

# Weakly-Supervised 3D Reconstruction of Clothed Humans via Normal Maps

Jane Wu<sup>1</sup>    Diego Thomas<sup>2</sup>    Ronald Fedkiw<sup>1,3</sup>  
<sup>1</sup>Stanford University    <sup>2</sup>Kyushu University    <sup>3</sup>Epic Games

janehwu@stanford.edu, thomas@ait.kyushu-u.ac.jp, fedkiw@cs.stanford.edu

## Abstract

We present a novel deep learning-based approach to the 3D reconstruction of clothed humans using weak supervision via 2D normal maps. Given a single RGB image or multiview images, our network infers a signed distance function (SDF) discretized on a tetrahedral mesh surrounding the body in a rest pose. Subsequently, inferred pose and camera parameters are used to generate a normal map from the SDF. A key aspect of our approach is the use of Marching Tetrahedra to (uniquely) compute a triangulated surface from the SDF on the tetrahedral mesh, facilitating straightforward differentiation (and thus backpropagation). Thus, given only ground truth normal maps (with no volumetric information ground truth information), we can train the network to produce SDF values from corresponding RGB images. Optionally, an additional multiview loss leads to improved results. We demonstrate the efficacy of our approach for both network inference and 3D reconstruction.

## 1. Introduction

Recent work on 3D human digitization has largely focused on the fully-supervised setting, where deep neural networks (DNNs) are trained to explicitly fit so-called ground truth 3D geometry [70, 88, 89, 102]. In such approaches, high-end capture setups (with 4D scanners or a large number of cameras) are typically used to obtain high-quality, multiview training data [65, 87]. Inferring 3D geometry and appearance from 2D information is a highly underconstrained problem; thus, it can be challenging for models trained on such high-quality data to generalize to the lower quality images typical of consumer-grade devices (such as phones and webcams). However, the ability to do so is crucial to the democratization of digital humans required for many applications in AR/VR, robotics, healthcare, etc.

Given a single (monocular) RGB image, we estimate the clothed body with a DNN that infers signed distance function (SDF) values for each vertex of a tetrahedral mesh surrounding the body (similar to [40, 41, 58, 86]). Implicit surfaces have become a common choice for 3D reconstruction

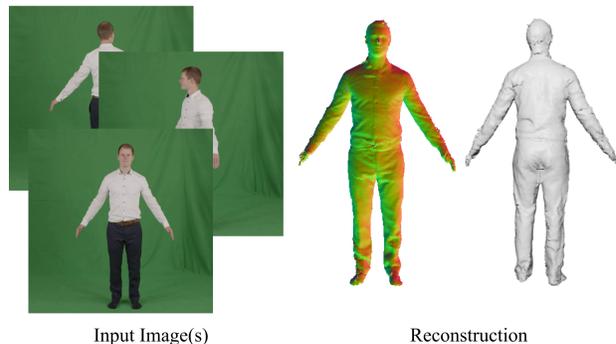


Figure 1. Given in-the-wild image(s) or video, our method is able to reconstruct clothed humans using inferred normal maps as the supervisory signal. Sample frames from the input video are shown to the left, and the predicted triangle mesh is shown to the right (the front-facing mesh is shaded with its normal map).

from images (see e.g. [11, 64, 70, 79, 88, 89, 102]), especially as neural radiance fields (NeRFs) [51] have gained popularity. Importantly, the tetrahedral mesh data structure enables our use of Marching Tetrahedra to uniquely compute a triangulated surface. The resulting algorithm is straightforward to differentiate, alleviating concerns associated with the nondeterministic nature of Marching Cubes (see e.g. [48, 66]) or the ray tracing of an implicit surface (see e.g. [11, 56, 82, 92]).

Our approach is a follow-up to [58], which also uses an explicit representation of an SDF on a tetrahedral mesh; however, we add a second explicit representation of the surface via a triangle mesh. We thus have access to two explicit versions of the neural SDF, and energies can be conveniently formulated for either the volume or the surface or both. While similar in spirit to [50], our method does not require the construction of a velocity in order to capture these energies with an evolving level set function; thus, we can control mesh based properties (e.g. area and dihedral angles) that would be lost when converting to a velocity field. On the other hand, the approach in [50] could be used to alleviate locking concerns in cases where discretizations on the triangle mesh and the tetrahedral mesh do not interact as

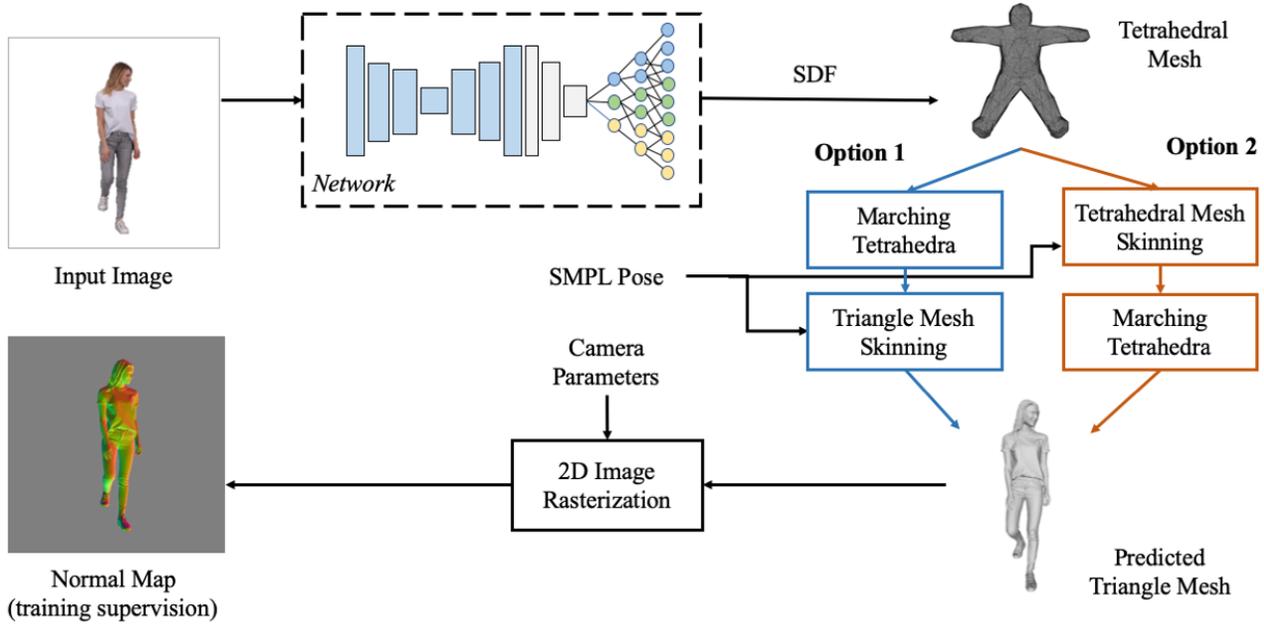


Figure 2. Given an RGB image, a graph-based DNN infers SDF values on a tetrahedral mesh parameterizing a volume surrounding/containing the body. Option 1: Marching Tetrahedra is used to compute a triangle mesh from this implicit surface, and the resulting mesh is skinned to an estimated pose. Option 2: Alternatively, the tetrahedral mesh can be skinned first before using Marching Tetrahedra. During training (only), a normal map (corresponding to inferred pose and camera parameters) is generated for the predicted clothed body mesh and compared with the ground truth.

expected due to differences in the degrees of freedom.

The main goal of our work is to provide weak supervision during DNN training via 2D normal maps [11, 19, 29, 55, 70, 74, 88, 89]. By formulating the optimization problem with respect to image-based normals, we aim to better represent fundamental correlations between 2D images and 3D geometry in order to facilitate subsequent democratization to consumer-grade devices. Given (inferred) pose and camera parameters, a normal map can be computed from the skinned triangulated surface. See Figure 2. While we initialize the SDF network parameters with the data from [58], the resulting model will tend to overfit (since it is trained on a limited quantity of 3D data) and thus generalize poorly to in-the-wild 2D images. Thus, after this initialization, we train the model in a weakly supervised manner using only ground truth normal maps.

To summarize, our contributions:

- We illustrate that our network can be used to reconstruct 3D geometry from sparse multiview RGB data obtained with consumer-grade cameras (and no ground truth 3D labels).
- We compute the correct gradients to differentiate through Marching Tetrahedra via a Lagrangian formulation, which enables differentiable mesh generation and thus end-to-end training with both volumetric and surface-based energies.

- We present a differentiable image rasterizer that: (1) allows us to use normal maps for weak supervision and (2) can efficiently compute normal maps from triangle meshes with over 300k triangles during network training.
- We formulate regularization energies that coerce inferred implicit surfaces to: (1) resemble true SDFs and (2) be locally smooth.
- We formulate silhouette energies defined to enforce 3D boundary matching.

## 2. Related Work

### 2.1. Human Shape Estimation

Various works use parametric body models such as SMPL [47] to estimate human body pose and shape without clothing [21, 36, 46, 54]. While existing methods are able to generalize to in-the-wild images, the inferred body mesh is often quite different from the underlying body shape and does not capture clothing.

