# Mesh Strikes Back: Fast and Efficient Human Reconstruction from RGB videos

Rohit Jena[1][*]    Pratik Chaudhari[1]    James C. Gee[1]    Ganesh Iyer[2]    Siddharth Choudhary[2]

Brandon M. Smith[2]

[1]University of Pennsylvania    [2]Amazon.com, Inc

## Abstract

*Human reconstruction and synthesis from monocular RGB videos is a challenging problem due to clothing, occlusion, texture discontinuities and sharpness, and frame-specific pose changes. Many methods employ deferred rendering, NeRFs and implicit methods to represent clothed humans, on the premise that mesh-based representations cannot capture complex clothing and textures from RGB, silhouettes, and keypoints alone. We provide a counter viewpoint to this fundamental premise by optimizing a SMPL+D mesh and an efficient, multi-resolution texture representation using only RGB images, binary silhouettes and sparse 2D keypoints. Experimental results demonstrate that our approach is more capable of capturing geometric details compared to visual hull, mesh-based methods. We show competitive novel view synthesis and improvements in novel pose synthesis compared to NeRF-based methods, which introduce noticeable, unwanted artifacts. By restricting the solution space to the SMPL+D model combined with differentiable rendering, we obtain dramatic speedups in compute, training times (up to 24x) and inference times (up to 192x). Our method therefore can be used as is or as a fast initialization to NeRF-based methods.*

## 1. Introduction

Our goal is to generate detailed, personalized, and animatable 3D human models. This has many downstream applications, including teleconferencing, entertainment, surveillance, and realistic synthetic data generation. 3D body scanners are the gold standard when it comes to 3D reconstructions. Results are accurate and realistic, but scanners tend to be expensive and ungainly. Moreover, they require additional postprocessing, such as registration and rigging, before they can be used. More recently, CV-based systems have demonstrated recovering realistic 3D human geometry and appearance from monocular images or videos. In the case of monocular *images*, a predictor can be learned
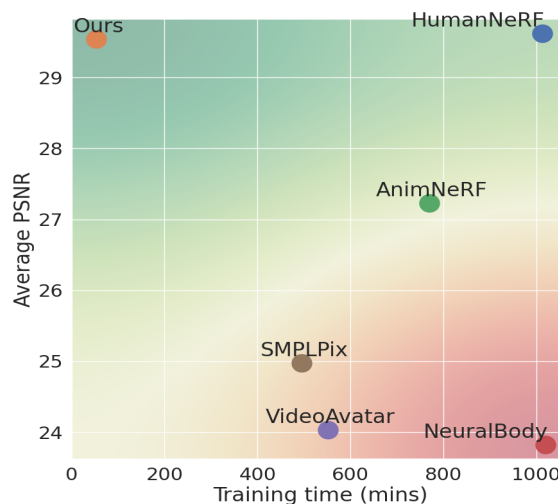
---

[*]Work done during an internship at Amazon



Figure 1. **Performance vs training time tradeoff**: Our method maximizes the performance to training time tradeoff compared to other methods employing meshes, NeRFs and deferred rendering.

from a dataset of real or synthetic humans [53, 69]. However, this is ill-conditioned because a single view is insufficient to estimate the entire 3D human geometry or appearance completely or accurately. Extending it to multi-view or 360° video input would require running inference for all frames and fusing per-frame texture and mesh information. 3D geometry and texture recovery is therefore formulated as an optimization problem. This is an alternative to expensive 3D scanning and motion capture pipelines.

Human-specific NeRFs have recently become popular under this latter, video-based formulation. There are two downsides to typical NeRFs: (1) the highly unconstrained solution space of NeRFs and the use of volume rendering leads to long training and rendering times *e.g*., up to a few days for training [12, 67, 30] and up to minutes for rendering images [22], and (2) dense viewpoints are required. On the contrary, mesh-based reconstruction methods can be faster and less compute-intensive because they do not optimize over a volume. The solution space of the mesh is highly constrained due to flexible yet accurate param-

eterized human body priors like SMPL [32]. Moreover, mesh-based methods may perform better than NeRFs under sparser viewpoints without expensive pretraining or strong regularization that misses geometric details ( [38]), due to the human shape prior.

However, mesh-based methods that rely on a visual hull or silhouette-based approach (*e.g.*, [5]) suffer from ambiguous shape recovery, *i.e.*, any concave surface is not captured by the visual hull and thus cannot be recovered. We show that silhouette-based optimization is highly ill-posed (Sec. 3.1). To overcome this ambiguity, recent methods augment silhouettes with depths or normals [69, 18]. However, obtaining ground truth depths or normals is expensive, and prediction can be error-prone. RGB images provides additional information that can be used across multiple frames to better recover 3D information. However, directly employing an RGB loss is non-trivial because of differentiability challenges in rasterization. NeRFs use the inherent 'softness' of the occupancy field (in the volume integral) to reason about self-occluded portions of the body, allowing pose-refinement [12]. We use a soft differentiable rendering pipeline [31] to emulate this behavior. This allows us to use 'analysis-by-synthesis' as part of our optimization problem. We adopt an approach that uses texture information as part of the optimization similar to NeRFs, but we parameterize the body similar to mesh-based reconstruction methods. Meshes are simpler and more efficient, which affords significant speedup (up to 24x in training time and 192x in inference time) compared to NeRFs.

However, a naive mesh-based optimization to simultaneously minimize RGB and silhouette losses does not work because of *moving targets*. The problem of *moving targets* occur when the RGB losses between the image and a partially learnt texture representation hinder the mesh deformation and vice versa. We propose a method to reduce the ill-posedness and mitigate the problem of moving targets in optimization using a two-stage optimization. We demonstrate that our 3D reconstruction results are similar in quality and accuracy to NeRFs, and significantly better than existing mesh-based reconstruction methods. In addition, we show competitive results in novel view and novel pose synthesis compared to NeRF-based methods.

In summary, this paper makes three main contributions: (1) To our knowledge, we present the first method to incorporate photogrammetric losses in the context of generating human avatars from monocular videos using a mesh representation (Sec. 3). This allows us to optimize texture and geometry using *analysis-by-synthesis* in our optimization without additional auxiliary inputs. (2) An efficient, multi-resolution texture representation using hash encoding capable of capturing fine details is proposed. Unlike texel-based representations, capacity is not wasted on uniform-texture regions. (3) To mitigate the moving targets prob-

lem, where partially learnt texture and geometry hinder each other's loss functions, a novel two-stage optimization is proposed. This ensures stability and optimal convergence.

