not intersect with foreground masks. For detailed definitions of loss terms and model architectures, please refer to Appendix A, B.

## 4   Experiments

We validate the generalization ability and reconstruction quality of our proposed method against several recent baselines [58, 60, 72]. As was done in [60], we consider a setup with 4 cameras positioned equally spaced around the human subject. For an ablation study on different design choices of our model, including ray sampling strategy, LBS networks, and number of initializations for root-finding, we refer readers to Appendix F.

**Datasets:** We use the ZJU-MoCap [60] dataset as our primary testbed because its setup includes 23 cameras which allows us to extract pseudo-ground-truth geometry to evaluate our model. More specifically, the dataset consists of 9 sequences captured with 23 calibrated cameras. We use the training/testing splits from Neural Body [60] for both the cameras and the poses. As one of our goals is learn to detailed geometry, we collect pseudo-ground-truth geometry for the training poses. We use all 23 cameras and apply NeuS with a background NeRF model [76], a state-of-the-art method for multi-view reconstruction. Note that we refrain from using the masks provided by Neural Body [60] as these masks are noisy and insufficient for accurate static scene reconstruction. We observe that geometry reconstruction with NeuS [76] fails when subjects wear black clothes or the environmental light is not bright enough. Therefore, we manually exclude bad reconstructions and discard sequences with less than 3 valid reconstructions. For completeness, we also tested our approach on the H36M dataset [26] and report a quantitative comparison to [52, 58] in Appendix G.

**Baselines:** We compare against three major baselines: Neural Body [60](NB), Ani-NeRF [58](AniN), and A-NeRF [72](AN). Neural Body diffuses per-SMPL-vertex latent codes into observation space as additional conditioning for NeRF models to achieve state-of-the-art novel view synthesis results on training poses.

Ani-NeRF learns a canonical NeRF model and a backward LBS network which predicts residuals to the deterministic SMPL-based backward LBS. Consequently, the LBS network needs to be re-trained for each test sequence. A-NeRF employs a deterministic backward mapping with bone-relative embeddings for query points and only uses keypoints and joint rotations instead of surface models (*i.e.* SMPL surface). For the detailed setups of these baselines, please refer to Appendix E.

**Benchmark Tasks:** We benchmark our approach on three tasks: generalization to unseen poses, geometry reconstruction, and novel-view synthesis. To analyze generalization ability, we evaluate the trained models on unseen testing poses. Due to the stochastic nature of cloth deformations, we quantify performance via perceptual similarity to the ground-truth images with the LPIPS [91] metric. We report PSNR and SSIM in Appendix G. We also encourage readers to check out qualitative comparison videos at https://neuralbodies.github.io/arah/.

For geometry reconstruction, we evaluate our method and baselines on the training poses. We report point-based L2 Chamfer distance (CD) and normal consistency (NC) wrt. the pseudo-ground-truth geometry. During the evaluation, we only keep the largest connected component of the reconstructed meshes. Note that is in favor of the baselines as they are more prone to producing floating blob artifacts. We also remove any ground-truth or predicted mesh points that are below an estimated ground plane to exclude outliers from the ground plane from the evaluation. For completeness, we also evaluate novel-view synthesis with PSNR, SSIM, and LPIPS using the poses from the training split.

Table 1: **Generalization to Unseen Poses**. We report LPIPS [91] on synthesized images under unseen poses from the testset of the ZJU-MoCap dataset [60] (*i.e.* all views except 0, 6, 12, and 18). Our approach consistently outperforms the baselines by a large margin. We report PSNR and SSIM Appendix G.

| Sequence | Metric | NB | AniN | AN | Ours |
|----------|--------|------|------|------|------|
| 313 | LPIPS ↓ | 0.126 | 0.115 | 0.209 | **0.092** |
| 315 | LPIPS ↓ | 0.152 | 0.167 | 0.232 | **0.105** |
| 377 | LPIPS ↓ | 0.119 | 0.153 | 0.165 | **0.093** |
| 386 | LPIPS ↓ | 0.171 | 0.187 | 0.241 | **0.127** |
| 387 | LPIPS ↓ | 0.135 | 0.145 | 0.162 | **0.099** |
| 390 | LPIPS ↓ | 0.163 | 0.173 | 0.226 | **0.126** |
| 392 | LPIPS ↓ | 0.135 | 0.169 | 0.183 | **0.106** |
| 393 | LPIPS ↓ | 0.132 | 0.155 | 0.175 | **0.104** |
| 394 | LPIPS ↓ | 0.150 | 0.171 | 0.199 | **0.111** |

Table 2: **Geometry Reconstruction**. We report L2 Chamfer Distance (CD) and Normal Consistency (NC) on the training poses of the ZJU-MoCap dataset [60]. Note that AniN and AN occasionally produce large background blobs that are connected to the body resulting in large deviations from the ground truth.

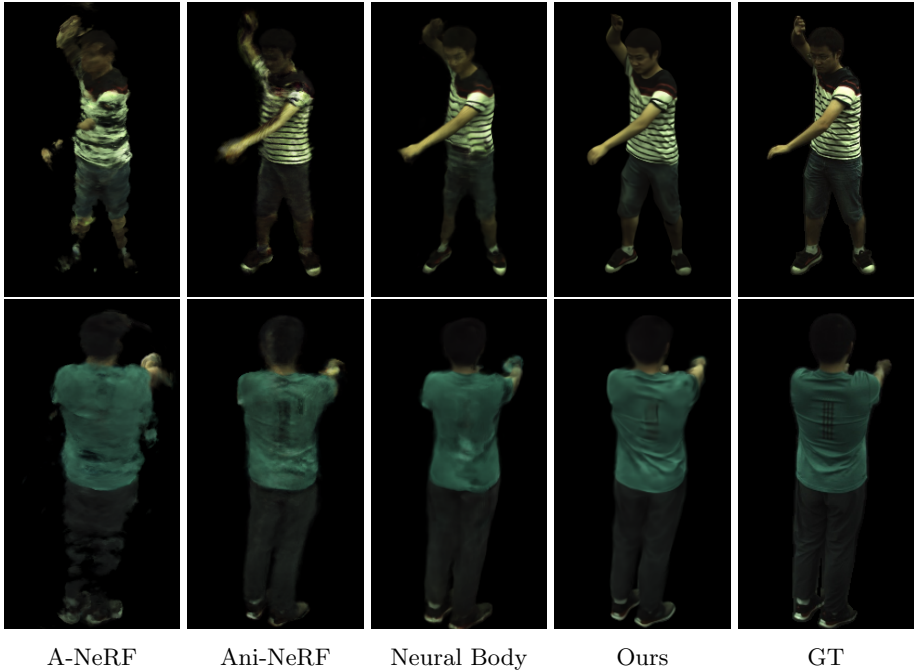| Sequence | Metric | NB | AniN | AN | Ours |
|----------|--------|-------|-------|-------|-------|
| 313 | CD ↓ | 1.258 | 1.242 | 9.174 | **0.707** |
| | NC ↑ | 0.700 | 0.599 | 0.691 | **0.809** |
| 315 | CD ↓ | 2.167 | 2.860 | 1.524 | **0.779** |
| | NC ↑ | 0.636 | 0.450 | 0.610 | **0.753** |
| 377 | CD ↓ | 1.062 | 1.649 | 1.008 | **0.840** |
| | NC ↑ | 0.672 | 0.541 | 0.682 | **0.786** |
| 386 | CD ↓ | 2.938 | 23.53 | 3.632 | **2.880** |
| | NC ↑ | 0.607 | 0.325 | 0.596 | **0.741** |
| 393 | CD ↓ | 1.753 | 3.252 | 1.696 | **1.342** |
| | NC ↑ | 0.600 | 0.481 | 0.605 | **0.739** |
| 394 | CD ↓ | 1.510 | 2.813 | 558.8 | **1.177** |
| | NC ↑ | 0.628 | 0.540 | 0.639 | **0.762** |

Fig. 3: **Generalization to Unseen Poses** on the testing poses of ZJU-MoCap. A-NeRF struggles with unseen poses due to the limited training poses and the lack of a SMPL surface prior. Ani-NeRF produces noisy images as it uses an inaccurate backward mapping function. Neural Body loses details, e.g. wrinkles, because its conditional NeRF is learned in observation space. Our approach generalizes well to unseen poses and can model fine details like wrinkles.

## 4.1   Generalization to Unseen Poses

We first analyze the generalization ability of our approach in comparison to the baselines. Given a trained model and a pose from the test set, we render images of the human subject in the given pose. We show qualitative results in Fig. 3 and quantitative results in Table 1. We significantly outperform the baselines both qualitatively and quantitatively. The training poses of the ZJU-MoCap dataset are extremely limited, usually comprising just 60-300 frames of repetitive motion. This limited training data results in severe overfitting for the baselines. In contrast, our method generalizes well to unseen poses, even when training data is limited.