In order to reconstruct humans wearing clothing from a single image [78], template-based approaches either rely on parametric models [12, 32, 38, 96, 103] or use person-specific meshes [25, 93]. For instance, GTA [96] projects SMPL onto a learned 3D triplane representation, and [103] constructs local implicit fields centered around locations on the SMPL-X model [63]. Limitations of template-based ap-

proaches to clothed human reconstruction include the output being constrained by the topology of the template as well as a reliance on accurate pose estimation. Template-free methods typically leverage 2D signals or 3D geometric representations to recover geometry. In [19], reconstruction is achieved by generating front and back depth images that are later combined into a 3D surface. [26] builds on this idea and proposes a coarse-to-fine reconstruction method leveraging both predicted depth and normal images. Inspired by shape-from-silhouette techniques, SiCloPe [55] recovers geometry by predicting silhouette images and 3D joint positions. [81, 101] predict volumetric occupancy on a uniform voxel grid directly, while [17] proposed learning a Fourier subspace of 3D occupancy; in both cases, Marching Cubes can be used to generate a triangle mesh. PIFu [69] and PIFuHD [70] infer 3D shape with neural implicit functions sampled onto a grid. Follow-up work [9, 88] leverages predicted normal maps to improve depth inference. PAMIR [102] extends PIFuHD to increase generalizability by regularizing the implicit function using semantic features from a parametric model. ICON [88] and ECON [89] leverage inferred front and back normal maps as an intermediate encoding of 3D geometry, but these methods still rely on ground truth 3D scan data during training.

Instead of a single input image, other works aim to construct animatable avatars from a sparse set of cameras [27, 71, 98, 104], video [14, 18, 23, 29–31, 60, 76, 85, 94], depth [15, 83, 90, 100], point clouds [35, 49], 4D capture [28], or scans [42, 45, 72]. Most similar to our work, Self-Recon [29] uses normal maps inferred from PIFuHD [70] to supervise network training, and SeSDF [7] can either take as input a single image or uncalibrated multiview images

Moving towards improved generalization, weakly supervised methods have also been explored for human pose estimation [34, 57, 62, 91], human body shape [53, 95], and garment template reconstruction (on the SMPL body) [13, 53].

## 2.2. Differentiable Marching Cubes / Tetrahedra

A number of recent works have proposed methods for backpropagating through Marching Cubes [48] / Marching Tetrahedra [80] by either leveraging properties of a point-to-SDF network [66, 73] (e.g. as in DeepSDF [61]) or by training a DNN for mesh generation [44]. MeshSDF [66] builds on DeepSDF [61], where a network  $f_\eta$  is trained to infer an SDF value at a location  $x$  conditioned on a latent shape code  $\eta$ . In the forward step,  $f_\eta(x)$  is computed for every vertex of a fixed voxel grid so that Marching Cubes can be used to generate a triangle mesh. The authors postulate that a small increase in the SDF values would move a triangle vertex in the normal direction, which is only true idealistically when there are no shocks/rarefactions in the SDF isocontours (see e.g. [59]); moreover, Marching Cubes does

not move vertices in such a manner, even when it produces a consistent set of vertices under perturbation. Their assumptions also necessitate  $f_\eta$  being a true SDF, even though it is only an approximation. Follow-up work in [73] presents a similar formulation using Marching Tetrahedra. See [50] for more discussion on the problematic assumptions in [66].

## 3. Weak Supervision via Normal Maps

In the context of human digitization, it can be challenging to train generalizable ML-based models with full 3D supervision (see e.g. [58, 69, 70, 102]). Robust generalization typically necessitates access to a large number of ground truth training examples; however, publicly available datasets of 3D scan data for clothed human bodies remain scarce (in part, because it is both expensive and complicated to obtain). Even if such data were more readily available, existing works typically require unskinning the data to a reference pose (see e.g. [24, 25, 58, 86]), which can lead to various complications: tangling, self-intersection, inversion, etc.

To alleviate dependency on labeled 3D data, we propose a weakly supervised approach using 2D normal maps as ground truth labels (only) during training. A 2D normal map defines an RGB value for each pixel, corresponding to the (camera or world space) unit normal that best represents the geometry rasterized to that pixel. A normal map can be approximated by casting a ray through the pixel center and subsequently interpolating normals to the ray-geometry intersection point, although a better estimate would be obtained by supersampling (similar to the way pixel color is computed). Importantly, difficult to handle occluded regions (such as the armpit) may be ignored (in contrast to full 3D supervision). Since there are a number of ways to obtain ground truth normal maps (besides utilizing 3D scan data), this approach vastly increasing the amount of data available for training. For example, one can utilize RGBD images [2, 20, 33, 97], stereo pairs [3, 39], and/or neural networks (including NeRFs [51]) trained to produce normal maps from RGB images [70, 88, 89]. Increasing the amount of training data (in this way) facilitates generalization to a much more representative and diverse set of people in clothing (as compared to using only a limited number of 3D scans).

Given inferred SDF values  $\hat{\phi}_k$  on the (fixed) tetrahedral mesh vertices  $u_k$ , Marching Tetrahedra is used to uniquely generate a triangle mesh with vertices  $v_i(\hat{\phi}_k)$ . Given an inferred pose  $\hat{\theta}$  and camera parameters  $\hat{c}$ , a normal map  $N(v_i(\hat{\phi}_k), \hat{\theta}, \hat{c})$  can be generated. The objective function to be minimized is then

$$\mathcal{L}(\hat{\phi}_k, \hat{\theta}, \hat{c}) = \left\| N(v_i(\hat{\phi}_k), \hat{\theta}, \hat{c}) - N_{GT} \right\| \quad (1)$$

where  $N_{GT}$  is a ground truth normal map.

## 4. Tetrahedral Mesh Framework

The 3D space surrounding and including the human body is parameterized via a tetrahedral mesh. First, a Cartesian grid based level set representation is generated for the SMPL template body [47] in the star pose (similar to [58, 86]). Then, a constant value is subtracted from the SDF values in order to inflate the zero level set so that its interior can contain a wide range of clothed body shapes. Subsequently, a tetrahedral mesh is generated for this interior region using red/green refinement [52, 77]. See Figure 3.

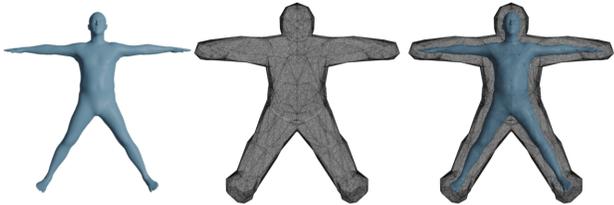


Figure 3. The tetrahedral mesh parameterizes a volume of air surrounding and including the body.

Given an input image, a DNN is trained to infer an implicit surface approximation to the clothed body, represented by SDF values  $\hat{\phi}_k$  on the tetrahedral mesh vertices  $u_k$  (similar to [58]). The DNN is composed of a CNN-based stacked hourglass encoder, followed by graph convolutional layers that progressively increase the resolution of the sampled SDF. This encodes an input image into a feature vector, which is then decoded into SDF values on the vertices of the tetrahedral mesh. Due to the large number of tetrahedral mesh vertices, the graph convolutional layers are only partially connected in order to significantly reduce memory usage.

### 4.1. Skinning

The tetrahedral mesh can be deformed via linear blend skinning (LBS) using per-tetrahedron-vertex, per-joint skinning weights  $w_{kj}$ . The skinning weights are assigned in the star pose by first finding the point on the SMPL template body mesh closest to each tetrahedral mesh vertex, and then barycentrically interpolating skinning weights to that point from the vertices of the SMPL template body mesh triangle that contains it.

Given pose parameters  $\theta$  with joint transformations  $T_j(\theta)$ , the skinned position of each tetrahedral mesh vertex is  $u_k(\theta) = \sum_j w_{kj} T_j(\theta) u_k^j$  where  $u_k^j$  is the location of  $u_k$  in the untransformed reference space of joint  $j$ . VIBE [37] is used to estimate the SMPL pose parameters  $\hat{\theta} \in \mathbb{R}^{72}$  for any given input image. During training, all layers of a pretrained VIBE model are frozen except for the final Gated Recurrent Unit (GRU) layer.

## 5. Marching Tetrahedra

Given SDF values  $\phi_k$  defined on tetrahedral mesh vertices  $u_k$ , Marching Tetrahedra can be implemented to compute a unique (non-ambiguous) triangle mesh with vertices  $v_i$ , making differentiation more straightforward as compared to the many cases (and non-uniqueness) that need to be considered for Marching Cubes. In order to avoid triangle vertices  $v_i$  coincident with a tetrahedron vertex  $u_k$ , all the  $\phi_k$  values are preprocessed (infinitesimally) changing those with  $|\phi_k| < \epsilon$  to  $\phi_k = \epsilon \text{sign}(\phi_k)$ , e.g. with  $\epsilon = 10^{-8}$ .