## 2. Related Work

### 2.1. Human reconstruction via prediction

Recovering 3D human shape and pose estimation using a parametric 3D body model such as SCAPE [44], SMPL [33], SMPL-X [41], STAR [39] or GHUM [70] is an active area of research in the CV community. Most of these approaches directly predict model parameters using a learned model [25, 27, 54, 53, 56, 55, 57, 64, 60, 61]. Recent approaches (*e.g.*, [28, 57, 53, 55]) have improved on prior methods by directly regressing body shape. However, these approaches can only represent the shape and pose of a minimally clothed body and fail to model complex topology due to clothing, hair, *etc*. BodyNet [64] and DeepHuman [76] attempt to predict volumetric representations of the human model from a single image. Implicit representations are an interesting alternative for representing high-fidelity 3D geometry without requiring the entire output volume be kept in memory. In contrast to explicit representations, implicit representations define a surface as a level set of a function. Recent methods, such as PiFu, PiFuHD and PHORHUM [48, 49, 7], learn an implicit surface representation estimated based on pixel aligned features and the depth of 3D locations. Some recent methods combine the benefits of both explicit and implicit representations to represent clothed people [16, 8, 75, 11, 21, 69]. These methods require 3D ground-truth supervision, limiting their applicability to a few datasets and their ability to generalize beyond in-distribution poses. Recently, implicit representations have been used to learn a generative model of 3D people in clothing [13, 6, 17, 50, 14]. However, these approaches require ground truth posed 3D meshes or RGB-D video sequence to learn a model [72, 76, 40, 3, 2, 1]. All of the methods share a common limitation, *i.e.*, predicting the human shape and texture from a single image is ill-posed, and incorrectly regressed predictions are not iteratively refined using auxiliary signals.

### 2.2. Human reconstruction via optimization

**Recovering pose**: Some of the earliest approaches fit model parameters via optimization at test time [5, 9, 29]. Bogo *et al*. [9] optimize SMPL parameters by minimizing the joint reprojection loss and prior terms, whereas [29, 28, 73] employ a likelihood term over the pose via a learned network, and perform iterative regression to minimize reprojection or multiview losses. Pavlakos *et al*. [42] reconstruct SMPL parameters of multiple humans from videos of TV shows. However, these methods only recover a pose with generic shapes and do not capture subject-specific de-

formations and textures.

**Recovering shape and texture**: Alldieck *et al.* [5] extended this approach to monocular video by fusing the unposed silhouettes from all frames to generate a consensus mesh. Each human silhouette, extracted from a video frame, defines a constraint on the human shape, which can be used to estimate deviations in shape. The main problem with this method is that shape is ambiguous, *i.e.*, concave surfaces are not captured. One would need to use auxiliary inputs like depths or normals to disambiguate the problem [69, 4], but depth prediction systems can introduce their own errors.

Mildenhall *et al.* [36, 68] pioneered NeRFs for representing static scenes with a color and density field without requiring any 3D ground truth supervision. Recently this approach has been extended to reconstruct clothed humans as well [43, 58, 12, 67, 24, 66, 30, 65, 23, 19, 62]. These approaches use SMPL as a prior to unpose the human body across multiple frames by transforming the rays from observation space to canonical space which is then rendered using a NeRF. PINA [18] learns a SDF and a learned deformation field to create an animatable avatar from an RGB-D image sequence. Chen *et al.* [15] used a polygon rasterization pipeline to speed up NeRF rendering as a post-process; however, their approach does not reduce NeRF training times. Some recent methods have improved the efficiency of scene-agnostic NeRFs (*e.g.* [71, 51, 37, 20, 46]). These methods demonstrate fast training and rendering times, but adapting them to include an animatable human shape prior requires expensive nearest neighbor operations that erode efficiency gains. However, most of these approaches are computationally expensive. In addition, the resulting representation may have poor generalization to OOD poses.

Our method falls in the category of shape and texture optimization, being most similar to [5, 12, 23, 67], where we aim to contrast NeRFs with its mesh-based equivalent formulation. In contrast to NeRFs, an explicit mesh representation combined with a carefully chosen optimization scheme enables our method to recover accurate geometry, while ensuring photo-realistic rendering, at substantially lower computational and time costs. We show that, contrary to conventional wisdom, our method, when employed with the correct optimization objective, can recover complex geometry (like loose clothing, skirts, long hair, hoodies, *etc.*). Because we constrain the solution space using a mesh, our method is computationally inexpensive, and can be run on a single consumer-grade GPU in about an hour. Our representation allows us to use differentiable rasterization pipelines, drastically reducing its training and inference time. Note that our goal is not to replace NeRF methods, which have asymptotically better performance due to more representational capacity. Rather, *we provide a computationally inexpensive alternative that can be used for real-time rendering and/or on-device applications, or to bootstrap NeRF opti-*
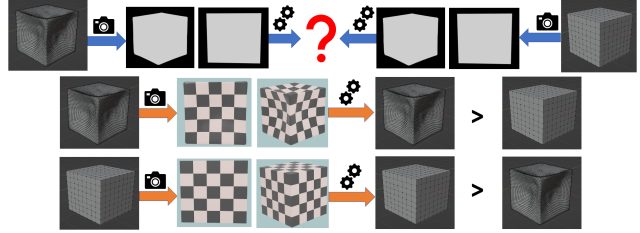


Figure 2. A toy example hightlighting that methods based on visual hulls alone cannot recover concavities in the underlying object. Optimization from visual hull is ill-conditioned, which is disambiguated by multi-view RGB consistency. This idea is used in NeRFs, and in this work we show that it is possible to use RGB and visual hull to optimize a mesh representation.

*mization.* Unlike [43, 30], which have reported failure cases for out-of-distribution poses, our method's accuracy is only limited by the artefacts of the skinning function. A comparison with relevant methods is provided in Fig. 1.

## 3. Method

This paper focuses on jointly recovering accurate geometry and realistic textures from a monocular RGB video of a person using a mesh representation. First, we illustrate the ambiguity in mesh reconstruction from visual hull only, and how multi-view RGB consistency can help disambiguate the problem (Section 3.1). This allows us to avoid using auxiliary inputs, such as depths, landmarks, and normals, which are either expensive to obtain or error-prone if predicted. Next, we discuss the shape and texture representations (Sec 3.2, 3.3) used in our optimization. To enable photogrammetric losses to backpropagate to a parameterized mesh representation, we describe the differentiable rendering pipeline (Section 3.4). Finally, we describe the forward model (Section 3.5), loss functions (Section 3.6) and training pipeline (Section 3.7) to jointly recover the clothed human geometry and texture of the subject.

### 3.1. A motivating toy problem

We use the following toy problem to illustrate the ambiguity of using visual hulls for mesh reconstruction. Consider a cube and the same cube but with all its faces dented inwards, shown in Fig. 2 (top). Rendering the visual hull of both objects gives us the same set of binary silhouettes for all camera angles, making it ambiguous for any optimization scheme to recover a unique mesh from the set of silhouettes. This ambiguity necessitates the use of auxiliary inputs, *e.g.*, depth or normals [18, 69, 47]. Now consider the same scenario, but with the same overlaid texture on both objects. In this case, the two objects have different renders in Fig. 2 (middle & bottom) from the same viewpoints, disambiguating the shape of the underlying object when optimized with a multi-view RGB consistency framework. This

| Method | RGB loss | Mask loss | KPS loss | Representation | Novel pose | Training time | GPU |
|---|---|---|---|---|---|---|---|
| VideoAvatar [5] | ✗ | ✓ | ✓ | SMPL+D | ✓ | 16 hours | NA |
| AnimNeRF [12] | ✓ | ✓ | ✗ | NeRF | ✓ | 26 hours | 48GB |
| Neuralbody [43] | ✓ | ✓ | ✗ | NeRF | ✗ | 14 hours | 48GB |
| HumanNeRF [67] | ✓ | ✓ | ✗ | NeRF | ✓ | 72 hours | 48GB |
| NeuralActor [30] | ✓ | ✓ | ✗ | NeRF | ✓ | 48 hours | 256GB |
| SCARF [19] | ✓ | ✓ | ✗ | NeRF+SMPLX-D | ✓ | 40 hours | 32GB |
| Ours | ✓ | ✓ | ✓ | SMPL+D | ✓ | **<1hour** | **5GB** |