We additionally animate our models trained on the ZJU-MoCap dataset using extreme out-of-distribution poses from the AMASS [43] and AIST++ [34] datasets. As shown in Fig. 5, even under extreme pose variation our approach produces plausible geometry and rendering results while all baselines show severe artifacts. We attribute the large improvement on unseen poses to our root-
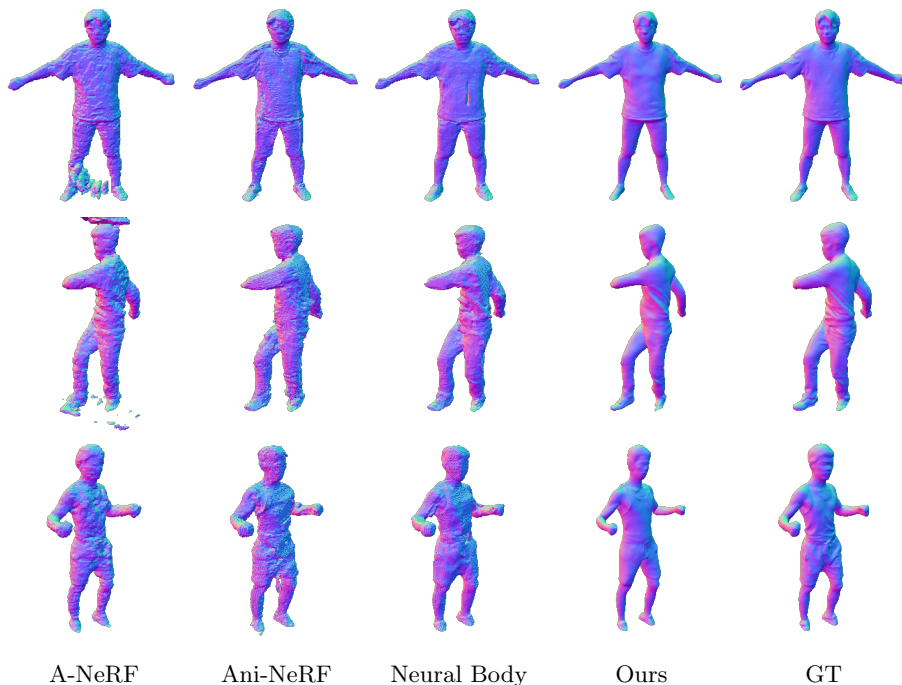
| A-NeRF | Ani-NeRF | Neural Body | Ours | GT |

Fig. 4: **Geometry Reconstruction**. Our approach reconstructs more fine-grained geometry than the baselines while preserving high-frequency details such as wrinkles. Note that we remove an estimated ground plane from all meshes.

finding-based backward skinning, as the learned forward skinning weights are constants per subject, while root-finding is a deterministic optimization process that does not rely on learned neural networks that condition on inputs from the observation space. More comparisons can be found in Appendix H.2, H.3.

## 4.2    Geometry Reconstruction on Training Poses

Next, we analyze the geometry reconstructed with our approach against reconstructions from the baselines. We compare to the pseudo-ground-truth obtained from NeuS [76]. We show qualitative results in Fig. 4 and quantitative results in Table 2. Our approach consistently outperforms existing NeRF-based human models on geometry reconstruction. As evidenced in Fig. 4, the geometry obtained with our approach is much cleaner compared to NeRF-based baselines, while preserving high-frequency details such as wrinkles.

## 4.3    Novel View Synthesis on Training Poses

Lastly, we analyze our approach for novel view synthesis on training poses. Table. 3 provides a quantitative comparison to the baselines. While not the main

Table 3: **Novel View Synthesis.** We report PSNR, SSIM, and LPIPS [91] for novel views of training poses of the ZJU-MoCap dataset [60]. Due to better geometry, our approach produces more consistent rendering results across novel views than the baselines. We include qualitative comparisons in Appendix H.1. Note that we crop slightly larger bounding boxes than Neural Body [60] to better capture loose clothes, *e.g.* sequence 387 and 390. Therefore, the reported numbers vary slightly from their evaluation.

| | 313 | | | 315 | | | 377 | | |
| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|---|---|---|---|---|---|
| NB | 30.5 | 0.967 | 0.068 | 26.4 | 0.958 | 0.079 | **28.1** | **0.956** | 0.080 |
| Ani-N | 29.8 | 0.963 | 0.075 | 23.1 | 0.917 | 0.138 | 24.2 | 0.925 | 0.124 |
| A-NeRF | 29.2 | 0.954 | 0.075 | 25.1 | 0.948 | 0.087 | 27.2 | 0.951 | 0.080 |
| Ours | **31.6** | **0.973** | **0.050** | **27.0** | **0.965** | **0.058** | 27.8 | **0.956** | **0.071** |
| | 386 | | | 387 | | | 390 | | |
| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| NB | 29.0 | **0.935** | 0.112 | 26.7 | 0.942 | 0.101 | **27.9** | 0.928 | 0.112 |
| Ani-N | 25.6 | 0.878 | 0.199 | 25.4 | 0.926 | 0.131 | 26.0 | 0.912 | 0.148 |
| A-NeRF | 28.5 | 0.928 | 0.127 | 26.3 | 0.937 | 0.100 | 27.0 | 0.914 | 0.126 |
| Ours | **29.2** | 0.934 | **0.105** | **27.0** | **0.945** | **0.079** | **27.9** | **0.929** | **0.102** |
| | 392 | | | 393 | | | 394 | | |
| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| NB | **29.7** | **0.949** | 0.101 | **27.7** | 0.939 | 0.105 | 28.7 | 0.942 | 0.098 |
| Ani-N | 28.0 | 0.931 | 0.151 | 26.1 | 0.916 | 0.151 | 27.5 | 0.924 | 0.142 |
| A-NeRF | 28.7 | 0.942 | 0.106 | 26.8 | 0.931 | 0.113 | 28.1 | 0.936 | 0.103 |
| Ours | 29.5 | 0.948 | **0.090** | **27.7** | **0.940** | **0.093** | **28.9** | **0.945** | **0.084** |

focus of this work, our approach also outperforms existing methods on novel view synthesis. This suggests that more faithful modeling of geometry is also beneficial for the visual fidelity of novel views. Particularly when few training views are available, NeRF-based methods produce blob/cloud artifacts. By removing such artifacts, our approach achieves high image fidelity and better consistency across novel views. Due to space limitations, we include further qualitative results on novel view synthesis in Appendix H.1.

## 5   Conclusion

We propose a new approach to create animatable avatars from sparse multi-view videos. We largely improve geometry reconstruction over existing approaches by modeling the geometry as articulated SDFs. Further, our novel joint root-finding algorithm enables generalization to extreme out-of-distribution poses. We discuss limitations of our approach in Appendix I.
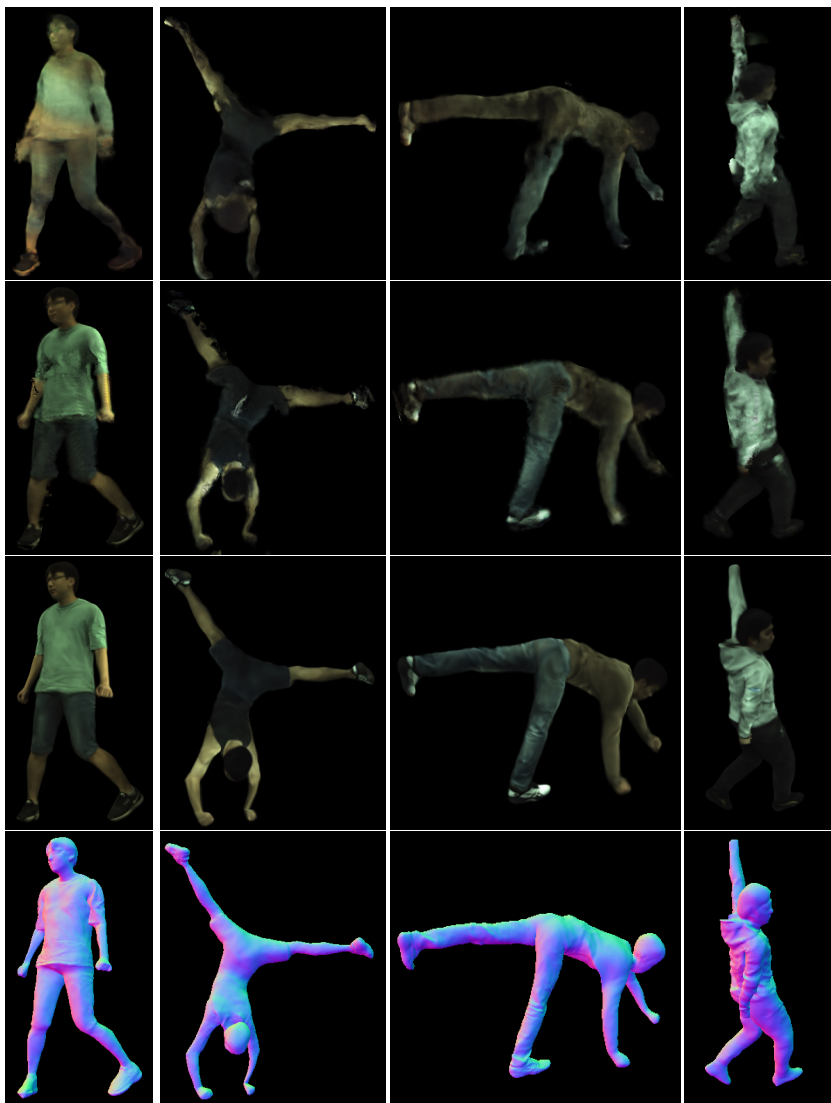
Fig. 5: **Qualitative Results on Out-of-distribution Poses** from the AMASS [43] and AIST++ [34] datasets. From top to bottom row: Neural Body, Ani-NeRF, our rendering, and our geometry. Note that Ani-NeRF requires re-training their backward LBS network on novel pose sequence. We did not show A-NeRF results as it already produces severe overfitting effects on ZJU-MoCap test poses. For more qualitative comparisons, please refer to Appendix H.3.