For each tetrahedron mesh edge  $e_i = \{u_{k_1}, u_{k_2}\}$  that includes a sign change, i.e.  $\text{sign}(\phi_{k_1}) \neq \text{sign}(\phi_{k_2})$ , a triangle vertex

$$v_i = \frac{-\phi_{k_2}}{\phi_{k_1} - \phi_{k_2}} u_{k_1} + \frac{\phi_{k_1}}{\phi_{k_1} - \phi_{k_2}} u_{k_2} \quad (2)$$

is defined using linear interpolation. Afterwards, triangles are constructed in a tetrahedron-by-tetrahedron manner by considering the two cases that can occur: either three edges of the tetrahedron contain triangle vertices and one triangle is constructed, or four edges contain triangle vertices and a quadrilateral is constructed and split into two triangles. Note that this typically arbitrary splitting of the quadrilateral can be made consistent for the sake of differentiation. Since the tetrahedral mesh does not change topology, the edges can be numbered in a fixed manner; then, one can consistently split a quadrilateral by connecting the triangle vertex on the lowest numbered edge to the triangle vertex on the highest numbered edge (or via a similar alternative strategy). The resulting triangle mesh is guaranteed to be watertight, and the vertices in each triangle are reordered (when necessary) to ensure that all face normals point outwards.

### 5.1. Ray-Tracing the Implicit Surface Directly

As an alternative to Marching Tetrahedra, consider casting a ray to find an intersection point with the implicit surface and subsequently using the normal vector defined (directly) by the implicit surface at that intersection point. A number of existing works consider such approaches in various ways, see e.g. [4, 11, 56, 82, 92]. Perturbations of the intersection point depend on perturbations of the  $\phi$  values on the vertices of the tetrahedron that the intersection point lies within. If a change in  $\phi$  values causes the intersection point to no longer be contained inside the tetrahedron, one would need to discontinuously jump to some other tetrahedron (which could be quite far away, if it even exists). A potential remedy for this would be to define a virtual implicit surface that extends out of the tetrahedron in a way that provides some sort of continuity (especially along silhouette boundaries).

Comparatively, our Marching Tetrahedra approach allows us to presume (for example) that the point of intersection remains fixed on the face of the triangle even as the

triangle moves. Since the implicit surface has no explicit parameterization, one is unable to similarly hold the intersection point fixed. The implicit surface utilizes an Eulerian point of view where the rays (which represent the discretization) are held fixed while the implicit surface moves (as  $\phi$  values change), in contrast to our Lagrangian discretization where the rays are allowed to move/bend in order to follow fixed intersection points during differentiation. A similar approach for an implicit surface would hold the intersection point inside the tetrahedron fixed even as  $\phi$  changes. Although such an approach holds potential due to the fact that implicit surfaces are amenable to computing derivatives off of the surface itself, the merging/pinching of isocontours created by convexity/concavity would likely lead to various difficulties. Furthermore, other issues would need to be addressed as well, e.g. the gradients (and thus normals) are only piecewise constant (and thus discontinuous) in the piecewise linear tetrahedral mesh basis.

## 5.2. Computing Gradients

According to Equation 2,

$$\frac{\partial v_i}{\partial(\phi_{k_1}, \phi_{k_2})} = \begin{bmatrix} \frac{\phi_{k_2}(u_{k_1} - u_{k_2})}{(\phi_{k_1} - \phi_{k_2})^2} & \frac{\phi_{k_1}(u_{k_2} - u_{k_1})}{(\phi_{k_1} - \phi_{k_2})^2} \end{bmatrix} \quad (3)$$

where dividing by  $(\phi_{k_1} - \phi_{k_2})^2$  can be problematic. The preprocess at the beginning of Section 5 guarantees that  $|\phi_{k_1} - \phi_{k_2}| \geq 2\epsilon$ , which means that the worst possible scenario for Equation 2 (when  $|\phi_{k_1}| = |\phi_{k_2}| = \epsilon$ ) still results in  $\mathcal{O}(1)$  coefficients for  $u_{k_1}$  and  $u_{k_2}$ ; however, the  $\phi$ -based coefficients in Equation 3 would be  $\mathcal{O}(1/\epsilon)$ . Thus, while  $\epsilon = 10^{-8}$  is sufficient for Equation 2, a larger value of  $\epsilon$  might be prudent when considering Equation 3.

## 5.3. Skinning

There are two options for the algorithm ordering between skinning and Marching Tetrahedra (the latter of which reverses the order in Figure 2). For skinning the triangle mesh, the skinned position of each triangle mesh vertex is  $v_i(\theta, \phi) = \sum_j w_{ij}(\phi) T_j(\theta) v_i^j(\phi)$  where  $v_i^j$  is the location of  $v_i$  in the untransformed reference space of joint  $j$ . Unlike in Section 4.1 where  $w_{kj}$  and  $u_k^j$  were fixed,  $w_{ij}$  and  $v_i^j$  both vary yielding three terms in the product rule.  $\partial v_i^j / \partial \phi$  is computed according to Equation 3, noting that  $u_{k_1}$  and  $u_{k_2}$  are fixed.  $w_{ij}(\phi)$  is defined similarly to Equation 2,

$$w_{ij} = \frac{-\phi_{k_2}}{\phi_{k_1} - \phi_{k_2}} w_{k_1j} + \frac{\phi_{k_1}}{\phi_{k_1} - \phi_{k_2}} w_{k_2j} \quad (4)$$

where  $w_{k_1j}$  and  $w_{k_2j}$  are fixed; similar to Equation 3,  $\partial w_{ij} / \partial \phi$  will contain  $\mathcal{O}(1/\epsilon)$  coefficients. For skinning the tetrahedral mesh, Equations 2 and 3 directly define  $v_i$  and  $\partial v_i / \partial \phi$  since the skinning is moved to the tetrahedral mesh vertices  $u_k$ . Then,  $\partial v_i / \partial u_k$  is computed according to Equation 2 in order to chain rule to skinning (i.e. to  $\partial u_k / \partial \theta$ ,

which is computed according to the equations in Section 4.1).

## 6. Image Rasterization

Given a skinned triangulated surface and parameters for a perspective camera model, a camera space normal map is computed using a right-handed coordinate system. We assume that the geometry is centered in the image, since images are cropped and rescaled during preprocessing. Normal maps made using different assumptions, or decoded and stored as RGB values, are readily transformed back into unit normals (in camera space) in order to match our assumptions.

### 6.1. Normals

Recall (from Section 5) that triangle vertices are reordered (if necessary) in order to obtain outward-pointing face normals. The area-weighted outward face normal is

$$n_f(v_1, v_2, v_3) = \frac{1}{2}(v_2 - v_1) \times (v_3 - v_1) \quad (5)$$

where

$$Area(v_1, v_2, v_3) = \frac{1}{2} \|(v_2 - v_1) \times (v_3 - v_1)\|_2 \quad (6)$$

is the area weighting. Area-averaged vertex unit normals  $\hat{n}_v$  are computed via

$$n_v = \sum_f n_f \quad \hat{n}_v = \frac{n_v}{\|n_v\|_2} \quad (7)$$

where  $f$  ranges over all the triangle faces that include vertex  $v$ . Note that one can drop the 1/2 in Equation 5, since it cancels out when computing  $\hat{n}_v$  in Equation 7.

### 6.2. Camera Model

The camera rotation and translation are used to transform each vertex  $v_g$  of the geometry to the camera view coordinate system (where the origin is located at the camera aperture), i.e.  $v_c = Rv_g + T$ . The normalized device coordinate system normalizes geometry in the viewing frustum (with  $z \in [n, f]$ ) so that all  $x, y \in [-1, 1]$  and all  $z \in [0, 1]$ . See Figure 4, left. Vertices are transformed into this coordinate system via

$$\begin{bmatrix} [v_{NDC}] z_c \\ z_c \end{bmatrix} = \begin{bmatrix} \frac{2n}{W} & 0 & 0 & 0 \\ 0 & \frac{2n}{H} & 0 & 0 \\ 0 & 0 & \frac{f}{f-n} & \frac{-fn}{f-n} \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} [v_c] \\ 1 \end{bmatrix} \quad (8)$$

where  $H = 2n \tan(\theta_{fov}/2)$  is the height of the image,  $\theta_{fov}$  is the field of view,  $W = Ha$  is the width of the image, and  $a$  is the aspect ratio. The screen coordinate system is

obtained by transforming the origin to the top left corner of the image, with  $+x$  pointing right and  $+y$  pointing down. See Figure 4, right. Vertices are transformed into this coordinate system via

$$\begin{bmatrix} [v'] \\ 1 \end{bmatrix} = \begin{bmatrix} -W/2 & 0 & 0 & W/2 \\ 0 & -H/2 & 0 & H/2 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} [v_{NDC}] \\ 1 \end{bmatrix} \quad (9)$$

or via

$$\begin{bmatrix} [v'] & z_c \\ z_c & \end{bmatrix} = \begin{bmatrix} -n & 0 & W/2 & 0 \\ 0 & -n & H/2 & 0 \\ 0 & 0 & \frac{f}{f-n} & \frac{-fn}{f-n} \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} [v_c] \\ 1 \end{bmatrix} \quad (10)$$

which is obtained by multiplying both sides of Equation 9 by  $z_c$  and substituting in Equation 8.

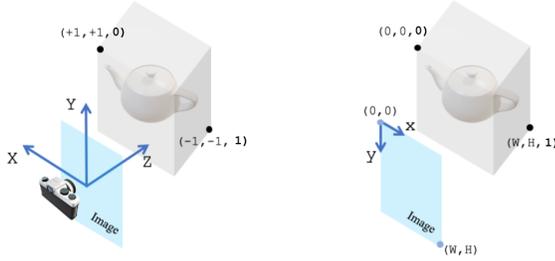


Figure 4. The normalized device (left) and screen (right) coordinate systems used during rasterization (based on Pytorch3D conventions<sup>1</sup>).