Table 1. Comparison of different methods for human reconstruction. Our simple yet clever use of NeRF-like losses with an SMPL+D representation bridges the gap between mesh-based and NeRF-based optimization with dramatic speedups and compute savings.

is the idea used in NeRFs [36] to recover the occupancy and radiance volume from images alone. Therefore, one can use multi-view RGB consistency as a surrogate to depth, normals, *etc*. The key to using RGB images for mesh optimization is to assign a unique RGB value to each point on the mesh, such that it can guide the mesh vertices to produce consistent renderings in all views. However, doing so introduces a 'moving target' problem (partially optimized mesh and RGB hinder each other's learning) which is non-trivial to optimize. Empirically, we find that carefully formulating the optimization problem (Sec 3.7) allows us to capture complex geometry, including hoodies, loose shirts, pants, skirts, and voluminous hair, better than prior work that uses visual hulls alone for optimization.

### 3.2. Geometry model

Optimizing a high-fidelity textured avatar from monocular or multi-camera RGB video requires us to learn geometry and corresponding texture on the learned geometry. Methods using neural rendering learn a canonical or conditional volume from scratch without using any structural priors [36]. Instead of learning a volume from scratch, we use the parametric SMPL human model [33]. We learn per-frame pose and camera parameters $\{(\boldsymbol{\theta}_i, \boldsymbol{R}_i, \boldsymbol{t}_i)\}_{i \in \{1..n\}}$ and a common shape parameter $\boldsymbol{\beta}$. Since the shape parameter is a low-dimensional embedding capturing human shape, we also learn a per-vertex offset matrix $\mathbf{D} \in \mathbb{R}^{V \times 3}$ where $V$ is the number of vertices in the mesh. Unlike [4] we show that the base SMPL+D model is enough to capture most geometric details (Fig. 6 and in Appendix) and we do not need to use a subdivided SMPL model.

### 3.3. Texture representation

Mesh-based representations generally use predefined UV coordinates for each vertex [5, 4], and the RGB value at any point on the face is determined by evaluating the UV coordinate of the point, and using interpolation from the RGB values in the UV map. This is converted into a texel-map [26, 31]. This representation, although simple, is not adaptive to the complexity of texture on the mesh. Meshes may have regions of low or high texture (*e.g.*, solid-color clothing versus faces, hair, or patterned clothing) which
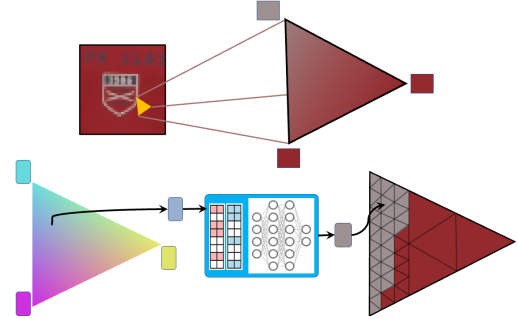


Figure 3. (**Top**) Traditional texturing methods use fixed per-vertex UV coordinates and linearly interpolate the color from the UV coordinates over the face leading to blurry texture. (**Bottom**) We use learnt per-vertex 3D texture coordinates (in cyan,magenta,yellow) and linearly interpolate the texture coordinates over the face. The interpolated texture coordinate is input into the multi-res hash encoding [37] which allows us to represent sharp textures within the face (illustrated by how the network treats the input space across multiple resolutions).

may require a more flexible texture representation. Bilinear interpolation also contributes to blurry textures , making it unsuitable for high-frequency or discontinuous textures. Moreover, deferred rendering [45, 63] may overfit to the UV distribution of the training frames and may not generalize to UV distributions of OOD poses (Fig. 4).

To alleviate these problems, we take inspiration from Müller *et al*. [37] who proposed a multi-resolution hash encoding of the input to capture high-frequency details. The main idea of this work is to use implicit hash collisions to average the gradients at a particular hash lookup entry, and therefore weigh the representation appropriately. We use the same idea to learn a high frequency texture representation. Human avatars have varying levels of texture complexity across their surface—regions such as faces, logos on clothing, *etc*. have high levels of detail, whereas regions such as skin and solid-color clothing have low levels of detail, and require little representation capacity to learn. The hash-encoding representation would learn this implicitly by finding the best embedding which leads to the least RGB reconstruction loss, effectively 'splitting' its representation capacity between high and low-frequency details. For a

given frame $i$ using shape $\boldsymbol{\beta}$, body pose $\boldsymbol{\theta}_i$ and deformation $\mathbf{D}$ we produce a mesh $\boldsymbol{M}_i$. This mesh is then projected to the image plane using camera extrinsics $(\boldsymbol{R}_i, \boldsymbol{t}_i)$ and a fixed intrinsics matrix using the projective transformation to mesh $\boldsymbol{m}_i$. We perform rasterization on the projected mesh to output an image containing face indices and barycentric coordinates of the face index for each pixel. Let face $f_j$ be rendered at pixel $p$ with barycentric coordinates $\{w_{j1}, w_{j2}, w_{j3}\}$ such that $w_{jk} \geq 0$ and $\sum_{k=1}^{3} w_{jk} = 1$. If the texture coordinates of $f_j$ are given by $\{t_{j1}, t_{j2}, t_{j3}\}$, then the rendered point on the face has the texture coordinate $t = \sum_{k=1}^{3} w_{jk} t_{jk}$. The texture coordinate is a low-dimensional input that is passed into the hash encoder and MLP to output the RGB color at pixel $p$. Conventionally, hardcoded 2D texture coordinates are used for vertices of each face, which are interpolated and used to lookup values in a UV map (typically an image). Due to UV unwrapping of a closed mesh, some mesh vertices have multiple texture coordinates. In contrast, we use a learned 3D texture coordinate for each vertex. This serves two purposes: (1) it sidesteps the need for UV unwrapping by moving the textures coordinates into 3D space, and (2) learnable coordinates allow us to expand or shrink the 3D coordinates of each vertex, which in turn accommodate different levels of detail for each face (See Appendix).

## 3.4. Differentiable rendering

A major advantage of using NeRF-like methods is that the volume rendering process is fully differentiable. Moreover, the gradients with respect to the integral (to compute the color) takes occlusions into account. This allows NeRFs to learn a robust occupancy volume and RGB colors to minimize rendering errors. In contrast, rasterization is an inherently non-differentiable operation with respect to the vertices of the mesh [31]. Many solutions have been proposed to approximate gradients of the vertices from gradients in rendered images (*e.g.*, [26, 31, 34]). We use SoftRas [31] due to its ability to flow gradients to the occluded and far-range vertices, allowing us to perform pose refinement via analysis-by-synthesis. Complex pose changes such as bringing an occluded limb into view, rotating joints, etc. can now be performed since SoftRas allows us to pass gradients into the occluded parts of the render. Empirically, we observe that SoftRas is helpful in updating body pose when a small amount of joint rotation is required. To learn the texture, we need the exact forward rasterization to map texture coordinates to RGB values. In SoftRas, we can set the softening parameters $\gamma = \sigma = 0$ and use the RGB loss to guide the texture learning. However, this leads to numerical instability and rendering artefacts in SoftRas. To mitigate this issue, we use NMR [26] to propagate texture gradients.