# References

1. Alldieck, T., Magnor, M., Bhatnagar, B.L., Theobalt, C., Pons-Moll, G.: Learning to reconstruct people in clothing from a single RGB camera. In: Proc. of CVPR (2019) 2, 3
2. Alldieck, T., Magnor, M., Xu, W., Theobalt, C., Pons-Moll, G.: Video based reconstruction of 3d people models. In: Proc. of CVPR (2018) 2, 3, 23
3. Alldieck, T., Xu, H., Sminchisescu, C.: imghum: Implicit generative models of 3d human shape and articulated pose. In: Proc. of CVPR (2021) 4
4. Alldieck, T., Zanfir, M., Sminchisescu, C.: Photorealistic monocular 3d reconstruction of humans wearing clothing. In: In Proc. of CVPR (2022) 3
5. Anguelov, D., Srinivasan, P., Koller, D., Thrun, S., Rodgers, J., Davis, J.: Scape: shape completion and animation of people. ACM Transasctions Graphics **24** (2005) 3, 6
6. Bhatnagar, B.L., Sminchisescu, C., Theobalt, C., Pons-Moll, G.: Combining implicit function learning and parametric models for 3d human reconstruction. In: Proc. of ECCV (2020) 2, 3
7. Bhatnagar, B.L., Sminchisescu, C., Theobalt, C., Pons-Moll, G.: Loopreg: Self-supervised learning of implicit surface correspondences, pose and shape for 3d human mesh registration. In: Proc. of NeurIPS (2020) 2, 3
8. Burov, A., Nießner, M., Thies, J.: Dynamic surface function networks for clothed human bodies. In: Proc. of ICCV (2021) 3
9. Chan, E., Monteiro, M., Kellnhofer, P., Wu, J., Wetzstein, G.: pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In: Proc. of CVPR (2021) 21
10. Chen, J., Zhang, Y., Kang, D., Zhe, X., Bao, L., Jia, X., Lu, H.: Animatable neural radiance fields from monocular rgb videos. arXiv preprint arXiv:2106.13629 (2021) 4
11. Chen, X., Jiang, T., Song, J., Yang, J., Black, M.J., Geiger, A., Hilliges, O.: gdna: Towards generative detailed neural avatars. In: In Proc. of CVPR (2022) 3
12. Chen, X., Zheng, Y., Black, M., Hilliges, O., Geiger, A.: Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. In: Proc. of ICCV (2021) 3, 6, 7, 8, 21, 22, 23, 26, 30
13. Chen, Z., Zhang, H.: Learning implicit fields for generative shape modeling. In: Proc. of CVPR (2019) 3
14. Chibane, J., Alldieck, T., Pons-Moll, G.: Implicit functions in feature space for 3d shape reconstruction and completion. In: Proc. of CVPR (2020) 3
15. Corona, E., Pumarola, A., Alenyà, G., Pons-Moll, G., Moreno-Noguer, F.: Smplicit: Topology-aware generative model for clothed people. In: In Proc. of CVPR (2021) 3
16. Dong, Z., Guo, C., Song, J., Chen, X., Geiger, A., Hilliges, O.: Pina: Learning a personalized implicit neural avatar from a single rgb-d video sequence. In: In Proc. of CVPR (2022) 3
17. Gropp, A., Yariv, L., Haim, N., Atzmon, M., Lipman, Y.: Implicit geometric regularization for learning shapes. In: Proc. of ICML (2020) 9, 20
18. Guan, P., Reiss, L., Hirshberg, D.A., Weiss, E., Black, M.J.: Drape: Dressing any person. ACM Transasctions Graphics **31**(4) (2012) 2, 3
19. Gundogdu, E., Constantin, V., Seifoddini, A., Dang, M., Salzmann, M., Fua, P.: Garnet: A two-stream network for fast and accurate 3d cloth draping. In: Proc. of ICCV (2019) 2, 3

20. Hart, J.C.: Sphere tracing: A geometric method for the antialiased ray tracing of implicit surfaces. The Visual Computer **12**(10) (1995) 7
21. Hasler, N., Stoll, C., Sunkel, M., Rosenhahn, B., Seidel, H.P.: A Statistical Model of Human Pose and Body Shape. Computer Graphics Forum **28**, 337–346 (2009) 3, 6
22. He, T., Collomosse, J., Jin, H., Soatto, S.: Geo-pifu: Geometry and pixel aligned implicit functions for single-view human reconstruction. In: Proc. of NeurIPS (2020) 3
23. He, T., Xu, Y., Saito, S., Soatto, S., Tung, T.: Arch++: Animation-ready clothed human reconstruction revisited. In: In Proc. of ICCV (2021) 3
24. Hu, T., Yu, T., Zheng, Z., Zhang, H., Liu, Y., Zwicker, M.: Hvtr: Hybrid volumetric-textural rendering for human avatars. arXiv preprint arXiv:2112.10203 (2021) 4
25. Huang, Z., Xu, Y., Lassner, C., Li, H., Tung, T.: ARCH: Animatable Reconstruction of Clothed Humans. In: Proc. of CVPR (2020) 2, 3
26. 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) 3, 9, 27
27. Jiang, B., Hong, Y., Bao, H., Zhang, J.: Selfrecon: Self reconstruction your digital avatar from monocular video. In: In Proc. of CVPR (2022) 4, 9
28. Joo, H., Simon, T., Sheikh, Y.: Total capture: A 3d deformation model for tracking faces, hands, and bodies. In: Proc. of CVPR (2018) 3
29. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: Proc. of ICLR (2015) 21, 24
30. Kwon, Y., Kim, D., Ceylan, D., Fuchs, H.: Neural human performer: Learning generalizable radiance fields for human performance rendering. In: Proc. of NeurIPS (2021) 3
31. Lähner, Z., Cremers, D., Tung, T.: Deepwrinkles: Accurate and realistic clothing modeling. In: Proc. of ECCV (2018) 2, 3
32. Li, R., Tanke, J., Vo, M., Zollhoefer, M., Gall, J., Kanazawa, A., Lassner, C.: Tava: Template-free animatable volumetric actors. In: In Proc. of ECCV (2022) 4
33. Li, R., Xiu, Y., Saito, S., Huang, Z., Olszewski, K., Li, H.: Monocular real-time volumetric performance capture. In: Proc. of ECCV (2020) 3
34. Li, R., Yang, S., Ross, D.A., Kanazawa, A.: Ai choreographer: Music conditioned 3d dance generation with aist++. In: Proc. of ICCV (2021) 11, 14, 29, 30, 34
35. Li, Y., Habermann, M., Thomaszewski, B., Coros, S., Beeler, T., Theobalt, C.: Deep physics-aware inference of cloth deformation for monocular human performance capture. In: Proc. of 3DV (2021) 2, 3
36. Li, Z., Yu, T., Pan, C., Zheng, Z., Liu, Y.: Robust 3d self-portraits in seconds. In: Proc. of CVPR (2020) 2
37. Li, Z., Yu, T., Zheng, Z., Guo, K., Liu, Y.: Posefusion: Pose-guided selective fusion for single-view human volumetric capture. In: Proc. of CVPR (2021) 2
38. 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) 3
39. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: SMPL: A skinned multi-person linear model. ACM Transasctions Graphics **34**(6) (2015) 3, 6
40. Ma, Q., Saito, S., Yang, J., Tang, S., Black, M.J.: SCALE: Modeling clothed humans with a surface codec of articulated local elements. In: Proc. of CVPR (2021) 3

41. Ma, Q., Yang, J., Ranjan, A., Pujades, S., Pons-Moll, G., Tang, S., Black, M.J.: Learning to dress 3D people in generative clothing. In: Proc. of CVPR (2020) 21
42. Ma, Q., Yang, J., Tang, S., Black, M.J.: The power of points for modeling humans in clothing. In: Proc. of ICCV (2021) 3
43. Mahmood, N., Ghorbani, N., Troje, N.F., Pons-Moll, G., Black, M.J.: AMASS: Archive of motion capture as surface shapes. In: Proc. of ICCV (2019) 11, 14
44. Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., Geiger, A.: Occupancy networks: Learning 3d reconstruction in function space. In: Proc. of CVPR (2019) 3
45. Michalkiewicz, M., Pontes, J.K., Jack, D., Baktashmotlagh, M., Eriksson, A.: Implicit surface representations as layers in neural networks. In: Proc. of ICCV (2019) 3
46. Mihajlovic, M., Bansal, A., Zollhoefer, M., Tang, S., Saito, S.: KeypointNeRF: Generalizing image-based volumetric avatars using relative spatial encoding of keypoints. In: In Proc. of ECCV (2022) 4
47. Mihajlovic, M., Saito, S., Bansal, A., Zollhoefer, M., Tang, S.: COAP: Compositional articulated occupancy of people. In: In Proc. of CVPR (2022) 3
48. Mihajlovic, M., Zhang, Y., Black, M.J., Tang, S.: LEAP: Learning articulated occupancy of people. In: Proc. of CVPR (2021) 3, 6
49. 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: Proc. of ECCV (2020) 2, 4, 8, 26
50. Nichol, A., Achiam, J., Schulman, J.: On first-order meta-learning algorithms. arXiv preprint arXiv:1803.02999 (2018) 21
51. Niemeyer, M., Mescheder, L., Oechsle, M., Geiger, A.: Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In: Proc. of CVPR (2020) 7
52. Noguchi, A., Sun, X., Lin, S., Harada, T.: Neural articulated radiance field. In: Proc. of ICCV (2021) 3, 4, 9, 28
53. Oechsle, M., Peng, S., Geiger, A.: Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In: Proc. of ICCV (2021) 4, 7
54. Osman, A.A.A., Bolkart, T., Black, M.J.: Star: Sparse trained articulated human body regressor. In: Proc. of ECCV (2020) 3, 6
55. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: Deepsdf: Learning continuous signed distance functions for shape representation. In: Proc. of CVPR (2019) 3
56. Patel, C., Liao, Z., Pons-Moll, G.: Tailornet: Predicting clothing in 3d as a function of human pose, shape and garment style. In: Proc. of CVPR (2020) 2, 3
57. Pavlakos, G., Choutas, V., Ghorbani, N., Bolkart, T., Osman, A.A.A., Tzionas, D., Black, M.J.: Expressive body capture: 3d hands, face, and body from a single image. In: Proc. of CVPR (2019) 3, 6
58. Peng, S., Dong, J., Wang, Q., Zhang, S., Shuai, Q., Zhou, X., Bao, H.: Animatable neural radiance fields for modeling dynamic human bodies. In: Proc. of ICCV (2021) 2, 3, 4, 9, 25, 27, 28
59. Peng, S., Zhang, S., Xu, Z., Geng, C., Jiang, B., Bao, H., Zhou, X.: Animatable neural implict surfaces for creating avatars from videos. arXiv preprint arXiv:2203.08133 (2022) 4
60. 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: Proc. of CVPR (2021) 2, 3, 4, 9, 10, 13, 21, 25, 27