### 6.3. Normal Map

For each pixel, a ray is cast from the camera aperture through the pixel center to find its first intersection with the triangulated surface at a point  $p$  in world space. Denoting  $v_1, v_2, v_3$  as the vertices of the intersected triangle, barycentric weights for the intersection point

$$\begin{aligned} \hat{\alpha}_1 &= \frac{Area(p, v_2, v_3)}{Area(v_1, v_2, v_3)} \\ \hat{\alpha}_2 &= \frac{Area(v_1, p, v_3)}{Area(v_1, v_2, v_3)} \\ \hat{\alpha}_3 &= \frac{Area(v_1, v_2, p)}{Area(v_1, v_2, v_3)} \end{aligned} \quad (11)$$

are used to compute a rotated (into screen space) unit normal from the unrotated vertex unit normals (see Equation 7) via

$$\hat{n} = R \frac{\hat{\alpha}_1 \hat{n}_{v_1} + \hat{\alpha}_2 \hat{n}_{v_2} + \hat{\alpha}_3 \hat{n}_{v_3}}{\|\hat{\alpha}_1 \hat{n}_{v_1} + \hat{\alpha}_2 \hat{n}_{v_2} + \hat{\alpha}_3 \hat{n}_{v_3}\|} \quad (12)$$

for the normal map. Note that dropping the denominators in Equation 11 does not change  $\hat{n}$ .

<sup>1</sup><https://pytorch3d.org/docs/cameras>

### 6.4. Scanline Rendering

After projecting a visible triangle into the screen coordinate system (via Equation 10), its projected area can be computed as

$$Area2D(v'_1, v'_2, v'_3) = -\frac{1}{2} \det \begin{pmatrix} x'_2 - x'_1 & y'_2 - y'_1 \\ x'_3 - x'_1 & y'_3 - y'_1 \end{pmatrix} \quad (13)$$

similar to Equation 6 (where the negative sign accounts for the fact that visible triangles have normals pointing towards the camera). When a projected triangle overlaps a pixel center  $p'$ , barycentric weights for  $p'$  are computed by using  $Area2D$  instead of  $Area$  in Equation 11. Notably, un-normalized world space barycentric weights can be computed from un-normalized screen space barycentric weights via  $\alpha_1 = z'_2 z'_3 \alpha'_1$ ,  $\alpha_2 = z'_1 z'_3 \alpha'_2$ ,  $\alpha_3 = z'_1 z'_2 \alpha'_3$  or

$$\begin{aligned} \alpha_1 &= z'_2 z'_3 Area2D(p', v'_2, v'_3) \\ \alpha_2 &= z'_1 z'_3 Area2D(v'_1, p', v'_3) \\ \alpha_3 &= z'_1 z'_2 Area2D(v'_1, v'_2, p') \end{aligned} \quad (14)$$

giving

$$\hat{n} = R \frac{\alpha_1 \hat{n}_{v_1} + \alpha_2 \hat{n}_{v_2} + \alpha_3 \hat{n}_{v_3}}{\|\alpha_1 \hat{n}_{v_1} + \alpha_2 \hat{n}_{v_2} + \alpha_3 \hat{n}_{v_3}\|} \quad (15)$$

as an (efficient) alternative to Equation 12. If more than one triangle overlaps  $p'$ , the closest one (i.e. the one with the smallest value of  $z' = \hat{\alpha}'_1 z'_1 + \hat{\alpha}'_2 z'_2 + \hat{\alpha}'_3 z'_3$  at  $p'$ ) is chosen.

### 6.5. Computing Gradients

For each pixel overlapped by the triangle mesh, the derivative of the normal (in Equation 15) with respect to the vertices of the triangle mesh is required, i.e.  $\partial \alpha_i / \partial v_g$  and  $\partial \hat{n}_{v_i} / \partial v_g$  are required.  $\partial \alpha_i / \partial v'$  can be computed from Equations 14 and 13,  $\partial v' / \partial v_c$  can be computed from Equation 10, and  $\partial v_c / \partial v_g$  can be computed from  $v_c = R v_g + T$ .  $\partial \hat{n}_{v_i} / \partial v_g$  can be computed from Equations 7 and 5.

## 7. SDF Regularization

Two regularization terms are utilized during neural network training in order to encourage: (1) the inferred  $\hat{\phi}$  values to resemble a true SDF and (2) smoothness (similar to [43, 68]). Notably, the smoothness regularizer behaves significantly better when  $\hat{\phi}$  is closer to a true SDF.

### 7.1. Eikonal Regularization

Given a tetrahedron  $t$  with vertices  $u_k = (x_k, y_k, z_k)$  and inferred  $\hat{\phi}_k$  values,  $\hat{\phi}$  can be linearly approximated within the tetrahedron by writing

$$\hat{\phi}_k = ax_k + by_k + cz_k + d \quad (16)$$

for each of the four vertices; then, the resulting  $4 \times 4$  linear system of equations can be solved to obtain the unknown coefficients  $(a, b, c, d)$  leading to

$$|\nabla \hat{\phi}_t| = \sqrt{a^2 + b^2 + c^2} \quad (17)$$

as the norm of the gradient. Summing over tetrahedra leads to

$$E_{1a} = \frac{1}{2} \sum_t (|\nabla \hat{\phi}_t| - 1)^2 \quad (18)$$

as the energy to be minimized. The problem with Equation 18 (and similar approaches, such as [6, 22]) is that the chain rule moves the square root in Equation 17 to the denominator, potentially leading to NaNs/overflow; notably, even an exact SDF has  $|\nabla \phi| = 0$  at both extrema and pinching/merging saddles, and an inferred  $\hat{\phi}$  can have  $|\nabla \hat{\phi}| = 0$  elsewhere as well. This can be avoided by instead using

$$E_{1b} = \frac{1}{2} \sum_t (|\nabla \hat{\phi}_t|^2 - 1)^2 \quad (19)$$

which still enforces  $|\nabla \hat{\phi}_t| = 1$ ; alternatively,

$$E_{1c} = \frac{1}{2} \sum_t \text{Volume}(t) (|\nabla \hat{\phi}_t|^2 - 1)^2 \quad (20)$$

scales the penalty on each tetrahedron by its volume.

## 7.2. Motion by Mean Curvature

In order to encourage smoothness, we define an energy that when minimized results in motion by mean curvature. Following [10, 99], the surface area can be calculated via

$$\int_{\Omega} |\nabla H(\phi(x, y, z))| dV \quad (21)$$

where  $H$  is a Heaviside function and  $V$  is the volume; thus, on our tetrahedral mesh, we minimize

$$E_2 = \sum_t |\nabla H(\hat{\phi})| \text{Volume}(t) \quad (22)$$

using a smeared-out Heaviside Function

$$H(\hat{\phi}) = \begin{cases} 0 & \hat{\phi} < -\epsilon_H \\ \frac{1}{2} + \frac{\hat{\phi}}{2\epsilon_H} + \frac{1}{2\pi} \sin\left(\frac{\pi \hat{\phi}}{\epsilon_H}\right) & -\epsilon_H \leq \hat{\phi} \leq \epsilon_H \\ 1 & \hat{\phi} > \epsilon_H \end{cases} \quad (23)$$

where  $\epsilon_H$ , chosen as 1.5 times the average tetrahedral mesh edge length, determines the bandwidth of numerical smearing (see [75]).  $|\nabla H(\hat{\phi})|$  is discretized by linearly approximating  $H(\hat{\phi})$  in each tetrahedron along the lines of Equation 16 in order to obtain coefficients  $(a, b, c, d)$  for use in the equivalent of Equation 17. In order to avoid division by small numbers, we ignore tetrahedra with  $|\nabla H(\hat{\phi})| < 10^{-8}$  in Equation 22 reasoning that  $|\nabla H(\hat{\phi})|$  is small enough and thus  $\hat{\phi}$  is smooth enough in such tetrahedra.

## 8. Silhouette Losses

Instead of striving to make the inverse rendering differentiable at silhouette boundaries (as in e.g. [4]), we introduce energies that force the silhouettes to match.

### 8.1. Shrinking

For pixels that overlap the inferred surface but not the ground truth surface, the interior of the inferred surface needs to shrink so that the corresponding triangles disappear. For each tetrahedron mesh edge containing a vertex of a problematic triangle, the edge’s parent tetrahedral mesh vertices are added to the set  $U_{\text{shrink}}$  if they have negative SDF values; then,

$$\mathcal{L}_{\text{shrink}} = \frac{1}{2} \sum_{k \in U_{\text{shrink}}} (\hat{\phi}_k - \epsilon_s)^2 \quad (24)$$

encourages those negative  $\hat{\phi}_k$  values to target a positive  $\epsilon_s = 5 \times 10^{-3}$ , which is chosen as half the average tetrahedral mesh edge length.