## 3.5. Forward model

In this section, we describe the forward model for a given frame $i$. The SMPL+D and camera parameters $(\boldsymbol{\beta}, \boldsymbol{\theta}_i, \mathbf{D}, \boldsymbol{R}_i, \boldsymbol{t}_i)$ are used to generate the mesh $\boldsymbol{M}_i$ which is projected into the image plane as $\boldsymbol{m}_i$. The projected mesh and texture network parameters are passed into NMR and SoftRas to give us an opaque and translucent RGB image respectively. For texture parameters $\phi$, we have:

$$\left[ \hat{\boldsymbol{I}}_{i,\text{NMR}}; \hat{\boldsymbol{I}}_{i,\text{SR}} \right] = \text{NMR}(\boldsymbol{m}_i; \phi), \text{SoftRas}(\boldsymbol{m}_i; \text{sg}(\phi), \sigma, \gamma) \tag{1}$$

where $\sigma, \gamma$ are the face blur and depth scale parameters of SoftRas, and sg is the stop-grad operator.

## 3.6. Loss functions

In this section, we describe the losses used to generate our results. Let the masked ground-truth image for frame $i$ be $\boldsymbol{I}_i$ and ground-truth binary foreground be $\boldsymbol{S}_i$.

**Image losses.** The RGB losses are given by:

$$\mathcal{L}_{i,\text{RGB}} = \|\boldsymbol{I}_i - \hat{\boldsymbol{I}}_{i,\text{NMR}}\|_1 + \|\boldsymbol{I}_i - \hat{\boldsymbol{I}}_{i,\text{SR}}\|_1 \tag{2}$$

We also obtain a binary silhouette using NMR and projected mesh $\boldsymbol{m}_i$, which we denote as $\hat{\boldsymbol{S}}_i$. The silhouette loss is defined as an IOU loss:

$$\mathcal{L}_{i,\text{Sil}} = 1 - \frac{\sum_p \hat{\boldsymbol{S}}_i(p) \cdot \boldsymbol{S}_i(p)}{\sum_p (\hat{\boldsymbol{S}}_i(p) + \boldsymbol{S}_i(p) - \hat{\boldsymbol{S}}_i(p) \cdot \boldsymbol{S}_i(p))}. \tag{3}$$

**Keypoint loss.** We add a keypoint loss to help guide the pose of the SMPL model. For simplicity, we use only keypoints corresponding to limbs (elbows, wrists, knees, ankles) and nose. From the mesh $\boldsymbol{M}_i$, we project keypoints $\hat{\boldsymbol{k}}_i$ and encourage their proximity to target 2D keypoints $\boldsymbol{k}_i$ via the loss:

$$\mathcal{L}_{i,\text{kps}} = \|\hat{\boldsymbol{k}}_i - \boldsymbol{k}_i\|_2^2. \tag{4}$$

We run HRNet [59] on each frame to produce $\boldsymbol{k}_i$ for people-snapshot and use the given keypoints for ZJU Mocap.

**Mesh regularization losses.** We add some regularization to the mesh representation. This is especially important for the free-form per-vertex offsets $\mathbf{D}$. Unlike previous works (*e.g.*, [5]), we observe that imposing a low-deformation loss reduces performance (Sec 4.4). This is because loose clothing need not necessarily correspond to a low deformation from the underlying SMPL model. We only encourage normal consistency of adjacent faces in the mesh. Let $\boldsymbol{M}_D$ be the mesh generated from the SMPL parameters $(\boldsymbol{\beta}, \mathbf{0}, \mathbf{D}, \mathbb{I}, \mathbf{0})$ and let $\boldsymbol{M}_0$ be generated from the SMPL parameters $(\boldsymbol{\beta}, \mathbf{0}, \mathbf{0}, \mathbb{I}, \mathbf{0})$. Let $f_j$ be the $j^{\text{th}}$ face of $\boldsymbol{M}_D$ and $f_j'$ be the $j^{\text{th}}$ face of $\boldsymbol{M}_0$. With some abuse of notation,

| | Novel view (people-snapshot) | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Subject ID** | **PSNR↑** | | | | | **SSIM(x100) ↑** | | | | | **LPIPS(x100) ↓** | | | | |
| | VA | SP | NB | AN | Ours | VA | SP | NB | AN | Ours | VA | SP | NB | AN | Ours |
| **m3c** | 22.91 | 22.94 | 23.98 | **29.43** | <u>29.40</u> | 93.16 | 92.56 | 96.12 | **97.11** | <u>96.24</u> | 4.87 | 6.89 | 7.24 | **1.85** | <u>2.65</u> |
| **m4c** | 22.63 | 21.43 | 22.84 | **27.50** | 26.31 | 93.22 | 92.66 | <u>94.81</u> | **95.87** | 94.27 | 6.00 | 8.04 | 10.93 | **3.77** | <u>5.30</u> |
| **f3c** | 22.10 | 21.80 | <u>23.19</u> | 22.96 | **27.25** | 94.35 | 93.95 | <u>95.83</u> | 94.56 | **96.21** | 5.43 | 5.61 | 10.22 | <u>4.61</u> | **3.47** |
| **f4c** | 23.49 | 22.64 | 22.18 | <u>29.03</u> | **29.61** | 93.99 | 93.27 | 95.63 | **96.88** | <u>96.61</u> | 4.12 | 5.92 | 8.52 | **2.10** | <u>2.42</u> |
| | Novel view (ZJU Mocap) | | | | | | | | | | | | | | |
| | VA | SP | NB | HN | Ours | VA | SP | NB | HN | Ours | VA | SP | NB | HN | Ours |
| **C377** | 24.48 | 27.28 | 24.81 | **30.86** | <u>30.58</u> | 93.12 | 94.64 | <u>97.17</u> | **97.45** | 96.51 | 9.01 | 5.82 | 5.71 | **2.58** | <u>4.48</u> |
| **C386** | 27.67 | 29.22 | 25.08 | **33.36** | <u>33.28</u> | 93.72 | 95.80 | <u>97.27</u> | **97.29** | 96.22 | 7.98 | 10.65 | 4.86 | **3.39** | <u>4.18</u> |
| **C387** | 23.30 | 24.27 | 23.60 | **28.58** | <u>28.07</u> | 92.58 | 95.08 | <u>96.08</u> | **96.10** | 94.73 | 9.18 | 12.08 | 6.88 | **3.96** | <u>6.40</u> |
| **C392** | 25.70 | 28.66 | 24.35 | **31.42** | <u>31.35</u> | 92.89 | 95.64 | **96.86** | <u>96.83</u> | 96.05 | 9.52 | 7.54 | 6.37 | **3.76** | <u>5.37</u> |
| **C393** | 23.45 | 24.83 | 24.17 | **28.89** | <u>28.35</u> | 92.30 | 93.60 | **96.35** | <u>95.83</u> | 93.88 | 10.99 | 12.77 | 6.66 | **4.22** | <u>6.43</u> |
| **C394** | 24.46 | 27.34 | 23.97 | <u>30.73</u> | **31.21** | 91.67 | 95.58 | **96.43** | <u>96.16</u> | 95.57 | 11.28 | 9.07 | 6.71 | **3.75** | <u>4.91</u> |
| | Novel pose (ZJU Mocap) | | | | | | | | | | | | | | |
| **C377** | 24.36 | 27.00 | 23.84 | **30.50** | <u>30.48</u> | 93.25 | 96.51 | <u>96.78</u> | **97.41** | 96.54 | 8.29 | 5.81 | 5.59 | **2.69** | <u>4.27</u> |
| **C386** | 28.34 | 30.38 | 23.26 | <u>33.55</u> | **34.03** | 93.84 | 96.60 | <u>96.46</u> | **97.20** | 96.16 | 7.43 | 9.79 | 5.50 | **3.41** | <u>4.00</u> |
| **C387** | 23.02 | 23.80 | 23.15 | **29.02** | <u>28.43</u> | 92.83 | 95.38 | <u>95.58</u> | **96.44** | 95.23 | 8.46 | 11.47 | 6.77 | **3.25** | <u>5.71</u> |
| **C392** | 25.83 | 29.12 | 22.46 | <u>31.43</u> | **32.22** | 92.98 | <u>96.44</u> | 95.97 | **96.89** | 96.37 | 9.40 | 7.26 | 7.03 | **3.70** | <u>5.01</u> |
| **C393** | 23.50 | 24.79 | 22.41 | **29.32** | <u>28.62</u> | 92.49 | 95.07 | <u>95.45</u> | **96.09** | 94.07 | 10.58 | 12.65 | 7.13 | **3.86** | <u>6.47</u> |
| **C394** | 24.33 | 26.99 | 22.19 | <u>30.20</u> | **30.36** | 91.72 | 95.64 | <u>95.43</u> | **95.95** | 95.10 | 11.09 | 8.44 | 7.29 | **4.07** | <u>5.25</u> |