61. Peng, S., Jiang, C.M., Liao, Y., Niemeyer, M., Pollefeys, M., Geiger, A.: Shape as points: A differentiable poisson solver. In: Proc. of NeurIPS (2021) 3
62. Perez, E., Strub, F., de Vries, H., Dumoulin, V., Courville, A.C.: Film: Visual reasoning with a general conditioning layer. In: Proc. of AAAI (2018) 21
63. Prokudin, S., Black, M.J., Romero, J.: SMPLpix: Neural avatars from 3D human models. In: Proc. WACV (2021) 4
64. Raj, A., Tanke, J., Hays, J., Vo, M., Stoll, C., Lassner, C.: Anr-articulated neural rendering for virtual avatars. In: Proc. of CVPR (2021) 3
65. Saito, S., , Huang, Z., Natsume, R., Morishima, S., Kanazawa, A., Li, H.: Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In: Proc. of ICCV (2019) 2, 3
66. Saito, S., Simon, T., Saragih, J., Joo, H.: Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In: Proc. of CVPR (2020) 2, 3, 5
67. Saito, S., Yang, J., Ma, Q., Black, M.J.: SCANimate: Weakly supervised learning of skinned clothed avatar networks. In: Proc. of CVPR (2021) 3, 6
68. Salimans, T., Kingma, D.P.: Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In: Proc. of NeurIPS (2016) 21
69. Santesteban, I., Thuerey, N., Otaduy, M.A., Casas, D.: Self-Supervised Collision Handling via Generative 3D Garment Models for Virtual Try-On. In: Proc. of CVPR (2021) 2, 3
70. Sitzmann, V., Martel, J.N., Bergman, A.W., Lindell, D.B., Wetzstein, G.: Implicit neural representations with periodic activation functions. In: Proc. of NeurIPS (2020) 21
71. Su, S.Y., Bagautdinov, T., Rhodin, H.: Danbo: Disentangled articulated neural body representations via graph neural networks. In: In Proc. of ECCV (2022) 4
72. Su, S.Y., Yu, F., Zollhoefer, M., Rhodin, H.: A-neRF: Articulated neural radiance fields for learning human shape, appearance, and pose. In: Proc. of NeurIPS (2021) 3, 4, 9, 25
73. Su, Z., Xu, L., Zheng, Z., Yu, T., Liu, Y., Fang, L.: Robustfusion: Human volumetric capture with data-driven visual cues using a rgbd camera. In: Proc. of ECCV (2020) 2
74. Tiwari, G., Sarafianos, N., Tung, T., Pons-Moll, G.: Neural-GIF: Neural generalized implicit functions for animating people in clothing. In: Proc. of ICCV (2021) 3
75. Tiwari, L., Bhowmick, B.: Deepdraper: Fast and accurate 3d garment draping over a 3d human body. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops (2021) 2, 3
76. Wang, P., Liu, L., Liu, Y., Theobalt, C., Komura, T., Wang, W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. In: Proc. of NeurIPS (2021) 4, 6, 7, 9, 12, 26
77. Wang, S., Geiger, A., Tang, S.: Locally aware piecewise transformation fields for 3d human mesh registration. In: In Proc. of CVPR (2021) 3
78. Wang, S., Mihajlovic, M., Ma, Q., Geiger, A., Tang, S.: Metaavatar: Learning animatable clothed human models from few depth images. In: Proc. of NeurIPS (2021) 3, 6, 7, 21, 22, 24, 25
79. Weng, C.Y., Curless, B., Srinivasan, P.P., Barron, J.T., Kemelmacher-Shlizerman, I.: Humannerf: Free-viewpoint rendering of moving people from monocular video. In: In Proc. of CVPR (2022) 4
80. Xiu, Y., Yang, J., Tzionas, D., Black, M.J.: ICON: Implicit Clothed humans Obtained from Normals. In: In Proc. of CVPR (2022) 3

81. Xu, H., Alldieck, T., Sminchisescu, C.: H-neRF: Neural radiance fields for rendering and temporal reconstruction of humans in motion. In: Proc. of NeurIPS (2021) 3, 4, 9, 25
82. Xu, H., Bazavan, E.G., Zanfir, A., Freeman, W.T., Sukthankar, R., Sminchisescu, C.: Ghum & ghuml: Generative 3d human shape and articulated pose models. In: Proc. of CVPR (2020) 3, 6
83. Xu, L., Su, Z., Han, L., Yu, T., Liu, Y., Fang, L.: Unstructuredfusion: Real-time 4d geometry and texture reconstruction using commercial rgbd cameras. IEEE Transactions on Pattern Analysis and Machine Intelligence **42** (2020) 2
84. Xu, T., Fujita, Y., Matsumoto, E.: Surface-aligned neural radiance fields for controllable 3d human synthesis. In: CVPR (2022) 4, 28
85. Yang, J., Franco, J.S., Hétroy-Wheeler, F., Wuhrer, S.: Analyzing clothing layer deformation statistics of 3d human motions. In: Proc. of ECCV (2018) 2, 3
86. Yariv, L., Gu, J., Kasten, Y., Lipman, Y.: Volume rendering of neural implicit surfaces. In: Proc. of NeurIPS (2021) 4, 6, 7, 8, 26
87. Yariv, L., Kasten, Y., Moran, D., Galun, M., Atzmon, M., Ronen, B., Lipman, Y.: Multiview neural surface reconstruction by disentangling geometry and appearance. In: Proc. of NeurIPS (2020) 4, 6, 7, 9, 20, 23, 24
88. Yu, T., Zheng, Z., Guo, K., Zhao, J., Dai, Q., Li, H., Pons-Moll, G., Liu, Y.: Doublefusion: Real-time capture of human performances with inner body shapes from a single depth sensor. In: Proc. of CVPR (2018) 2
89. Zakharkin, I., Mazur, K., Grigorev, A., Lempitsky, V.: Point-based modeling of human clothing. In: Proc. of ICCV (2021) 3
90. Zhang, C., Pujades, S., Black, M.J., Pons-Moll, G.: Detailed, accurate, human shape estimation from clothed 3d scan sequences. In: Proc. of CVPR (2017) 2, 3
91. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proc. of CVPR (2018) 10, 13
92. Zheng, Y., Abrevaya, V.F., Bühler, M.C., Chen, X., Black, M.J., Hilliges, O.: I M Avatar: Implicit morphable head avatars from videos. In: In Proc. of CVPR (2022) 4, 8
93. Zheng, Z., Yu, T., Liu, Y., Dai, Q.: Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. IEEE Transactions on Pattern Analysis and Machine Intelligence pp. 1–1 (2021). https://doi.org/10.1109/TPAMI.2021.3050505 3
94. Zuo, X., Wang, S., Sun, Q., Gong, M., Cheng, L.: Self-supervised 3d human mesh recovery from noisy point clouds. arXiv preprint arXiv:2107.07539 (2021) 3

# A    Loss Definition

In Section 3.5 of the main paper, we define the loss terms as follows

$$\mathcal{L} = \lambda_C \cdot \mathcal{L}_C + \lambda_E \cdot \mathcal{L}_E + \lambda_O \cdot \mathcal{L}_O + \lambda_I \cdot \mathcal{L}_I + \lambda_S \cdot \mathcal{L}_S \qquad (A.1)$$

In this section, we elaborate on how each loss term is defined. Let $I_p \in [0, 1]^3$ denote the ground-truth RGB value of a pixel $p$. Further, let $P$ denote the set of all pixels sampled from an image.

**RGB Color Loss:** The RGB color loss is defined as

$$\mathcal{L}_C = \frac{1}{|P|} \sum_{p \in P} \left| f_{\sigma_c}(\hat{\mathbf{x}}^{(p)}, \mathbf{n}^{(p)}, \mathbf{v}^{(p)}, \mathbf{z}, \mathcal{Z}) - I_p \right| \tag{A.2}$$

**Eikonal Regularization:** We sample 1024 points, denoted as $\hat{\mathbf{X}}_{\text{eik}}$, in the range $[-1, 1]^3$ in canonical space, and compute Eikonal loss [17] as follows:

$$\mathcal{L}_E = \frac{1}{|P|} \sum_{\hat{\mathbf{x}} \in \hat{\mathbf{X}}_{\text{eik}}} \left| \|\nabla_{\hat{\mathbf{x}}} f_{\sigma_f}(\hat{\mathbf{x}})\|_2 - 1 \right| \tag{A.3}$$

**Off-surface Point Loss:** In canonical space, we sample 1024 points whose distance to the canonical SMPL mesh is greater than 20cm. Let $\hat{\mathbf{X}}_{\text{off}}$ denote these sampled points, we compute the off-surface point loss as

$$\mathcal{L}_O = \frac{1}{|P|} \sum_{\hat{\mathbf{x}} \in \hat{\mathbf{X}}_{\text{off}}} \exp\left(-1e^2 \cdot f_{\sigma_f}(\hat{\mathbf{x}})\right) \tag{A.4}$$

**Inside Point Loss:** In canonical space, we sample 1024 points that are inside the canonical SMPL mesh and whose distance to the SMPL surface is greater than 1cm. Let $\hat{\mathbf{X}}_{\text{in}}$ denote these sampled points, we compute the inside point loss as

$$\mathcal{L}_I = \frac{1}{|P|} \sum_{\hat{\mathbf{x}} \in \hat{\mathbf{X}}_{\text{in}}} \text{sigmoid}\left(5e^3 \cdot f_{\sigma_f}(\hat{\mathbf{x}})\right) \tag{A.5}$$

**Skinning Loss:** Finally, in canonical space, we sample 1024 points on the canonical SMPL surface, $\hat{\mathbf{X}}_S$, and regularize the forward LBS network with the corresponding SMPL skinning weights $\mathbf{W} = \{\mathbf{w}\}$:

$$\mathcal{L}_S = \frac{1}{|P|} \sum_{\substack{\hat{\mathbf{x}} \in \hat{\mathbf{X}}_S \\ \mathbf{w} \in \mathbf{W}}} \sum_{i=1}^{i=24} \left| f_{\sigma_\omega}(\hat{\mathbf{x}})_i - \mathbf{w}_i \right| \tag{A.6}$$

We set $\lambda_C = 3e^1, \lambda_E = 5e^1, \lambda_O = 1e^2, \lambda_I = \lambda_S = 10$ throughout all experiments.

**Mask Loss:** As described in Section 3.5 of the main paper, our volume rendering formulation does not need explicit mask loss. Here we describe the mask loss from [87] which we use in the ablation study on surface rendering (Section F). Given the camera ray $\mathbf{r}^{(p)} = (\mathbf{c}, \mathbf{v}^{(p)})$ of a specific pixel $p$, we first define $S(\alpha, \mathbf{c}, \mathbf{v}^{(p)}) = \text{sigmoid}(-\alpha \min_{d \geq 0} f_{\sigma_f}(LBS_{\sigma_\omega}^{-1}(\mathbf{c} + d\mathbf{v}^{(p)})))$, *i.e.* the Sigmoid of the minimal SDF along a ray. In practice we sample 100 $ds$ uniformly between $[d_{\min}, d_{\max}]$ along the ray, where $d_{\min}$ and $d_{\max}$ are determined by the bounding box of the registered SMPL mesh. $\alpha$ is a learnable scalar parameter.