### 8.2. Expanding

For pixels that overlap the ground truth surface but not the inferred surface, the interior of the inferred surface needs to expand. In order to determine where this expansion should occur, the implicit surface is temporarily inflated by changing the sign of the SDF at every tetrahedral mesh vertex with both  $\hat{\phi} > 0$  and a one-ring neighbor with  $\hat{\phi} < 0$  (e.g. by setting  $\hat{\phi}_{\text{temp}} = -\epsilon_s$  at those vertices). Next, the pixels that previously overlapped the ground truth surface but not the inferred surface and now overlap both the ground truth surface and the new inflated surface are identified. For each tetrahedron mesh edge containing a vertex of a triangle corresponding to one of these pixels, the edge’s parent tetrahedral mesh vertices are added to the set  $U_{\text{expand}}$  if they had positive SDF values before inflation. At this point, all of the temporary  $\hat{\phi}_{\text{temp}}$  values are discarded and the original  $\hat{\phi}$  values are restored. Then,

$$\mathcal{L}_{\text{expand}} = \frac{1}{2} \sum_{k \in U_{\text{expand}}} (\hat{\phi}_k + \epsilon_s)^2 \quad (25)$$

encourages the positive  $\hat{\phi}_k$  values to target  $-\epsilon_s$ .

## 9. Experiments

We first demonstrate (in Section 9.1) that our network has the ability to reconstruct clothed humans when ground truth camera parameters and normal maps are known. In Section 9.2, we demonstrate that the network can be trained to reconstruct 3D geometry with increasing efficacy as the number of sparse views increases. Subsequently (in Section 9.3), we extend this process to real-world RGB data



Figure 5. After training from 8 camera views, the input image in the first column results in the geometry shown in the second column. Note that the geometry is shown from novel views. For the sake of comparison, the ground truth geometry is shown from the same novel views. See also Figure 6.

(with no ground truth information) in order to demonstrate the ability to reconstruct 3D geometry using only network-inferred normal maps. For the sake of comparison, we also present (in Section 9.4) the results we obtained using available implementations of other methods for single view and multiview reconstruction.

### 9.1. Network Efficacy

Given ground truth 3D data from RenderPeople [67], we show that our network has the capacity and flexibility to reconstruct clothed humans from either a single image or multiple images. Regardless of the number of input images, the network is trained by minimizing the normal map loss (Equation 1), SDF regularization losses (Equations 20 and 22), and silhouette losses (Equations 24 and 25). In the multiview case, each image is considered individually (i.e. we treat multiview as a collection of single view examples). Figure 5 shows an example of the results obtained by training our network on 8 camera views surrounding the person (as compared to the ground truth).

### 9.2. Geometry Reconstruction

To quantitatively evaluate the accuracy of our inferred results, we define a normal map error as

$$e_{normal} = \frac{1}{W \times H} \sum_p \left( \frac{1}{2} (1 - \hat{n}_p \cdot n_p) \right)^2 \quad (26)$$

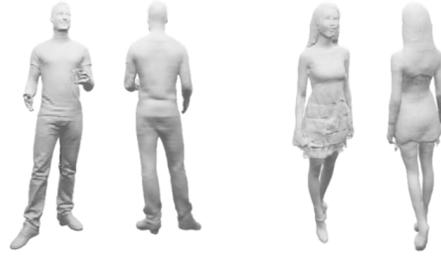


Figure 6. PIFuHD results, inferred using the input image in Figure 5 and shown from the same novel views (as in Figure 5). We stress that these results were obtained using inference from a single image, and so one would not expect the same efficacy (especially from novel views); however, these images do help to calibrate what one might expect from state-of-the-art inference. The conclusion is that our network has the ability to output high-quality reconstructed geometry.

where the ground truth and predicted normals at pixel  $p$  are  $n_p$  and  $\hat{n}_p$ , respectively, and  $\hat{n}_p \cdot n_p \in [-1, 1]$  is replaced with  $-1$  for pixels where the predicted and ground truth silhouettes do not overlap. Note that normal maps do not uniquely determine scale/depth; thus, the reconstructed objects could erroneously move closer/further from the camera becoming smaller/larger in scale (while also undergoing distortion, since this scale variance is not self-similar). In order to monitor this, we define a depth map error as

$$e_{depth} = \frac{1}{W \times H} \sum_p \left( \hat{d}_p - d_p \right)^2 \quad (27)$$

where  $(\hat{d}_p - d_p)$  is replaced with the thickness of the tetrahedral mesh (0.2 meters) for pixels where the predicted and ground truth silhouettes do not overlap.

Given ground truth 3D data from RenderPeople [67], we show how our network reconstructs 3D geometry with increasing efficacy as the number of sparse views increases. Figure 7 shows the inferred 3D geometry from a novel view, and Table 1 shows how per-pixel normal and depth errors decrease as the number of training views increases. When the network is trained on only one view, there are no constraints on the side/back of the person; hence, the predicted geometry has a high degree of noise when rendered from novel views. When trained with 5 views, the ground truth geometry is recovered with high accuracy.



Figure 7. Models trained on an increasing number of camera views (inference from a novel view is shown). Each set of images show: (left) predicted triangle mesh and its normal map, (middle) normal errors (blue is zero, red is max), and (right) depth errors. See also Figure 8.

# Views	Normals Error	Depth Error (m)
1	0.0317	0.0012
3	0.0299	0.0011
5	0.0060	0.0002

Table 1. Quantitative metrics for three unseen test views spaced between the training views.

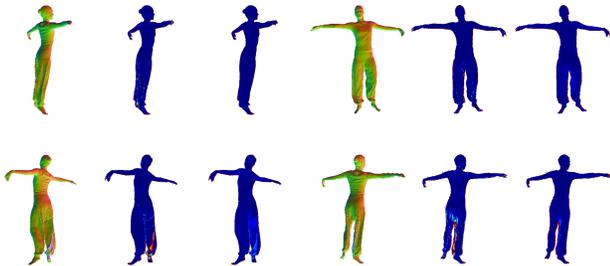


Figure 8. Using the 5-view network from Figure 7, we illustrate inference for four more novel views.

### 9.3. Geometry Reconstruction from RGB Images

Here, we illustrate that our network can be used to reconstruct 3D geometry from monocular uncalibrated RGB images, without requiring any pretraining on scanned data (or any other informed initialization of the network parameters). However, we do utilize a pretrained pix2pix network [84] (introduced in PIFuHD [70]) to infer ground truth normal maps and note that pix2pix was trained on 3D ground truth geometry. We do not consider this a severe limitation both because normal maps are easier to infer than 3D geometry and because there are other ways to obtain normal maps.

First, we captured monocular video footage of a person in a static pose; then, a sparse number of frames were extracted and preprocessed by removing the background using [8] and cropping to a square image. The resulting images were then passed into pix2pix to obtain “ground truth” normal maps. See Figure 9. Since estimated camera param-

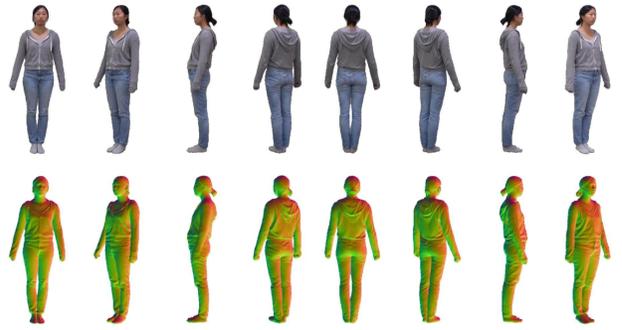


Figure 9. RGB images and corresponding normal maps. Odd views were used for training (starting with first column), and even views are used for testing.

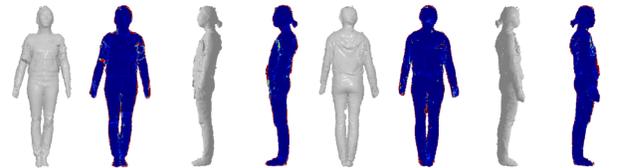


Figure 10. Final per-view meshes alongside their corresponding normal map errors (blue is zero, red is max).



Figure 11. Final triangle mesh shown from training (odd) and novel (even) views.

eters will be prone to error, we refine a rough initialization iteratively. At each iteration, we train the network and use Marching Tetrahedra to create a mesh inferred off of the image for (and overfit to) each view; then, we use ICP [5] to rigidly align all the meshes. Although one could delete all the triangles and remesh the point cloud, we obtained better results by updating each camera to match the ICP rigid transform of its corresponding mesh. The updated camera positions are then used to iteratively repeat the entire process. Once the camera parameters converge, the network can be trained with an additional loss that encourages 3D consistency. For a given camera view  $c_0$ , this loss is defined as

$$\mathcal{L}(\hat{\phi}_k) = \sum_{c \neq c_0} \left\| \hat{\phi}_k - \hat{\phi}_k(c) \right\| \quad (28)$$

where  $\hat{\phi}_k(c)$  refers to the inferred SDF values obtained from using view  $c$ 's image.