Table 2. Results on novel view and pose synthesis on people-snapshot and ZJU datasets. Best results are in **bold** and 2nd best is <u>underlined</u>.

for two faces $f_j$ and $f_k$, we denote $|f_j \cap f_k|$ as the number of vertices that are shared between both faces. The normal consistency loss is then given as:

$$\mathcal{L}_{NC} = \sum_{|f_j \cap f_k|=2} (1 - \hat{n}_{f_i} \cdot \hat{n}_{f_j}), \qquad (5)$$

where $\hat{n}_f$ is the outward normal of face $f$. We also encourage each face in the mesh to have the same area with and without the deformation. This respects the relative sizes of faces corresponding to different regions in the mesh. If $A_f$ represents the unsigned area of a face $f$, the face area loss is given as:

$$\mathcal{L}_{FA} = \sum_{j} \left( \frac{A_{f_j}}{A_{f'_j}} + \frac{A_{f'_j}}{A_{f_j}} \right). \qquad (6)$$

We prefer this loss instead of the L2 loss $\|A_{f_j} - A_{f'_j}\|_2^2$ because the gradients of L2 loss are small when the area $A_{f_j}$ approaches 0. On the other hand, we want to penalize shrinkage or expansion equally. The loss we propose is of the form $x + \frac{1}{x}$ which achieves its minima at $x = 1$ for positive $x$. The total loss is:

$$\mathcal{L}_f = \sum_{i} \left( \lambda_{\text{RGB}} \mathcal{L}_{i,\text{RGB}} + \lambda_{\text{Sil}} \mathcal{L}_{i,\text{Sil}} + \lambda_{\text{kps}} \mathcal{L}_{i,\text{kps}} \right)$$
$$+ \lambda_{NC} \mathcal{L}_{NC} + \lambda_{FA} \mathcal{L}_{FA} \qquad (7)$$

Unlike [12], we do not add any other regularization terms on $\beta$ or temporal pose consistency or deviation terms.

## 3.7. Two-stage training

Given a forward model to render RGB and silhouettes from mesh parameters $\beta, \theta_i, \mathbf{D}, \mathbf{R}_i, t_i$ and texture parameters $\phi$, an intuitive way to learn all the parameters is to jointly optimize them. Let $\Theta = \{\phi^*, \beta^*, \mathbf{D}^*, \{\theta_i^*, \mathbf{R}_i^*, t_i^*\}_{i \in \{1...n\}}\}$ be the set of all optimizable parameters. The optimization is of the form: $\Theta^* = \arg\min_{\Theta} \mathcal{L}_f$. This optimization is still highly underconstrained. Meshes obtained using this procedure are jagged with blurry textures because texture and deformation parameters locally optimize their own parameters (examples in Appendix) leading to the *moving target* problem. To alleviate this problem, we propose a two-stage training procedure. In the first stage, we use a *per-face* RGB value as a 'base' texture color instead of the full texture network. This allows for a coarse alignment of the mesh vertices for each frame. Constraining the color of each face to just one optimizes the mesh parameters to place it in the best possible location and scale to minimize RGB and silhouette losses, thus ensuring photogrammetric consistency in the optimization. In the second stage, the deformations, shape, and per-face RGB values are fixed, and the texture network is trained with per-frame pose refinement. This allows for fine-grained alignment of the poses for each frame. Implementation details are provided in the Appendix.

## 4. Results

To evaluate our method, we look at the following aspects of geometry and texture recovery: (1) novel-view and pose

Figure 4. Novel pose synthesis (on **m3c** and **f3c**) from left to right: SMPLpix, VideoAvatar, NeuralBody, AnimNeRF, Ours. Neural-Body does not generalize to novel poses, AnimNeRF introduces cloud and tearing artifacts. Our method preserves detailed texture and doesn't produce artifacts. More results in Appendix.
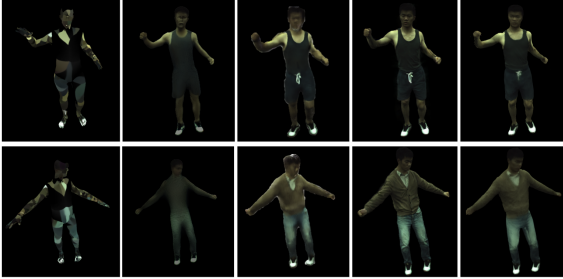


Figure 5. **Novel view on ZJU dataset, Left to right**: VideoAvatar, SMPLPix, NeuralBody, HumanNeRF, Ours. Our method performs dramatically better than VA/SMPLPix and is comparable to NeRF methods at significant training and inference speedups.

synthesis, (2) training/inference time and compute requirements, (3) geometry reconstruction. For experiments (1-2), we use People Snapshot [5] and ZJU Mocap [43] datasets. For (3), we use the Self-Recon synthetic dataset [23]. For people-snapshot, we follow the same subjects and experiment setup as [12], and for ZJU Mocap we use the same set of subjects as [67]. For ZJU Mocap, we use frames 0-450 in cameras 1,7,13,19 for training, and the rest of the frames for novel pose reconstruction. We use frames 0-450 from cameras 5,10,15,20 for novel view reconstruction. We choose baselines across a spectrum of representation choices: SMPLPix [45] (**SP**) which uses deferred rendering, VideoAvatar [5] (**VA**) which performs SMPL+D optimization, NeuralBody [43] (**NB**), HumanNeRF [67] (**HN**) and AnimNeRF [12] (**AN**) which are SOTA NeRF methods.