Let $O_p \in \{0, 1\}$ denote the foreground mask value (0 indicates background and 1 indicates foreground) of a pixel $p$. Further, let $P_{in}$ denote the set of pixels

for which ray-intersection with the iso-surface of neural SDF is found and $O_p = 1$, while $P_{out} = P \setminus P_{in}$ is the set of pixels for which no ray-intersection with the iso-surface of neural SDF is found or $O_p = 0$. The mask loss is defined as

$$\mathcal{L}_M = \frac{1}{\alpha|P|} \sum_{p \in P_{out}} \text{BCE}(O_p, S(\alpha, \mathbf{c}, \mathbf{v}^{(p)})))$$  (A.7)

where $\text{BCE}(\cdot)$ denotes binary cross entropy loss. We set the weight of $\mathcal{L}_M$ to be $3e^3$ and add this loss term to Eq. (A.1) for our surface rendering baseline in Section F.

# B    Network Architectures

In this section, we describe detailed network architectures for the forward LBS network $f_{\sigma_\omega}$, the SDF network $f_{\sigma_f}$ and the color network $f_{\sigma_c}$ introduced in Sections 3.1-3.2 of the main paper.

## B.1    Forward LBS Network

We use the same forward LBS network as [12], which consists of 4 hidden layers with 128 channels and weight normalization [68]. It uses Softplus activation with $\beta = 100$. $f_{\sigma_\omega}$ only takes query points in canonical space as inputs and does not have any conditional inputs.

To initialize this forward LBS network, we meta learn the network on skinning weights of canonical meshes from the CAPE [41] dataset. Specifically, we use Reptile [50] with 24 inner steps. The inner learning rate is set to $1e^{-4}$ while the outer learning rate is set to $1e^{-5}$. Adam [29] optimizer is used for both the inner and the outer loop. We train with a batch size of 4 for 100k steps of the outer loop. We use the resulting model as the initialization for our per-subject optimization on the ZJU-MoCap [60] dataset.

## B.2    Canonical SDF Network

We describe our canonical SDF network in Fig. B.1. The hypernetwork (top) and neural SDF (middle) are initialized with MetaAvatar [78] pre-trained on the CAPE dataset. Note that the SDF network from MetaAvatar can be trained with canonical meshes only and does not need any posed meshes as supervision. Each MLP of the hypernetwork (top) consists of one hidden layer with 256 channels and uses ReLU activation. The neural SDF (middle) consists of 5 hidden layers with 256 channels and uses a periodic activation [70]. In addition to the MetaAvatar SDF, we add a mapping network [9,62] which consists of 2 hidden layers with 256 channels and a ReLU activation. It maps the per-frame latent code $\mathcal{Z}$ to scaling factors and offsets that modulate the outputs from each layer of the neural SDF. We initialize the last layer of the mapping network to predict scaling factors with value 1 and offsets with value 0.
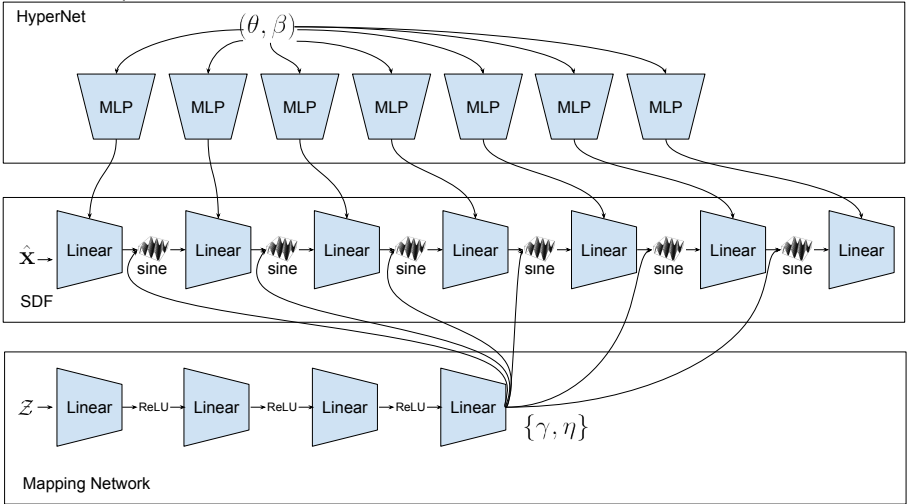
SDF Network $f_{\sigma_f}(\hat{\mathbf{x}}, \theta, \beta, \mathcal{Z})$



Fig. B.1: **Network Architecture for the SDF Network.** Our SDF network builds upon MetaAvatar [78] which uses a hypernetwork (top) that conditions on local body poses and shape $(\theta, \beta)$, and predicts the parameters of a neural SDF with periodic activation (middle). Since MetaAvatar does not model per-frame latent codes, we add a mapping network (bottom) that maps the per-frame latent code $\mathcal{Z}$ to scaling factors $\{\gamma\}$ and offsets $\{\eta\}$ which are used to modulate the outputs from each linear layer of the neural SDF, except for the last layer.

### B.3    Canonical Color Network

We describe our canonical color network in Fig. B.2. The network consists of 4 hidden layers with 256 channels and ReLU activation. The inputs to the network are also concatenated with activations of the third layer and fed into the fourth layer together.

## C    Implicit Gradients

In this section, we describe how to compute gradients of the root-finding solutions wrt. the forward LBS network and the SDF network. In the main paper, we use our novel joint root-finding algorithm to find the surface point and sample points around the surface point; these sampled points, along with the surface point, are mapped to canonical space via iterative root-finding [12]. Section C.1 describes how to differentiate through these points to compute gradients wrt. the forward LBS network. Section C.2 describes how to compute gradients wrt. the forward LBS network and the SDF network given the surface point and its correspondence. Section C.1 is used for volume rendering, which is described in

Color Network $f_{\sigma_c}(\hat{\mathbf{x}}, \mathbf{n}, \mathbf{v}, \mathbf{z}, \mathcal{Z})$
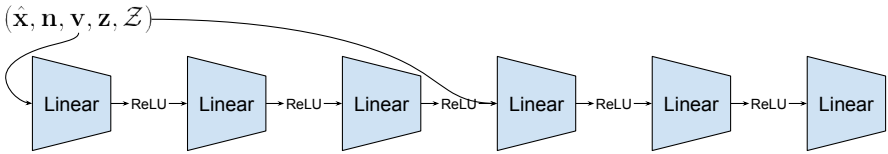
$(\hat{\mathbf{x}}, \mathbf{n}, \mathbf{v}, \mathbf{z}, \mathcal{Z})$ → Linear → ReLU → Linear → ReLU → Linear → ReLU → Linear → ReLU → Linear → ReLU → Linear

Fig. B.2: **Network Architecture for the Color Network.** The color network takes canonicalized query points $\hat{\mathbf{x}}$, normal vectors $\mathbf{n}$, viewing directions $\mathbf{v}$, an SDF feature $\mathbf{z}$, and a per-frame latent code $\mathcal{Z}$ as inputs.

Section 3.4 of the main paper. Section C.2 is used for surface rendering, which is one of our ablation baselines in Section F.

## C.1    Implicit Gradients for Forward LBS

Here we follow [12] and describe how to compute implicit gradients for the forward LBS network given samples on camera rays and their canonical correspondences. Denoting sampled points in observation space as $\bar{\mathbf{X}} = \{\bar{\mathbf{x}}\}_{i=1}^{N}$, and their canonical correspondences obtained by iterative root-finding [12] as $\hat{\mathbf{X}}^* = \{\hat{\mathbf{x}}^*\}_{i=1}^{N}$, they should satisfy the following condition

$$LBS_{\sigma_\omega}(\hat{\mathbf{x}}^{*(i)}) - \bar{\mathbf{x}}^{(i)} = 0, \quad \forall i = 1, \ldots, N \tag{C.1}$$

As done in [87], by applying implicit differentiation, we obtain a differentiable point sample $\hat{\mathbf{x}}$ as

$$\hat{\mathbf{x}} = \hat{\mathbf{x}}^* - (\mathbf{J}^*)^{-1} \cdot \left( LBS_{\sigma_\omega}(\hat{\mathbf{x}}^{*(i)}) - \bar{\mathbf{x}}^{(i)} \right) \tag{C.2}$$

where $\mathbf{J}^* = \frac{\partial LBS_{\sigma_\omega}}{\partial \hat{\mathbf{x}}}(\hat{\mathbf{x}}^*)$. $\hat{\mathbf{x}}^*$ and $\mathbf{J}^*$ are detached from the computational graph such that no gradient will flow through them. These differentiable samples can be used as inputs to the SDF and color networks. Gradients wrt. $\sigma_\omega$ are computed from photometric loss Eq. (A.2) via standard back-propagation. Taking the derivative wrt. $\sigma_\omega$ for both sides of Eq. (C.2) results in the same analytical gradient defined in Eq. (14) of [12].

**Pose and Shape Optimization:** We note that implicit gradients can also be back-propagated to SMPL parameters $\{\theta, \beta\}$ as the SMPL model is fully differentiable. We found pose and shape optimization particularly helpful when SMPL estimations are noisy, *e.g.* those estimated from monocular videos. In Fig. C.1 we show a qualitative sample on the People Snapshot dataset [2] where the pose is improved while the resulting model also achieves better visual quality.