The network obtained from the aforementioned process (to improve camera extrinsics) will tend to be less detailed

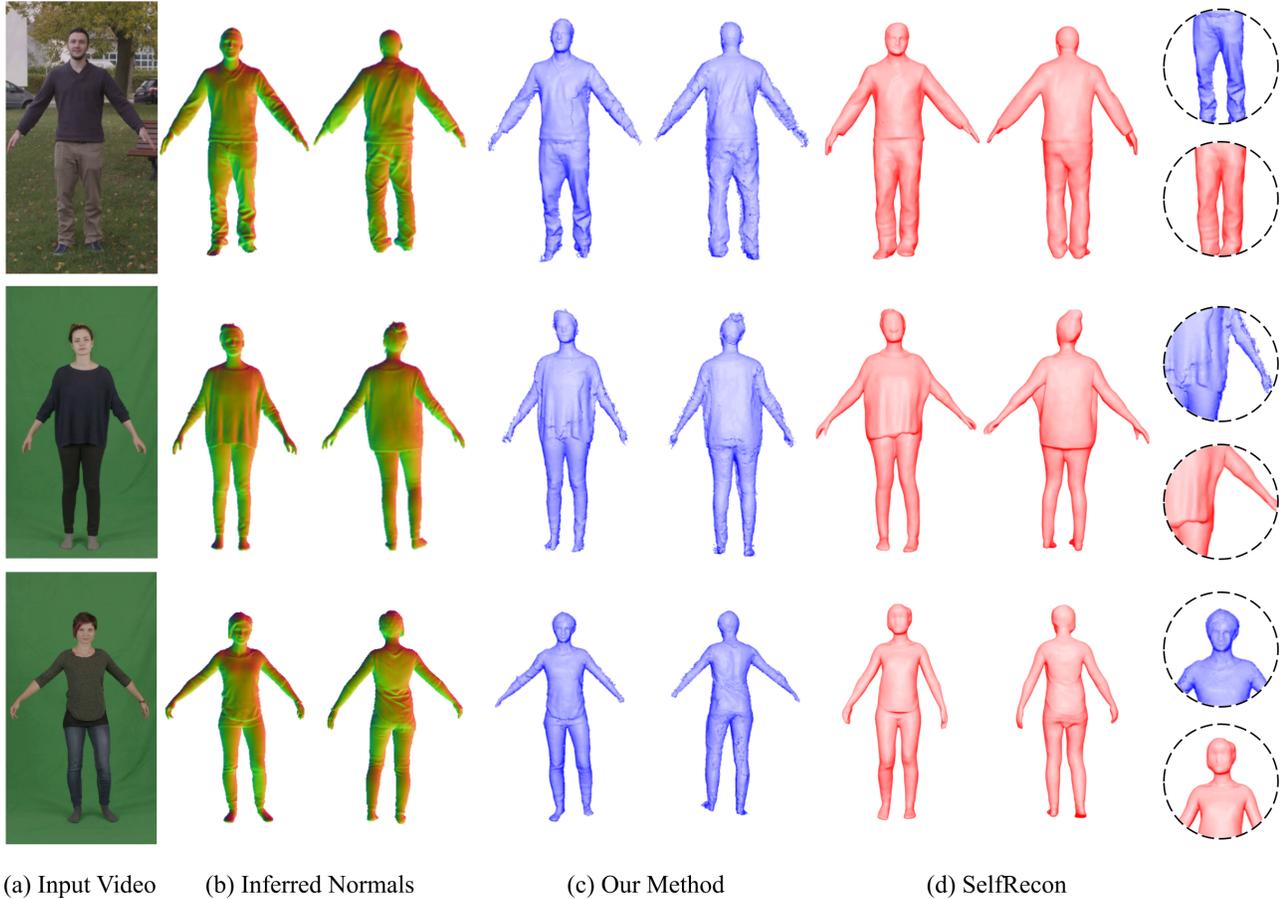


Figure 12. Predicted 3D geometry using our method (c) for videos from the PeopleSnapshot dataset (one frame is shown in (a)). The results from SelfRecon [29] are shown in (d). Note that the geometry is shown from novel views.

on the back side of the mesh, since only the front side can be seen in any given input image; thus, after improving camera extrinsics, we proceed as follows. Each view is fine-tuned with a regularizer that aims to keep  $\phi$  close to that which was obtained using Equation 28; then, we delete any visible triangles that are not consistent with the normal map (within some tolerance). See Figure 10. Since these are (actual) triangle meshes, it is trivial to load them into a suitable computer graphics application and align/resize the meshes in order to combine them into a single unified mesh. See Figure 11.

#### 9.4. Comparisons

We quantitatively compare our reconstruction method to existing single view [70, 89] and multiview [1, 29] reconstruction approaches using monocular videos from the People Snapshot Dataset [1]. Each video was captured with a fixed camera, and the subjects were asked to rotate while holding an A-pose. We trained our network on four frames per video (front, back, and two side views) and subsequently deleted

any visible triangles that are not consistent with the normal maps (within some tolerance, as in Section 9.3). SelfRecon [29] and VideoAvatar [1] were trained on all video frames. For the single view approaches [70, 89], we took the mesh predicted using the front or back-facing frame (whichever is closer to the test view) and scaled/rigidly aligned it using ICP to fit the corresponding SelfRecon mesh. Table 2 compares the results obtained with each method to the PIFuHD inferred normal map. SelfRecon has slightly more error and lacks detail compared to our approach (particularly around the face and wrinkles in the clothing). See Figure 12. Notably, the runtime of our approach on a single NVIDIA 3090 GPU is at least  $50\times$  faster than SelfRecon, which takes over a day of training (per video) to achieve their published results (our network is trained for about 20 minutes).

## 10. Conclusion

Although image-based reconstruction can be solved as an inverse problem, regularization is required in order to ad-

Method	Average Normals Error	STD
ECON [89]	0.1918	0.376
PIFuHD [70]	0.1680	0.358
VideoAvatar [1]	0.1342	0.325
SelfRecon [29]	0.0213	0.137
<b>Our Method</b>	<b>0.0207</b>	<b>0.111</b>

Table 2. Normal map errors (computed over predicted foreground pixels) corresponding to four test views, averaged over the examples shown in Figure 12.

dress issues with noise. Parameterized models (such as SMPL [47] or 3DMM [16]) provide for such regularization. We choose a neural network to parameterize our reconstruction where regularization is provided by having a limited number of network parameters. Our network aims to convert images from any view direction into a unique implicit surface, regardless of the view direction (similar in spirit to how the human brain process visual input); in fact, our eyes discern relative distance (similar to normal maps) more proficiently than they discern raw distance.

In summary, we present a weakly-supervised method for clothed human reconstruction by leveraging 2D normal maps as the supervisory signal during neural network training. In order to train a learned model that can infer high-frequency cloth and body geometry without any ground truth 3D data, our proposed approach builds on strong geometric priors for modeling and rendering. Our results reinforce the notion that less training data is required to train networks that infer normal maps than to train networks that infer 3D geometry (in agreement with ECON [89]). This means that working to improve the efficacy of network-inferred normal maps (and using the results for 3D reconstruction, as in Section 9.3) is likely to be more productive than working to obtain (via expensive 3D scanning) the excessive amount of ground truth data required to train a network to inference 3D geometry directly. Moreover, the process outlined in Section 9.3 provides an alternative mechanism (significantly cheaper than 3D scanning) for acquiring the ground truth data required to train a network to inference 3D geometry directly.

## Acknowledgements

Research supported in part by ONR N00014-19-1-2285, ONR N00014-21-1-2771. We would like to thank Reza and Behzad at ONR for supporting our efforts into machine learning. This work was also supported by JSPS KAKENHI Grant Number JP23H03439. J. W. was supported in part by the Gerald J. Lieberman Graduate Fellowship, the NSF Mathematical Sciences Postdoctoral Fellowship, and the UC President’s Postdoctoral Fellowship.

## References

- [1] Thiemo Alldieck, Marcus Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. Video based reconstruction of 3d people models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8387–8397, 2018. 10, 11
- [2] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *Proceedings of the IEEE Conference on computer Vision and Pattern Recognition*, pages 3686–3693, 2014. 3
- [3] Abhishek Badki, Alejandro Troccoli, Kihwan Kim, Jan Kautz, Pradeep Sen, and Orazio Gallo. Bi3d: Stereo depth estimation via binary classifications. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1600–1608, 2020. 3
- [4] Sai Praveen Bangaru, Michaël Gharbi, Fujun Luan, Tzu-Mao Li, Kalyan Sunkavalli, Milos Hasan, Sai Bi, Zexiang Xu, Gilbert Bernstein, and Fredo Durand. Differentiable rendering of neural sdfs through reparameterization. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022. 4, 7
- [5] Paul J Besl and Neil D McKay. Method for registration of 3-d shapes. In *Sensor fusion IV: control paradigms and data structures*, pages 586–606. Spie, 1992. 9
- [6] Hongrui Cai, Wanquan Feng, Xuetao Feng, Yan Wang, and Juyong Zhang. Neural surface reconstruction of dynamic scenes with monocular rgb-d camera. *arXiv preprint arXiv:2206.15258*, 2022. 7
- [7] Yukang Cao, Kai Han, and Kwan-Yee K Wong. Sesdf: Self-evolved signed distance field for implicit 3d clothed human reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4647–4657, 2023. 3
- [8] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 9
- [9] Kennard Yanting Chan, Guosheng Lin, Haiyu Zhao, and Weisi Lin. Integratedpifu: Integrated pixel aligned implicit function for single-view human reconstruction. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part II*, pages 328–344. Springer, 2022. 3
- [10] Tony F Chan and Luminita A Vese. Active contours without edges. *IEEE Transactions on image processing*, 10(2):266–277, 2001. 7
- [11] Xu Chen, Tianjian Jiang, Jie Song, Jinlong Yang, Michael J Black, Andreas Geiger, and Otmar Hilliges. gdna: Towards generative detailed neural avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20427–20437, 2022. 1, 2, 4
- [12] Enric Corona, Mihai Zanfir, Thiemo Alldieck, Eduard Gabriel Bazavan, Andrei Zanfir, and Cristian Sminchisescu. Structured 3d features for reconstructing