### 4.1. Novel view synthesis

We evaluate novel view synthesis by holding back a certain set of test frames of the subject to check the quality of reconstruction for those frames (similar to [12]). For ZJU Mocap, we evaluate frames 0-450 from cameras 5,10,15,20.

| FID score ↓ | | | | |
|---|---|---|---|---|
| **Method** | **m3c** | **m4c** | **f3c** | **f4c** |
| **SMPLpix** | 199.04 | 210.27 | 211.09 | 212.16 |
| **Neural body** | 402.47 | 357.77 | 328.26 | 358.54 |
| **VideoAvatar** | 189.82 | _186.03_ | _206.07_ | 162.79 |
| **AnimNeRF** | _183.03_ | 200.73 | 237.79 | **150.80** |
| **Ours** | **178.71** | **184.26** | **203.42** | _159.31_ |
| VGGFace2 ↑ | | | | |
| **SMPLpix** | .4808 | .6472 | .6222 | .4853 |
| **NeuralBody** | .3713 | .3743 | .4301 | .0000 |
| **VideoAvatar** | .8135 | .8799 | .8976 | .8417 |
| **AnimNeRF** | **.9079** | **.8974** | **.9452** | **.9259** |
| **Ours** | _.8766_ | _.8926_ | _.9380_ | _.8948_ |

Table 3. Quantitative analysis of texture quality of novel poses.

Results are shown in Table 2. Our method performs very competitively with AnimNeRF and HumanNeRF in terms all three metrics (PSNR, SSIM, LPIPS), and outperforms all other baselines by a substantial margin. On ZJU Mocap, our method competes with HumanNeRF consistently on the PSNR and LPIPS metrics, showing that our method can faithfully reconstruct accurate and realistic rendering. NeuralBody and SMPLpix are trained on the training poses only, and fail to generalize to novel views, even if the view deviations are very small from the training frames. VideoAvatar performs a multi-stage optimization, where the geometry is optimized from silhouettes only, and then a texture optimization step is performed. Errors from the mesh optimization propagate to the texture, leading to a low quality texture, and subsequently a low quality render. Our method recovers intricate details like loose clothing, loose pants, hoodies, hair, and skirts (Fig. 6). More qualitative results are in the Appendix.

### 4.2. Novel pose synthesis

We use a set of held-back frames for novel pose synthesis in the ZJU dataset. Results are in Tab. 2. SMPLpix and Neural Body do not generalize well because novel poses that are not seen during training results in a distribution shift in the inputs of the frameworks. VideoAvatar has an uncanny valley effect in the faces of its rendered outputs. HumanNeRF achieves highly realistic results capturing nuances in body geometry and texture.

**OOD pose rendering**: However, we notice that the pose distribution is not very different from those in training frames. Moreover, the people-snapshot dataset doesn't contain frames with other poses than the A-pose. Therefore, we also compare the realism of textures and faces in an a set of OOD poses. We curate a set of poses from the AMASS dataset [35]. We compare the realism of the models by evaluating the Fréchet Inception Distance (FID score) [52] of the input frames with novel pose renders, and
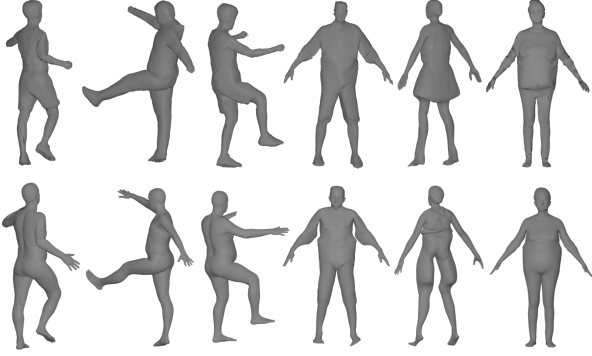
Figure 6. Geometry reconstruction quality of our method (top) and VideoAvatar (bottom) on ZJU Mocap, people-snapshot, and Self-Recon datasets. Our method captures loose fitting pants (1,3) and clothes (2,4,6), hoodie (2), hair (1,2,3), skirts (6).

| Metric | Method | F1 | F2 | F3 | M1 | M2 |
|--------|--------|------|------|------|------|------|
| Chamfer | VideoAvatar | 1.47 | 1.05 | 1.41 | 1.20 | 1.08 |
| | Ours | **1.15** | **0.96** | **1.05** | **1.01** | **0.93** |
| P2S | VideoAvatar | 0.59 | 0.46 | 0.57 | 0.54 | 0.45 |
| | Ours | **0.42** | **0.40** | **0.37** | **0.40** | **0.35** |

Table 5. Reconstruction loss (cm) on Self-Recon synthetic dataset.

comparing the face texture of the rendered images with that of the input frames. Since we use the same novel poses for all methods, the differences in FID must come from texture quality. To evaluate texture quality of faces, we use face identification as a proxy task. We use MTCNN [74] to detect faces from images and VGGFace2 [10] to generate a template feature vector for each method. We use the face similarity metric between the template of the method and that of the input data, as proposed in [10]. Results in Tab. 3 and Fig. 4 shows that our method preserves the subject identity significantly more than VideoAvatar, showing that our method can recover accurate texture with a mesh. AnimNeRF reconstructs the texture well in the parts on the surface, but introduces cloud and tearing artifacts, especially when regions around the unseen areas (armpits and thighs) are stretched too much. Our method doesn't have such an issue, since our representation is based on a mesh. Moreover, the texture distortion for novel poses is virtually non-existent. More qualitative results are shown in Appendix.

### 4.3. Training/inference time and compute

| Method | Training time (min) | Inference time (sec) | GPU usage (GB-hrs) |
|--------|---------------------|----------------------|---------------------|
| HN [67] | 1013.35 | 3.51 | 399.09 |
| AN [12] | 769.26 | 11.54 | 591.38 |
| NB [43] | 1020.27 | 0.77 | 62.22 |
| Ours | **43.31** | **0.06** | **4.11** |

Table 4. Averaged total training and per image inference time (in minutes) and total GPU usage (GB-hr). Our model train upto 24x faster than HumanNeRF on the same data while using upto 4x less compute. Results for individual subjects are in the Appendix.

We compare the training and inference time and compute required our method against NeRF methods in Tab. 4. NeRFs have achieved huge successes in representing scenes faithfully with very accurate rendering. However, they take

a prohibitively long time to train a single scene. Although several improvements have been proposed for static scenes, the main bottleneck for AnimNeRF is the KNN step for each sample along the ray, and unposing the transformation to the canonical space. This is a computationally expensive step since it has to be done for each point from each ray independently. Moreover, volume rendering leads to a significantly higher inference time and compute requirements [30]. In contrast, unposing is trivial using our method because the rasterization consists of the face index with barycentric coordinates.