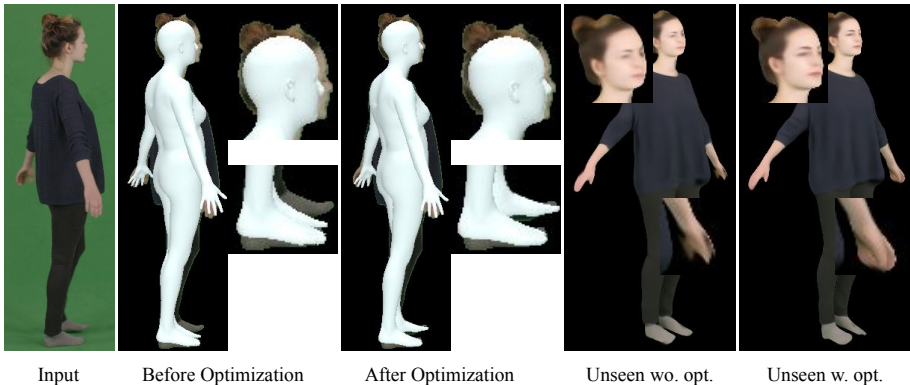| Input | Before Optimization | After Optimization | Unseen wo. opt. | Unseen w. opt. |

Fig. C.1: **Result of Pose and Shape Optimization.** We can improve the noisy SMPL estimations on training poses with implicit gradients and improve the rendering quality on unseen poses (see Unseen w. opt.).

## C.2    Implicit Gradients for Joint Root-finding

Now we derive implicit gradients for our joint root-finding algorithm. We denote the joint vector-valued function of the ray-surface intersection and forward LBS as $g_{\sigma_f, \sigma_\omega}(\hat{\mathbf{x}}, d)$. The joint root-finding problem is

$$g_{\sigma_f, \sigma_\omega}(\hat{\mathbf{x}}, d) = \begin{bmatrix} f_{\sigma_f}(\hat{\mathbf{x}}) \\ LBS_{\sigma_\omega}(\hat{\mathbf{x}}) - (\mathbf{c} + \mathbf{v} \cdot d) \end{bmatrix} = \mathbf{0} \tag{C.3}$$

with a slight abuse of notation, we denote the iso-surface point as $\hat{\mathbf{x}}^*$ and their corresponding depth in observation space as $d^*$. We follow [87] and use implicit differentiation to obtain a differentiable point sample $\hat{\mathbf{x}}$ and a depth sample $d$:

$$\begin{bmatrix} \hat{\mathbf{x}} \\ d \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{x}}^* \\ d^* \end{bmatrix} - (\mathbf{J}^*)^{-1} \cdot g_{\sigma_f, \sigma_\omega}(\hat{\mathbf{x}}^*, d^*) \tag{C.4}$$

where $\mathbf{J}^*$ is defined as

$$\mathbf{J}^* = \begin{bmatrix} \frac{\partial f_{\sigma_f}}{\partial \hat{\mathbf{x}}}(\hat{\mathbf{x}}^*) & 0 \\ \frac{\partial LBS_{\sigma_\omega}}{\partial \hat{\mathbf{x}}}(\hat{\mathbf{x}}^*) & -\mathbf{v} \end{bmatrix} \tag{C.5}$$

Similar to Section C.1, these differentiable samples can be used as inputs to the SDF and color networks and gradients wrt. $\sigma_f, \sigma_\omega$ can be computed from the photometric loss Eq. (A.2).

## D    Implementation Details

We use Adam [29] to optimize our models and the per-frame latent codes $\{\mathcal{Z}\}$. We initialize the SDF network with MetaAvatar [78] and set the learning rate

to $1e^{-6}$ as suggested in [78]. For the remaining models and the latent codes, we use a learning rate of $1e^{-4}$. We apply weight decay with a weight of 0.05 to the per-frame latent codes.

We train our models with a batch size of 4 and 2048 rays per batch, with 1024 rays sampled from the foreground mask and 1024 rays sampled from the background. As mentioned in Section 3.4 of the main paper, we sample 16 near and 16 far surface points for rays that intersect with a surface and 64 points for rays that do not intersect with a surface. Our model is trained for 250 epochs (except for sequence 313 which we trained for 1250 epochs, due to its training frames being much fewer than other sequences), which corresponds to 60k-80k iterations depending on the amount of training data. This takes about 1.5 days on 4 NVIDIA 2080 Ti GPUs. During training, we follow [81] and add normally distributed noise with zero mean and a standard deviation of 0.1 to the input $\theta$ of the SDF network. This noise ensures that the canonical SDF does not fail when given extreme out-of-distribution poses. We also augment the input viewing directions to the color network during training. We do so by randomly applying roll/pitch/yaw rotations sampled from a normal distribution with zero mean and a standard deviation of 45° to the viewing direction, but reject augmentation in which the angle between the estimated surface normal and the negated augmented viewing direction is greater than 90 degrees.

For inference, we follow [58, 60] and crop an enlarged bounding box around the projected SMPL mesh on the image plane and render only pixels inside the bounding box. For unseen test poses we follow the practice of [58, 60] and use the latent code $\mathcal{Z}$ of the last training frame as the input. The rendering time of a $512 \times 512$ image is about 10-20 seconds, depending on the bounding box size of the person. In this process, the proposed joint root-finding algorithm takes about 1 second.

# E    Implementation Details for Baselines

In this section, we describe the implementation details of the baselines from the main paper, *i.e.* Neural Body [60], Ani-NeRF [58], and A-NeRF [72].

## E.1    Neural Body

For quantitative evaluation, we use the official results provided by the Neural Body website. For generating rendering results and geometries, we use the official code of Neural Body and their pre-trained models without modification.

## E.2    Animatable NeRF (Ani-NeRF)

For quantitative evaluation, we use the official code and pre-trained models when possible, *i.e.* for sequences 313, 315, 377, and 386. For the remaining sequences that the official code does not provide pre-trained models, we train models using

the default hyperparameters that were applied to sequences 313, 315, 377, and 386.

We note that when reconstructing geometry on the training poses, Neural Body and Ani-NeRF compute visual hulls from ground-truth masks of training views and set density values outside the visual hulls to 0. This removes extraneous geometry blobs from reconstructions by Neural Body and Ani-NeRF. When testing on unseen poses, we disable the mask usage, as, by definition of the task, we do not have any image as input.

### E.3   A-NeRF

For A-NeRF, we follow the author's suggestions to 1) use a bigger foreground mask for ray sampling, 2) enable background estimation in the official code, and 3) use L2 loss instead of L1 loss. The learned models give reasonable novel view synthesis results on training poses (Fig. H.1) but cannot generalize to unseen poses (Fig. H.2). We hypothesize that this is because training poses on the ZJU-MoCap dataset are extremely limited, and A-NeRF uses only keypoints instead of surface models to construct their conditional inputs to NeRF networks. The lack of a surface model makes it easy for A-NeRF to confuse background and foreground, resulting in obvious floating blob artifacts. These artifacts are amplified when training poses are limited, making the generalization result of A-NeRF on the ZJU-MoCap dataset the worst among the baselines.

## F   Ablation Study

In this section, we ablate on ray sampling strategies as well as canonicalization strategies. We conduct an ablation on sequence 313. Metrics on all novel views of training poses are reported.

### F.1   Ablation on Ray Sampling Strategies

We compare our proposed ray sampling strategy to surface rendering and uniform sampling with 64 samples on the novel view synthesis task (Fig F.1). As discussed in the main paper, we did not use more sophisticated hierarchical sampling strategies [49,76,86] due to the computational cost of running the iterative root-finding [12] on dense samples and the memory cost for running additional forward/backward passes through the LBS network.

### F.2   Ablation on Learned forward LBS

In this subsection, we replace our learned forward LBS with (1) a backward LBS network that conditions on local body poses $\theta$, and (2) a deterministic LBS with nearest neighbor SMPL skinning weights. For the learned backward LBS, we always canonicalize the query points using the SMPL global translation and rotation before querying the LBS network. We also sample points on
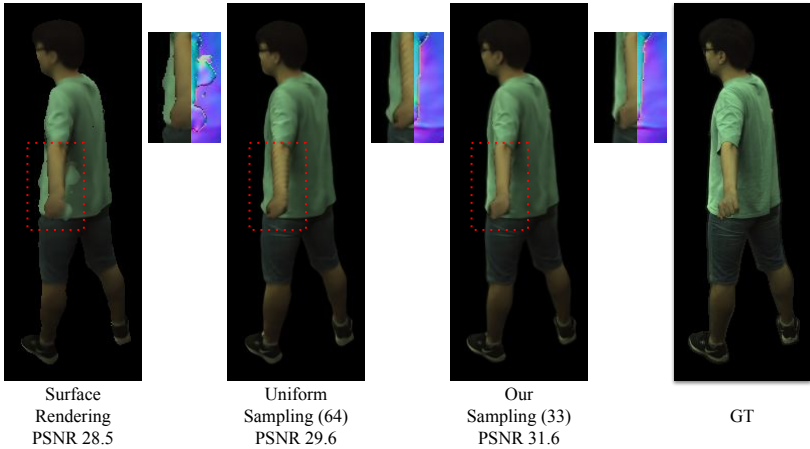
| Surface<br>Rendering<br>PSNR 28.5 | Uniform<br>Sampling (64)<br>PSNR 29.6 | Our<br>Sampling (33)<br>PSNR 31.6 | GT |

Fig. F.1: **Ablation on ray sampling strategies**. We observe severe geometric artifacts with models trained with surface rendering. A simple uniform sampling strategy (as used in [58, 60]) produces stratified artifacts due to the discretized sampling. In contrast, our proposed approach does not suffer from these problems and achieves better result.

the transformed SMPL meshes and supervise the backward LBS network with corresponding skinning weights using Eq. (A.6). We show qualitative results in Fig. F.2.

### F.3   Ablation on Root-finding Initialization

To ablate the effect of multiple initializations for root-finding, we add additional initializations from the nearest 2 SMPL bones but do not observe any noticeable change in metrics. We report PSNR/SSIM/LPIPS as: single initialization - 31.6/0.973/0.050, 2 more initializations: 31.5/0.972/0.049. Also, adding more initializations for root-finding drastically increases memory/time consumption, we thus decide to use only a single initialization for root-finding in our approach.

## G   Additional Quantitative Results

We present complete evaluation metrics including PSNR, SSIM, LPIPS on the test poses of the ZJU-MoCap [60] dataset in Table G.1.

We also report quantitative results on the H36M dataset [26], following the testing protocols proposed by [58] in Table G.2.