- reliable and animatable avatars. *arXiv preprint arXiv:2212.06820*, 2022. 2
- [13] Luca De Luigi, Ren Li, Benoît Guillard, Mathieu Salzmann, and Pascal Fua. Drapenet: Garment generation and self-supervised draping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1451–1460, 2023. 3
- [14] Junting Dong, Qi Fang, Yudong Guo, Sida Peng, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Totalsefscan: Learning full-body avatars from self-portrait videos of faces, hands, and bodies. *Advances in Neural Information Processing Systems*, 35:13654–13667, 2022. 3
- [15] Zheng Dong, Ke Xu, Ziheng Duan, Hujun Bao, Weiwei Xu, and Rynson Lau. Geometry-aware two-scale pifu representation for human reconstruction. *Advances in Neural Information Processing Systems*, 35:31130–31144, 2022. 3
- [16] Bernhard Egger, William AP Smith, Ayush Tewari, Stefanie Wuhler, Michael Zollhofer, Thabo Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, et al. 3d morphable face models: past, present, and future. *ACM Transactions on Graphics (TOG)*, 39(5):1–38, 2020. 11
- [17] Qiao Feng, Yebin Liu, Yu-Kun Lai, Jingyu Yang, and Kun Li. Fof: learning fourier occupancy field for monocular real-time human reconstruction. *Advances in Neural Information Processing Systems*, 35:7397–7409, 2022. 3
- [18] Yao Feng, Jinlong Yang, Marc Pollefeys, Michael J Black, and Timo Bolkart. Capturing and animation of body and clothing from monocular video. *arXiv preprint arXiv:2210.01868*, 2022. 3
- [19] Valentin Gabeur, Jean-Sébastien Franco, Xavier Martin, Cordelia Schmid, and Gregory Rogez. Moulding humans: Non-parametric 3d human shape estimation from single images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2232–2241, 2019. 2, 3
- [20] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3828–3838, 2019. 3
- [21] Shubham Goel, Georgios Pavlakos, Jathushan Rajasegaran, Angjoo Kanazawa\*, and Jitendra Malik\*. Humans in 4D: Reconstructing and tracking humans with transformers. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- [22] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. *arXiv preprint arXiv:2002.10099*, 2020. 7
- [23] Chen Guo, Tianjian Jiang, Xu Chen, Jie Song, and Otmar Hilliges. Vid2avatar: 3d avatar reconstruction from videos in the wild via self-supervised scene decomposition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12858–12868, 2023. 3
- [24] Marc Habermann, Weipeng Xu, Michael Zollhofer, Gerard Pons-Moll, and Christian Theobalt. Livecap: Real-time human performance capture from monocular video. *ACM Transactions on Graphics (TOG)*, 38(2):14, 2019. 3
- [25] Marc Habermann, Weipeng Xu, Michael Zollhofer, Gerard Pons-Moll, and Christian Theobalt. Deepcap: Monocular human performance capture using weak supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5052–5063, 2020. 2, 3
- [26] Sang-Hun Han, Min-Gyu Park, Ju Hong Yoon, Ju-Mi Kang, Young-Jae Park, and Hae-Gon Jeon. High-fidelity 3d human digitization from single 2k resolution images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12869–12879, 2023. 3
- [27] Yang Hong, Juyong Zhang, Boyi Jiang, Yudong Guo, Ligang Liu, and Hujun Bao. Stereopifu: Depth aware clothed human digitization via stereo vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 535–545, 2021. 3
- [28] Mustafa Işık, Martin Rünz, Markos Georgopoulos, Taras Khakhulin, Jonathan Starck, Lourdes Agapito, and Matthias Nießner. Humanrf: High-fidelity neural radiance fields for humans in motion. *arXiv preprint arXiv:2305.06356*, 2023. 3
- [29] Boyi Jiang, Yang Hong, Hujun Bao, and Juyong Zhang. Selfrecon: Self reconstruction your digital avatar from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5605–5615, 2022. 2, 3, 10, 11
- [30] Tianjian Jiang, Xu Chen, Jie Song, and Otmar Hilliges. Instantavatar: Learning avatars from monocular video in 60 seconds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16922–16932, 2023.
- [31] Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel, and Anurag Ranjan. Neuman: Neural human radiance field from a single video. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*, pages 402–418. Springer, 2022. 3
- [32] Sai Sagar Jinka, Astitva Srivastava, Chandradeep Pokhariya, Avinash Sharma, and PJ Narayanan. Sharp: Shape-aware reconstruction of people in loose clothing. *International Journal of Computer Vision*, 131(4):918–937, 2023. 2
- [33] Jinyoung Jun, Jae-Han Lee, Chul Lee, and Chang-Su Kim. Monocular human depth estimation via pose estimation. *IEEE Access*, 9:151444–151457, 2021. 3
- [34] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7122–7131, 2018. 3
- [35] Hyomin Kim, Hyeonsoo Nam, Jungeon Kim, Jaesik Park, and Seungyong Lee. Laplacianfusion: Detailed 3d clothed-human body reconstruction. *ACM Transactions on Graphics (TOG)*, 41(6):1–14, 2022. 3
- [36] Jeonghwan Kim, Mi-Gyeong Gwon, Hyunwoo Park, Hyukmin Kwon, Gi-Mun Um, and Wonjun Kim. Sampling is matter: Point-guided 3d human mesh reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vi-*

- sion and Pattern Recognition*, pages 12880–12889, 2023. [2](#)
- [37] Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. Vibe: Video inference for human body pose and shape estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5253–5263, 2020. [4](#)
- [38] Nikos Kolotouros, Georgios Pavlakos, Dinesh Jayaraman, and Kostas Daniilidis. Probabilistic modeling for human mesh recovery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11605–11614, 2021. [2](#)
- [39] Uday Kusupati, Shuo Cheng, Rui Chen, and Hao Su. Normal assisted stereo depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2189–2199, 2020. [3](#)
- [40] Minjae Lee, David Hyde, Michael Bao, and Ronald Fedkiw. A skinned tetrahedral mesh for hair animation and hair-water interaction. *IEEE transactions on visualization and computer graphics*, 25(3):1449–1459, 2018. [1](#)
- [41] Minjae Lee, David Hyde, Kevin Li, and Ronald Fedkiw. A robust volume conserving method for character-water interaction. In *Proceedings of the 18th annual ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pages 1–12, 2019. [1](#)
- [42] Zhe Li, Zerong Zheng, Hongwen Zhang, Chaonan Ji, and Yebin Liu. Avatarcap: Animatable avatar conditioned monocular human volumetric capture. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part I*, pages 322–341. Springer, 2022. [3](#)
- [43] Zhaoshuo Li, Thomas Müller, Alex Evans, Russell H Taylor, Mathias Unberath, Ming-Yu Liu, and Chen-Hsuan Lin. Neuralangelo: High-fidelity neural surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8456–8465, 2023. [6](#)
- [44] Yiyi Liao, Simon Donne, and Andreas Geiger. Deep marching cubes: Learning explicit surface representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2916–2925, 2018. [3](#)
- [45] Siyou Lin, Hongwen Zhang, Zerong Zheng, Ruizhi Shao, and Yebin Liu. Learning implicit templates for point-based clothed human modeling. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part III*, pages 210–228. Springer, 2022. [3](#)
- [46] Wu Liu, Qian Bao, Yu Sun, and Tao Mei. Recent advances of monocular 2d and 3d human pose estimation: a deep learning perspective. *ACM Computing Surveys*, 55(4):1–41, 2022. [2](#)
- [47] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):1–16, 2015. [2](#), [4](#), [11](#)
- [48] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. *ACM siggraph computer graphics*, 21(4):163–169, 1987. [1](#), [3](#)
- [49] Qianli Ma, Jinlong Yang, Siyu Tang, and Michael J Black. The power of points for modeling humans in clothing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10974–10984, 2021. [3](#)
- [50] Ishit Mehta, Manmohan Chandraker, and Ravi Ramamoorthi. A level set theory for neural implicit evolution under explicit flows. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part II*, pages 711–729. Springer, 2022. [1](#), [3](#)
- [51] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020. [1](#), [3](#)
- [52] Neil Molino, Robert Bridson, Joseph Teran, and Ronald Fedkiw. A crystalline, red green strategy for meshing highly deformable objects with tetrahedra. In *IMR*, pages 103–114. Citeseer, 2003. [4](#)
- [53] Gyeongsik Moon, Hyeongjin Nam, Takaaki Shiratori, and Kyoung Mu Lee. 3d clothed human reconstruction in the wild. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part II*, pages 184–200. Springer, 2022. [3](#)
- [54] Hyeongjin Nam, Daniel Sungho Jung, Yeonguk Oh, and Kyoung Mu Lee. Cyclic test-time adaptation on monocular video for 3d human mesh reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14829–14839, 2023. [2](#)
- [55] Ryota Natsume, Shunsuke Saito, Zeng Huang, Weikai Chen, Chongyang Ma, Hao Li, and Shigeo Morishima. Siclope: Silhouette-based clothed people. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4480–4490, 2019. [2](#), [3](#)
- [56] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3504–3515, 2020. [1](#), [4](#)
- [57] Mohamed Omran, Christoph Lassner, Gerard Pons-Moll, Peter Gehler, and Bernt Schiele. Neural body fitting: Unifying deep learning and model based human pose and shape estimation. In *2018 international conference on 3D vision (3DV)*, pages 484–494. IEEE, 2018. [3](#)
- [58] Hayato Onizuka, Zehra Hayirci, Diego Thomas, Akihiro Sugimoto, Hideaki Uchiyama, and Rin-ichiro Taniguchi. Tetratsdf: 3d human reconstruction from a single image with a tetrahedral outer shell. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6011–6020, 2020. [1](#), [2](#), [3](#), [4](#)
- [59] Stanley Osher, Ronald Fedkiw, and K Piechor. Level set methods and dynamic implicit surfaces. *Appl. Mech. Rev.*, 57(3):B15–B15, 2004. [3](#)
- [60] Anqi Pang, Xin Chen, Haimin Luo, Minye Wu, Jingyi Yu, and Lan Xu. Few-shot neural human performance rendering from sparse rgbd videos. *arXiv preprint arXiv:2107.06505*, 2021. [3](#)