### 4.4. Geometry reconstruction

We quantitatively compare the effectiveness of our method to recover the underlying geometry from the set of images using the Self-Recon dataset [23], which consists of renderings of 5 human subjects with their ground truth meshes. We compute the average Chamfer distance and Point-to-Surface (P2S) measures. Our comparison with VideoAvatar [5], the other mesh-based optimization method is shown in Tab. 5. Note that our lower distances show that even a low dimensional mesh can capture complex details like loose clothing, hair, etc. with the right optimization and training scheme. Qualitative results on all three datasets are in Fig. 6 and Appendix.

### 5. Conclusion

There are myriad applications for detailed, personalized, and animatable 3D human models, and they become increasingly practical as generation times, data acquisition, and hardware requirements decrease. One promising direction is to bootstrap the training process of human-specific NeRFs [67, 12] with the geometry and texture learnt from our model. Since our model directly provides 3D coordinates and RGBs in the canonical space, volume rendering is not needed during this pretraining step. While NeRFs have demonstrated great versatility, *i.e.*, they are scene-agnostic, they are prohibitive for reconstructing well-defined objects with strong priors (*e.g.* faces, human avatars) given their steep training requirements. For these applications, we argue that our mesh-based system is competitive and, in some cases, favorable. We have demonstrated that our approach is capable of generating results with very competitive performance to SOTA human-specific NeRFs [12, 67], but in a tiny fraction of the training time and compute.

## References

[1] Axyz dataset. http://secure.axyz-design.com/.

[2] Renderpeople dataset. http://renderpeople.com/.

[3] Twindom dataset. http://web.twindom.com/.

[4] T. Alldieck, M. Magnor, W. Xu, C. Theobalt, and G. Pons-Moll. Detailed human avatars from monocular video. In *2018 International Conference on 3D Vision (3DV)*, pages 98–109, Los Alamitos, CA, USA, sep 2018. IEEE Computer Society.

[5] T. Alldieck, M. Magnor, W. Xu, C. Theobalt, and G. Pons-Moll. Video based reconstruction of 3d people models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8387–8397, Jun 2018. CVPR Spotlight Paper.

[6] T. Alldieck, H. Xu, and C. 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.

[7] T. Alldieck, M. Zanfir, and C. Sminchisescu. Photorealistic monocular 3d reconstruction of humans wearing clothing. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[8] B. L. Bhatnagar, C. Sminchisescu, C. Theobalt, and G. Pons-Moll. Combining implicit function learning and parametric models for 3d human reconstruction. In *European Conference on Computer Vision (ECCV)*. Springer, aug 2020.

[9] F. Bogo, A. Kanazawa, C. Lassner, P. Gehler, J. Romero, and M. J. Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In B. Leibe, J. Matas, N. Sebe, and M. Welling, editors, *European Conference on Computer Vision (ECCV)*, pages 561–578. Springer International Publishing, 2016.

[10] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman. Vggface2: A dataset for recognising faces across pose and age. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*, pages 67–74. IEEE, 2018.

[11] Y. Cao, G. Chen, K. Han, W. Yang, and K.-Y. K. Wong. Jiff: Jointly-aligned implicit face function for high quality single view clothed human reconstruction. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[12] J. Chen, Y. Zhang, D. Kang, X. Zhe, L. Bao, X. Jia, and H. Lu. Animatable neural radiance fields from monocular rgb videos, 2021.

[13] X. Chen, T. Jiang, J. Song, J. Yang, M. J. Black, A. Geiger, and O. Hilliges. gdna: Towards generative detailed neural avatars. *arXiv*, 2022.

[14] X. Chen, Y. Zheng, M. J. Black, O. Hilliges, and A. Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. In *International Conference on Computer Vision (ICCV)*, 2021.

[15] Z. Chen, T. Funkhouser, P. Hedman, and A. Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. *arXiv preprint arXiv:2208.00277*, 2022.

[16] E. Corona, G. Pons-Moll, G. Alenyà, and F. Moreno-Noguer. Learned vertex descent: A new direction for 3d human model fitting. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[17] B. Deng, J. Lewis, T. Jeruzalski, G. Pons-Moll, G. Hinton, M. Norouzi, and A. Tagliasacchi. Neural articulated shape approximation. In *The European Conference on Computer Vision (ECCV)*. Springer, August 2020.

[18] Z. Dong, C. Guo, J. Song, X. Chen, A. Geiger, and O. Hilliges. PINA: Learning a personalized implicit neural avatar from a single rgb-d video sequence. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[19] Y. Feng, J. Yang, M. Pollefeys, M. J. Black, and T. Bolkart. Capturing and animation of body and clothing from monocular video, 2022.

[20] S. J. Garbin, M. Kowalski, M. Johnson, J. Shotton, and J. Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *ICCV*, 2021.

[21] T. He, Y. Xu, S. Saito, S. Soatto, and T. Tung. Arch++: Animation-ready clothed human reconstruction revisited. *arXiv*, abs/2108.07845, 2021.

[22] P. Hedman, P. P. Srinivasan, B. Mildenhall, J. T. Barron, and P. Debevec. Baking neural radiance fields for real-time view synthesis. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5855–5864, Los Alamitos, CA, USA, oct 2021. IEEE Computer Society.

[23] B. Jiang, Y. Hong, H. Bao, and J. Zhang. Selfrecon: Self reconstruction your digital avatar from monocular

video. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[24] W. Jiang, K. M. Yi, G. Samei, O. Tuzel, and A. Ranjan. Neuman: Neural human radiance field from a single video. In *European Conference on Computer Vision (ECCV)*, 2022.

[25] A. Kanazawa, M. J. Black, D. W. Jacobs, and J. Malik. End-to-end recovery of human shape and pose. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[26] H. Kato, Y. Ushiku, and T. Harada. Neural 3d mesh renderer. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[27] M. Kocabas, N. Athanasiou, and M. J. Black. Vibe: Video inference for human body pose and shape estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[28] N. Kolotouros, G. Pavlakos, M. J. Black, and K. Daniilidis. Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.

[29] N. Kolotouros, G. Pavlakos, D. Jayaraman, and K. Daniilidis. Probabilistic modeling for human mesh recovery. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.

[30] L. Liu, M. Habermann, V. Rudnev, K. Sarkar, J. Gu, and C. Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. *ACM Trans. Graph. (ACM SIGGRAPH Asia)*, 2021.

[31] S. Liu, T. Li, W. Chen, and H. Li. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2019.

[32] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, Oct. 2015.

[33] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. Smpl: A skinned multi-person linear model. *ACM Trans. Graph.*, 34(6), oct 2015.

[34] M. M. Loper and M. J. Black. Opendr: An approximate differentiable renderer. In D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, editors, *European Conference on Computer Vision (ECCV)*, pages 154–169, Cham, 2014. Springer International Publishing.

[35] N. Mahmood, N. Ghorbani, N. F. Troje, G. Pons-Moll, and M. J. Black. AMASS: Archive of motion capture as surface shapes. In *International Conference on Computer Vision*, pages 5442–5451, Oct. 2019.