Table G.1: **Complete evaluation results on novel pose synthesis.** PSNR, SSIM, LPIPS are reported for the test poses of the ZJU-MoCap dataset.

| | 313 | | | 315 | | | 377 | | |
|---|---|---|---|---|---|---|---|---|---|
| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| NB | 24.1 | 0.908 | 0.126 | 19.8 | 0.867 | 0.152 | 24.2 | 0.917 | 0.119 |
| Ani-N | 23.9 | 0.909 | 0.115 | 19.2 | 0.855 | 0.167 | 22.6 | 0.900 | 0.153 |
| A-NeRF | 22.0 | 0.855 | 0.209 | 18.7 | 0.810 | 0.232 | 22.6 | 0.890 | 0.165 |
| Ours | **24.4** | **0.914** | **0.092** | **20.0** | **0.881** | **0.105** | **25.5** | **0.933** | **0.093** |
| | 386 | | | 387 | | | 390 | | |
| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| NB | 26.1 | 0.894 | 0.171 | 22.7 | 0.902 | 0.135 | 24.2 | 0.882 | 0.164 |
| Ani-N | 25.5 | 0.884 | 0.187 | 23.1 | 0.906 | 0.145 | 23.9 | 0.887 | 0.173 |
| A-NeRF | 24.8 | 0.858 | 0.241 | 22.4 | 0.885 | 0.162 | 22.6 | 0.846 | 0.226 |
| Ours | **27.0** | **0.910** | **0.127** | **24.2** | **0.917** | **0.099** | **24.8** | **0.896** | **0.126** |
| | 392 | | | 393 | | | 394 | | |
| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| NB | 26.0 | 0.916 | 0.135 | 23.5 | 0.900 | 0.132 | 24.1 | 0.888 | 0.150 |
| Ani-N | 24.3 | 0.900 | 0.169 | 23.8 | 0.897 | 0.155 | 24.1 | 0.887 | 0.171 |
| A-NeRF | 23.7 | 0.886 | 0.183 | 22.1 | 0.875 | 0.175 | 22.7 | 0.861 | 0.199 |
| Ours | **26.2** | **0.927** | **0.106** | **24.4** | **0.915** | **0.104** | **25.2** | **0.908** | **0.111** |

Table G.2: **Evaluation results on the H36M dataset.** Numbers of NARF [52] and Ani-N [58] are reported in [84].

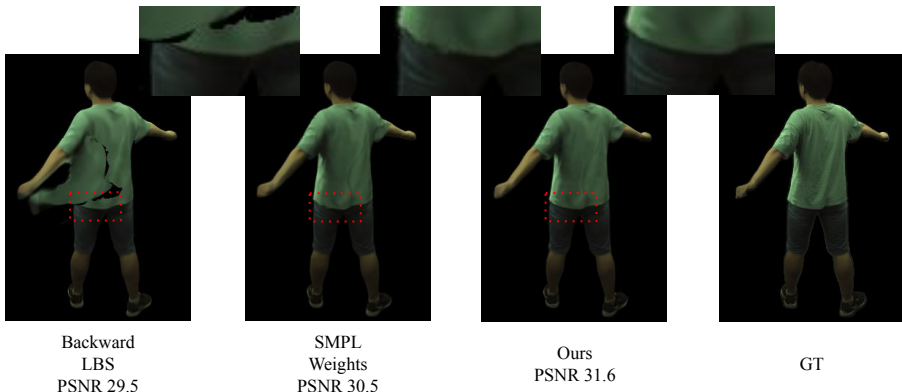| | Training Poses | | | | | | Unseen Poses | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR ↑ | | | SSIM ↑ | | | PSNR ↑ | | | SSIM ↑ | | |
| | NARF | Ani-N | Ours | NARF | Ani-N | Ours | NARF | Ani-N | Ours | NARF | Ani-N | Ours |
| S1 | 21.41 | 22.05 | **24.45** | 0.891 | 0.888 | **0.919** | 20.19 | 21.37 | **23.08** | 0.864 | 0.868 | **0.899** |
| S5 | **25.24** | 23.27 | 24.54 | 0.914 | 0.892 | **0.918** | **23.91** | 22.29 | 22.79 | **0.891** | 0.875 | 0.890 |
| S6 | 21.47 | 21.13 | **24.61** | 0.871 | 0.854 | **0.903** | 22.47 | 22.59 | **24.04** | 0.883 | 0.884 | **0.900** |
| S7 | 21.36 | 22.50 | **24.31** | 0.899 | 0.890 | **0.919** | 20.66 | 22.22 | **22.58** | 0.876 | 0.878 | **0.891** |
| S8 | 22.03 | 22.75 | **24.02** | 0.904 | 0.898 | **0.921** | 21.09 | 21.78 | **22.34** | 0.887 | 0.882 | **0.896** |
| S9 | 25.11 | 24.72 | **26.20** | 0.906 | 0.908 | **0.924** | 23.61 | 23.72 | **24.36** | 0.881 | 0.886 | **0.894** |
| S11 | 24.35 | 24.55 | **25.43** | 0.902 | 0.902 | **0.921** | 23.95 | 23.91 | **24.78** | 0.885 | 0.889 | **0.902** |
| Average | 23.00 | 23.00 | **24.79** | 0.898 | 0.890 | **0.918** | 22.27 | 22.55 | **23.42** | 0.881 | 0.880 | **0.896** |

Fig. F.2: **Ablation on Learned LBS networks**. Backward LBS has difficulties with learning skinning weights for points far from the surface, resulting in artifacts under specific poses. Canonicalization with deterministic SMPL weights results in discretized artifacts on the cloth surface. In contrast, our approach does not suffer from these problems.

## H    Additional Qualitative Results

### H.1    Qualitative Results on ZJU-MoCap Training Poses

We present additional qualitative results on ZJU-MoCap training poses in Fig. H.1. Due to better geometry constraints, our approach better captures cloth wrinkles, textures, and face details. We also avoid extraneous color blobs under novel views which all baselines suffer from.

### H.2    Additional Qualitative Results on ZJU-MoCap Test Poses

We show additional qualitative results on ZJU-MoCap test poses in Fig. H.2. Similar to the results presented in the main paper, A-NeRF and Neural Body do not generalize to these within-distribution poses. Ani-NeRF produces noisy rendering due to its inaccurate backward LBS network. Note that since these results are pose extrapolations, it is not possible to reproduce the exact color and texture of ground-truth images. Still, our approach does not suffer from the artifacts that baselines have demonstrated, resulting in better metrics, especially for LPIPS (Table G.1). We present more qualitative results in the supplementary video.

### H.3    Additional Qualitative Results on Out-of-distribution Poses

We show additional qualitative results on out-of-distribution poses [34] in Fig. H.3. We present more results in the supplementary video.

### H.4    Closest Training Poses to Out-of-distribution Poses

To further demonstrate the generalization ability of our approach, we also visualize the closest training pose from the ZJU-MoCap dataset to out-of-distribution test poses from the AIST++ dataset and the AMASS dataset in Fig. H.4. To find the closest training pose to a test pose, we convert local poses (*i.e.* all pose vectors excluding global orientation) to a matrix representation and find the closest training pose with nearest neighbor search using the converted matrix representation.

### H.5    Qualitative Results on Models Trained with Monocular Videos

In this subsection, we present models trained on monocular videos. For this monocular setup, we use only the first camera of the ZJU-MoCap dataset to train our models. We do not modify our approach and all hyperparameters remain the same as the multi-view setup. We train each model for 500 epochs on 500 frames of selected sequences in which the subjects do repetitive motions while rotating roughly 360 degrees. We animate the trained model with out-of-distribution poses from AIST++ [34]. Qualitative results are shown in Fig. H.5. Even under this extreme setup, our approach can still learn avatars with plausible geometry/appearance and the avatars still generalize to out-of-distribution poses. For the complete animation sequences, please see our supplementary video.

## I    Limitations

As reported in Section D, our approach is relatively slow at inference time. The major bottlenecks are the iterative root-finding [12] and the volume rendering.

Another limitation is that neural rendering-based reconstruction methods tend to overfit the geometry to the texture, resulting in a reconstruction bias. As shown in Fig. I.1, while NeRF-based baselines are unable to recover detailed wrinkles, SDF-based rendering (ours and NeuS) wrongfully reconstructs the stripes on the shirt as part of the geometry. Note that A-NeRF and Ani-NeRF also suffer from this kind of bias. Neural Body demonstrates less overfitting effects. We hypothesize that this is because the structured latent codes in Neural Body are local in space and thus give the color network more flexibility, making the density network less prone to overfitting. Still, Neural Body gives noisy reconstructions and cannot generalize to unseen poses. Resolving this reconstruction bias while maintaining a clean geometry is an interesting avenue for future research.

| A-NeRF | Ani-NeRF | Neural Body | Ours | GT |

Fig. H.1: **Novel View Synthesis Results** on the training poses of ZJU-MoCap.