- [61] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019. 3
- [62] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3d human pose and shape from a single color image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 459–468, 2018. 3
- [63] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10975–10985, 2019. 2
- [64] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9054–9063, 2021. 1
- [65] Gerard Pons-Moll, Sergi Pujades, Sonny Hu, and Michael J Black. ClothCap: Seamless 4d clothing capture and retargeting. *ACM Transactions on Graphics (ToG)*, 36(4):1–15, 2017. 1
- [66] Edoardo Remelli, Artem Lukoianov, Stephan Richter, Benoît Guillard, Timur Bagautdinov, Pierre Baque, and Pascal Fua. MeshSDF: Differentiable iso-surface extraction. *Advances in Neural Information Processing Systems*, 33: 22468–22478, 2020. 1, 3
- [67] RenderPeople. Renderpeople, 2018. 8
- [68] Radu Alexandru Rosu and Sven Behnke. PermutoSDF: Fast multi-view reconstruction with implicit surfaces using permutohedral lattices. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8466–8475, 2023. 6
- [69] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2304–2314, 2019. 3
- [70] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. PifuHD: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 84–93, 2020. 1, 2, 3, 9, 10, 11
- [71] Ruizhi Shao, Zerong Zheng, Hongwen Zhang, Jingxiang Sun, and Yebin Liu. DiffStereo: High quality human reconstruction via diffusion-based stereo using sparse cameras. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*, pages 702–720. Springer, 2022. 3
- [72] Kaiyue Shen, Chen Guo, Manuel Kaufmann, Juan Jose Zarate, Julien Valentin, Jie Song, and Otmar Hilliges. X-avatar: Expressive human avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16911–16921, 2023. 3
- [73] Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis. *Advances in Neural Information Processing Systems*, 34, 2021. 3
- [74] David Smith, Matthew Loper, Xiaochen Hu, Paris Mavroidis, and Javier Romero. Facsimile: Fast and accurate scans from an image in less than a second. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5330–5339, 2019. 2
- [75] Mark Sussman, Peter Smereka, and Stanley Osher. A level set approach for computing solutions to incompressible two-phase flow. *Journal of Computational physics*, 114(1):146–159, 1994. 7
- [76] Gusi Te, Xiu Li, Xiao Li, Jinglu Wang, Wei Hu, and Yan Lu. Neural capture of animatable 3d human from monocular video. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VI*, pages 275–291. Springer, 2022. 3
- [77] Joseph Teran, Neil Molino, Ronald Fedkiw, and Robert Bridson. Adaptive physics based tetrahedral mesh generation using level sets. *Engineering with computers*, 21(1): 2–18, 2005. 4
- [78] Yating Tian, Hongwen Zhang, Yebin Liu, and Limin Wang. Recovering 3d human mesh from monocular images: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 2
- [79] Garvita Tiwari, Nikolaos Sarafianos, Tony Tung, and Gerard Pons-Moll. Neural-gif: Neural generalized implicit functions for animating people in clothing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11708–11718, 2021. 1
- [80] Graham M Treece, Richard W Prager, and Andrew H Gee. Regularised marching tetrahedra: improved iso-surface extraction. *Computers & Graphics*, 23(4):583–598, 1999. 3
- [81] Gul Varol, Duygu Ceylan, Bryan Russell, Jimei Yang, Ersin Yumer, Ivan Laptev, and Cordelia Schmid. BodyNet: Volumetric inference of 3d human body shapes. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 20–36, 2018. 3
- [82] Delio Vicini, Sébastien Speierer, and Wenzel Jakob. Differentiable signed distance function rendering. *ACM Transactions on Graphics (TOG)*, 41(4):1–18, 2022. 1, 4
- [83] Shaofei Wang, Marko Mihajlovic, Qianli Ma, Andreas Geiger, and Siyu Tang. MetaAvatar: Learning animatable clothed human models from few depth images. *Advances in Neural Information Processing Systems*, 34:2810–2822, 2021. 3
- [84] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional GANs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 9
- [85] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-Shlizerman. Humanerf: Free-viewpoint rendering of moving people from

- monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16210–16220, 2022. 3
- [86] Jane Wu, Zhenglin Geng, Hui Zhou, and Ronald Fedkiw. Skinning a parameterization of three-dimensional space for neural network cloth. *arXiv preprint arXiv:2006.04874*, 2020. 1, 3, 4
- [87] Donglai Xiang, Fabian Prada, Timur Bagautdinov, Weipeng Xu, Yuan Dong, He Wen, Jessica Hodgins, and Chenglei Wu. Modeling clothing as a separate layer for an animatable human avatar. *ACM Transactions on Graphics (TOG)*, 40(6):1–15, 2021. 1
- [88] Yuliang Xiu, Jinlong Yang, Dimitrios Tzionas, and Michael J Black. Icon: Implicit clothed humans obtained from normals. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13296–13306, 2022. 1, 2, 3
- [89] Yuliang Xiu, Jinlong Yang, Xu Cao, Dimitrios Tzionas, and Michael J Black. Econ: Explicit clothed humans optimized via normal integration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 512–523, 2023. 1, 2, 3, 10, 11
- [90] Yuxuan Xue, Bharat Lal Bhatnagar, Riccardo Marin, Nikolaos Sarafianos, Yuanlu Xu, Gerard Pons-Moll, and Tony Tung. Nsf: Neural surface fields for human modeling from monocular depth. In *ICCV*, 2023. 3
- [91] Wei Yang, Wanli Ouyang, Xiaolong Wang, Jimmy Ren, Hongsheng Li, and Xiaogang Wang. 3d human pose estimation in the wild by adversarial learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5255–5264, 2018. 3
- [92] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33:2492–2502, 2020. 1, 4
- [93] Tao Yu, Zerong Zheng, Yuan Zhong, Jianhui Zhao, Qionghai Dai, Gerard Pons-Moll, and Yebin Liu. Simulcap: Single-view human performance capture with cloth simulation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2
- [94] Zhengming Yu, Wei Cheng, Xian Liu, Wayne Wu, and Kwan-Yee Lin. Monohuman: Animatable human neural field from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16943–16953, 2023. 3
- [95] Andrei Zanfir, Eduard Gabriel Bazavan, Hongyi Xu, William T Freeman, Rahul Sukthankar, and Cristian Sminchisescu. Weakly supervised 3d human pose and shape reconstruction with normalizing flows. In *European Conference on Computer Vision*, pages 465–481. Springer, 2020. 3
- [96] Zechuan Zhang, Li Sun, Zongxin Yang, Ling Chen, and Yi Yang. Global-correlated 3d-decoupling transformer for clothed avatar reconstruction. *arXiv preprint arXiv:2309.13524*, 2023. 2
- [97] Chaoqiang Zhao, Qiyu Sun, Chongzhen Zhang, Yang Tang, and Feng Qian. Monocular depth estimation based on deep learning: An overview. *Science China Technological Sciences*, 63(9):1612–1627, 2020. 3
- [98] Fuqiang Zhao, Wei Yang, Jiakai Zhang, Pei Lin, Yingliang Zhang, Jingyi Yu, and Lan Xu. Humannerf: Efficiently generated human radiance field from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7743–7753, 2022. 3
- [99] Hong-Kai Zhao, Tony Chan, Barry Merriman, and Stanley Osher. A variational level set approach to multiphase motion. *Journal of computational physics*, 127(1):179–195, 1996. 7
- [100] Ruichen Zheng, Peng Li, Haoqian Wang, and Tao Yu. Learning visibility field for detailed 3d human reconstruction and relighting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 216–226, 2023. 3
- [101] Zerong Zheng, Tao Yu, Yixuan Wei, Qionghai Dai, and Yebin Liu. Deephuman: 3d human reconstruction from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7739–7749, 2019. 3
- [102] Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai. Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. *IEEE transactions on pattern analysis and machine intelligence*, 2021. 1, 3
- [103] Zerong Zheng, Xiaochen Zhao, Hongwen Zhang, Boning Liu, and Yebin Liu. Avatarrex: Real-time expressive full-body avatars. *arXiv preprint arXiv:2305.04789*, 2023. 2
- [104] Tiansong Zhou, Tao Yu, Ruizhi Shao, and Kun Li. Hdhuman: High-quality human performance capture with sparse views. *CoRR*, 2022. 3