[36] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision (ECCV)*, 2020.

[37] T. Müller, A. Evans, C. Schied, and A. Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022.

[38] M. Niemeyer, J. T. Barron, B. Mildenhall, M. S. M. Sajjadi, A. Geiger, and N. Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5480–5490, June 2022.

[39] A. A. A. Osman, T. Bolkart, and M. J. Black. STAR: Sparse trained articulated human body regressor. In *European Conference on Computer Vision (ECCV)*, volume LNCS 12355, pages 598–613, Aug. 2020.

[40] P. Patel, C.-H. P. Huang, J. Tesch, D. T. Hoffmann, S. Tripathi, and M. J. Black. AGORA: Avatars in geography optimized for regression analysis. In *Proceedings IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2021.

[41] G. Pavlakos, V. Choutas, N. Ghorbani, T. Bolkart, A. A. A. Osman, D. Tzionas, and M. J. Black. Expressive body capture: 3D hands, face, and body from a single image. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 10975–10985, 2019.

[42] G. Pavlakos, E. Weber, M. Tancik, and A. Kanazawa. The one where they reconstructed 3d humans and environments in tv shows. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVII*, pages 732–749. Springer, 2022.

[43] S. Peng, Y. Zhang, Y. Xu, Q. Wang, Q. Shuai, H. Bao, and X. Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

[44] L. Pishchulin, S. Wuhrer, T. Helten, C. Theobalt, and B. Schiele. Building statistical shape spaces for 3d

human modeling. *Pattern Recogn.*, 67(C):276–286, jul 2017.

[45] S. Prokudin, M. J. Black, and J. Romero. Smplpix: Neural avatars from 3d human models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1810–1819, 2021.

[46] C. Reiser, S. Peng, Y. Liao, and A. Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *ICCV*, 2021.

[47] J. Romero, D. Tzionas, and M. J. Black. Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, Nov. 2017.

[48] S. Saito, Z. Huang, R. Natsume, S. Morishima, A. Kanazawa, and H. Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. *arXiv preprint arXiv:1905.05172*, 2019.

[49] S. Saito, T. Simon, J. Saragih, and H. Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[50] S. Saito, J. Yang, Q. Ma, and M. 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.

[51] Sara Fridovich-Keil and Alex Yu, M. Tancik, Q. Chen, B. Recht, and A. Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022.

[52] M. Seitzer. pytorch-fid: FID Score for PyTorch. https://github.com/mseitzer/pytorch-fid, August 2020. Version 0.2.1.

[53] A. Sengupta, I. Budvytis, and R. Cipolla. Synthetic training for accurate 3d human pose and shape estimation in the wild. In *British Machine Vision Conference (BMVC)*, September 2020.

[54] A. Sengupta, I. Budvytis, and R. Cipolla. Hierarchical Kinematic Probability Distributions for 3D Human Shape and Pose Estimation from Images in the Wild. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021.

[55] A. Sengupta, I. Budvytis, and R. Cipolla. Probabilistic 3d human shape and pose estimation from multiple unconstrained images in the wild. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16094–16104, June 2021.

[56] A. Sengupta, I. Budvytis, and R. Cipolla. Probabilistic estimation of 3d human shape and pose with a semantic local parametric model. In *British Machine Vision Conference (BMVC)*, pages 16094–16104, November 2021.

[57] B. M. Smith, V. Chari, A. Agrawal, J. M. Rehg, and R. Sever. Towards accurate 3d human body reconstruction from silhouettes. In *International Conference on 3D Vision (3DV)*, pages 279–288, 2019.

[58] S.-Y. Su, F. Yu, M. Zollhöfer, and H. Rhodin. A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose. In *Advances in Neural Information Processing Systems*, 2021.

[59] K. Sun, B. Xiao, D. Liu, and J. Wang. Deep high-resolution representation learning for human pose estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5686–5696, 2019.

[60] Y. Sun, Q. Bao, W. Liu, Y. Fu, B. Michael J., and T. Mei. Monocular, one-stage, regression of multiple 3d people. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.

[61] Y. Sun, W. Liu, Q. Bao, Y. Fu, T. Mei, and M. J. Black. Putting people in their place: Monocular regression of 3d people in depth. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[62] G. Te, X. Li, X. Li, J. Wang, W. Hu, and Y. Lu. Neural capture of animatable 3d human from monocular video. In *European Conference on Computer Vision (ECCV)*, 2022.

[63] J. Thies, M. Zollhöfer, and M. Nießner. Deferred neural rendering: Image synthesis using neural textures. *Acm Transactions on Graphics (TOG)*, 38(4):1–12, 2019.

[64] G. Varol, D. Ceylan, B. Russell, J. Yang, E. Yumer, I. Laptev, and C. Schmid. BodyNet: Volumetric inference of 3D human body shapes. In *European Conference on Computer Vision (ECCV)*, 2018.

[65] S. Wang, K. Schwarz, A. Geiger, and S. Tang. Arah: Animatable volume rendering of articulated human sdfs. In *European Conference on Computer Vision (ECCV)*, 2022.

[66] C.-Y. Weng, B. Curless, and I. Kemelmacher-Shlizerman. Vid2actor: Free-viewpoint animatable person synthesis from video in the wild. *arXiv preprint arXiv:2012.12884*, 2020.

[67] C.-Y. Weng, B. Curless, P. P. Srinivasan, J. T. Barron, and I. Kemelmacher-Shlizerman. HumanNeRF: Free-viewpoint rendering of moving people from monocular video. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16210–16220, June 2022.

[68] Y. Xie, T. Takikawa, S. Saito, O. Litany, S. Yan, N. Khan, F. Tombari, J. Tompkin, V. Sitzmann, and S. Sridhar. Neural fields in visual computing and beyond. *Computer Graphics Forum*, 2022.

[69] Y. Xiu, J. Yang, D. Tzionas, and M. J. Black. ICON: Implicit Clothed humans Obtained from Normals. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13296–13306, June 2022.

[70] H. Xu, E. G. Bazavan, A. Zanfir, W. T. Freeman, R. Sukthankar, and C. Sminchisescu. Ghum & ghuml: Generative 3d human shape and articulated pose models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6184–6193, 2020.

[71] A. Yu, R. Li, M. Tancik, H. Li, R. Ng, and A. Kanazawa. PlenOctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021.

[72] T. Yu, Z. Zheng, K. Guo, P. Liu, Q. Dai, and Y. Liu. Function4d: Real-time human volumetric capture from very sparse consumer rgbd sensors. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR2021)*, June 2021.

[73] H. Zhang, Y. Tian, X. Zhou, W. Ouyang, Y. Liu, L. Wang, and Z. Sun. Pymaf: 3d human pose and shape regression with pyramidal mesh alignment feedback loop. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.

[74] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE signal processing letters*, 23(10):1499–1503, 2016.

[75] Z. Zheng, T. Yu, Y. Liu, and Q. Dai. Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction, 2021.

[76] Z. Zheng, T. Yu, Y. Wei, Q. Dai, and Y. Liu. Deep-human: 3d human reconstruction from a single image. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.