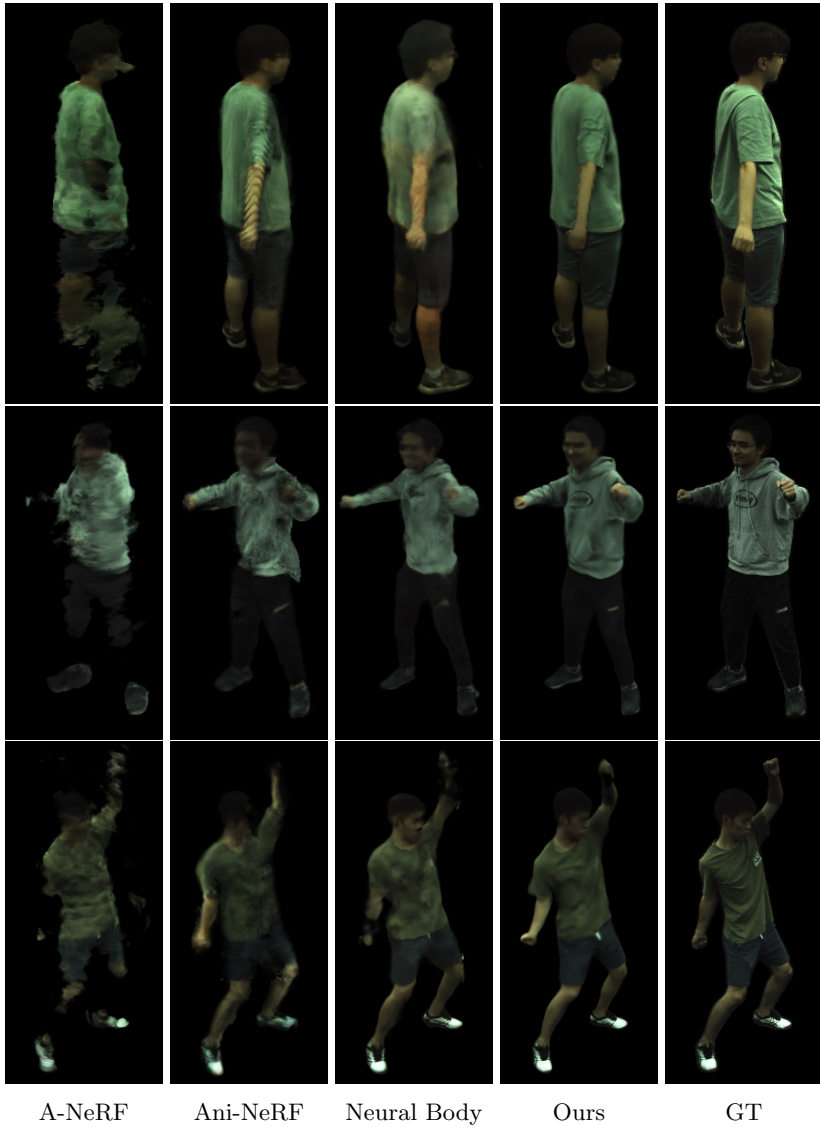|  A-NeRF | Ani-NeRF | Neural Body | Ours | GT |

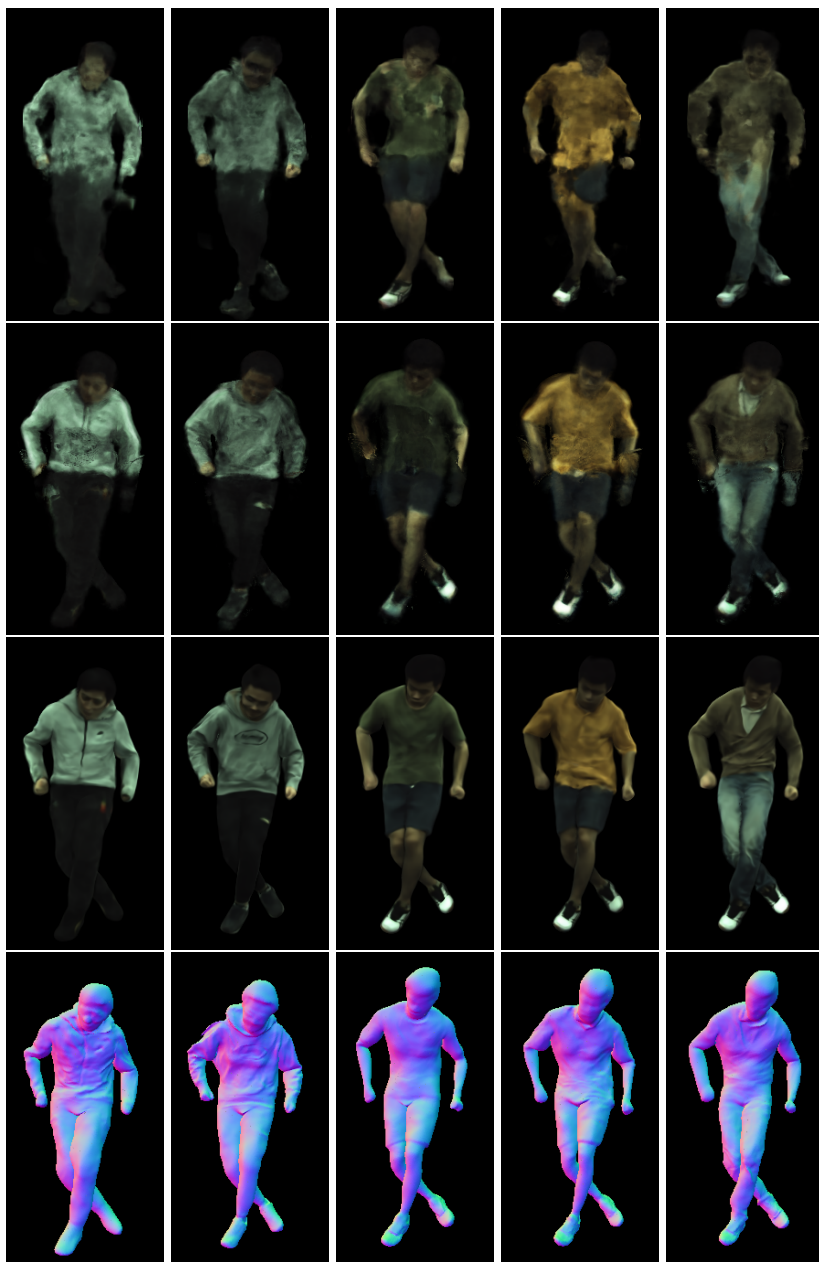Fig. H.2: **Additional Generalization Results on ZJU-MoCap Test Poses.**

Fig. H.3: **Additional Generalization Results on Out-of-distribution Poses.** From top to bottom: Neural Body, Ani-NeRF, ours, and our geometry.
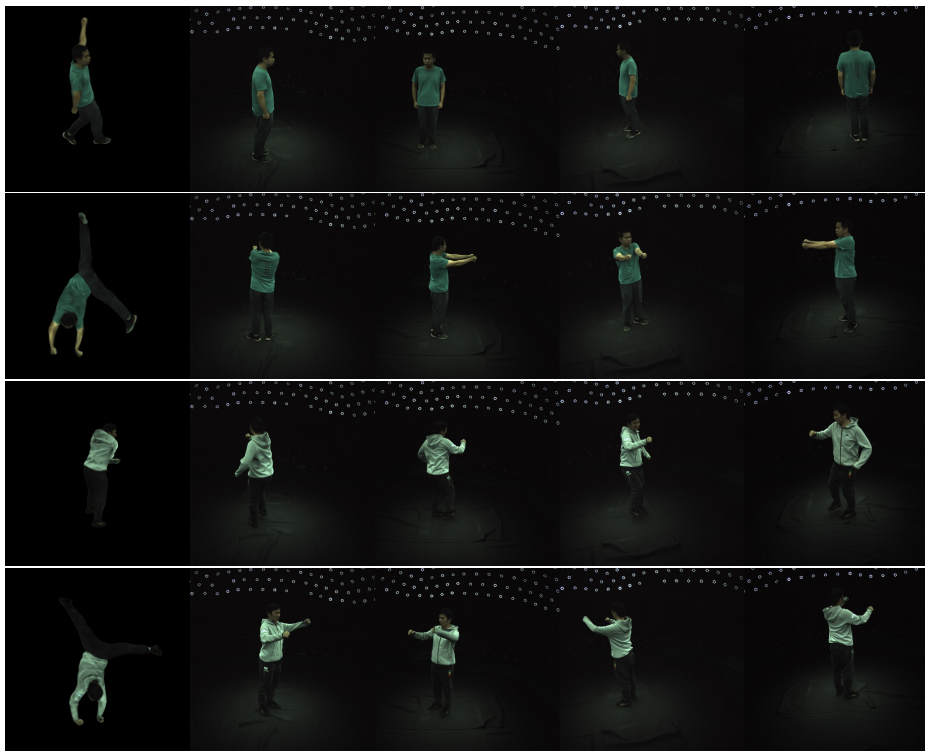
Fig. H.4: **Closest Training Poses to Out-of-distribution Test Poses.** We show rendering results of out-of-distribution poses on the left-most column, while demonstrating 4 training images of the closest training pose to each of the test poses.
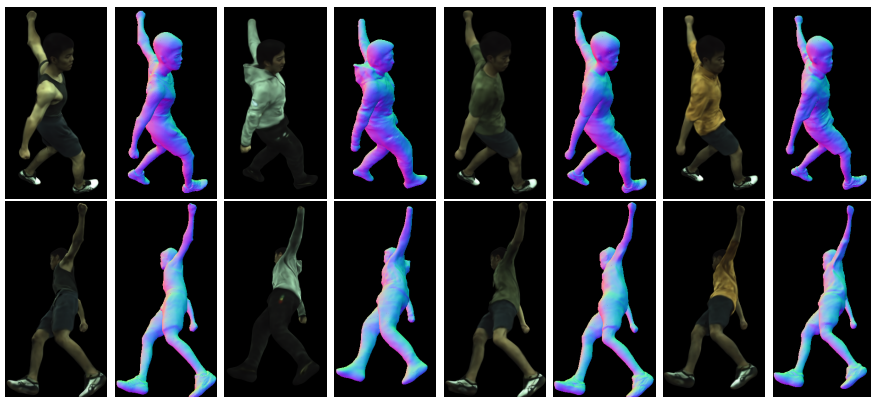


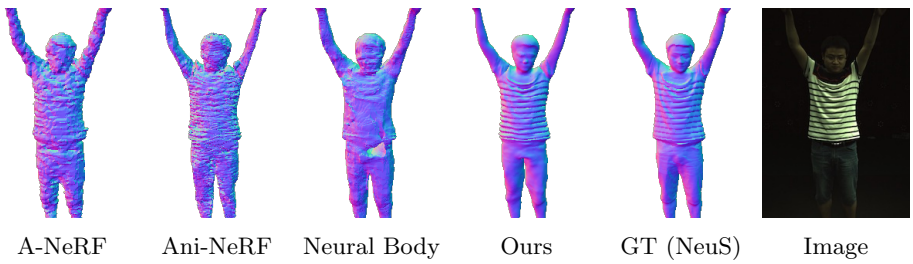Fig. H.5: **Generalization to AIST++ [34] Poses with Models Trained from Monocular Videos.**

| A-NeRF | Ani-NeRF | Neural Body | Ours | GT (NeuS) | Image |

Fig. I.1: **Shape-Appearance Ambiguity**. The Neural Rendering-based reconstruction tends to bake complex textures into the geometry, resulting in a biased geometry reconstruction.