

# Capturing and Animation of Body and Clothing from Monocular Video

YAO FENG, Max Planck Institute for Intelligent Systems, Germany and ETH Zürich, Switzerland

JINLONG YANG, Max Planck Institute for Intelligent Systems, Germany

MARC POLLEFEYS, ETH Zürich, Switzerland

MICHAEL J. BLACK, Max Planck Institute for Intelligent Systems, Germany

TIMO BOLKART, Max Planck Institute for Intelligent Systems, Germany

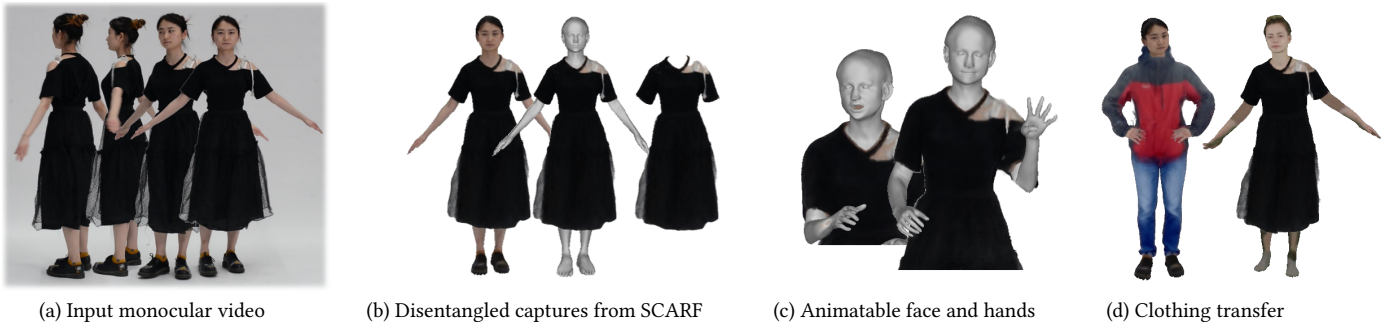


Fig. 1. Given a monocular video (a), our method (SCARF) builds an avatar where the body and clothing are disentangled (b). The body is represented by a traditional mesh, while the clothing is captured by an implicit neural representation. SCARF enables animation with detailed control over the face and hands (c) as well as clothing transfer between subjects (d).

While recent work has shown progress on extracting clothed 3D human avatars from a single image, video, or a set of 3D scans, several limitations remain. Most methods use a holistic representation to jointly model the body and clothing, which means that the clothing and body cannot be separated for applications like virtual try-on. Other methods separately model the body and clothing, but they require training from a large set of 3D clothed human meshes obtained from 3D/4D scanners or physics simulations. Our insight is that the body and clothing have different modeling requirements. While the body is well represented by a mesh-based parametric 3D model, implicit representations and neural radiance fields are better suited to capturing the large variety in shape and appearance present in clothing. Building on this insight, we propose SCARF (Segmented Clothed Avatar Radiance Field), a hybrid model combining a mesh-based body with a neural radiance field. Integrating the mesh into the volumetric rendering in combination with a differentiable rasterizer enables us to optimize SCARF directly from monocular videos, without any 3D supervision. The hybrid modeling enables SCARF to (i) animate the clothed body avatar by changing body poses (including hand articulation and facial expressions), (ii) synthesize novel views of the avatar, and (iii) transfer clothing between avatars in virtual try-on applications. We demonstrate that SCARF reconstructs clothing with higher visual quality than existing methods, that the clothing deforms with changing body pose and body shape, and that clothing can be successfully transferred between avatars of different subjects. The code and models are available at <https://github.com/YadiraF/SCARF>.

CCS Concepts: • **Computing methodologies** → *Shape modeling*.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

SA '22 Conference Papers, December 6–9, 2022, Daegu, Republic of Korea

© 2022 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-9470-3/22/12.

<https://doi.org/10.1145/3550469.3555423>

## ACM Reference Format:

Yao Feng, Jinlong Yang, Marc Pollefeys, Michael J. Black, and Timo Bolkart. 2022. Capturing and Animation of Body and Clothing from Monocular Video. In *SIGGRAPH Asia 2022 Conference Papers (SA '22 Conference Papers)*, December 6–9, 2022, Daegu, Republic of Korea. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3550469.3555423>

## 1 INTRODUCTION

Realistic avatar creation is one of the key enablers of the metaverse, and it supports many applications in virtual presence, fitness, digital fashion, and entertainment. Traditional ways to build avatars require either complex capture systems or manual design by artists, both of which are time-consuming and inefficient for large-scale avatar creation. To address this, previous work explores more practical ways to create avatars directly from single RGB images or monocular videos, which are more accessible to consumers.

The majority of work (e.g., [Choutas et al. 2020; Feng et al. 2021a; Kanazawa et al. 2018; Kolotouros et al. 2019; Pavlakos et al. 2019; Rong et al. 2021; Zanfir et al. 2021]) creates 3D human body avatars from images by estimating parameters of statistical 3D mesh models such as SCAPE [Angelov et al. 2005], Adam [Joo et al. 2018], SMPL/SMPL-X [Loper et al. 2015; Pavlakos et al. 2019], GHUM [Xu et al. 2020], or STAR [Osman et al. 2020], or implicit surface models like imGHUM [Alldieck et al. 2021] and LEAP [Mihajlovic et al. 2021]. As these models are trained from minimally clothed body scans, they are unable to capture clothing shape and appearance variations, which require a more flexible representation.

Methods that recover clothed bodies from images are instead trained with a large set of 3D clothed human scans [Saito et al. 2019, 2020; Xiu et al. 2022], or optimize the clothed avatar directly from multi-view images or videos [Chen et al. 2021b; Jiang et al. 2022; Liu

et al. 2021b; Peng et al. 2021a, 2022, 2021b; Xu et al. 2021]. To handle the complex topology of different clothing types, these methods model the body and clothing with a holistic implicit representation. Hence, hands and faces are typically poorly reconstructed and are not articulated. Additionally, holistic models of the body and clothing do not permit virtual try-on applications, which require the body and clothing to be represented separately. While neural radiance fields (NeRF) is able to model the head well (e.g., [Hong et al. 2022]), it remains unclear how to effectively combine such a part-based model with a clothed body representation.

Some methods treat the body and clothing separately with a layered representation, where clothing is modeled as a layer on top of the body [Corona et al. 2021; Jiang et al. 2020; Xiang et al. 2021; Zhu et al. 2020]. These methods require large datasets of 3D clothing scans for training, but still lack generalization to diverse clothing types. Furthermore, given an RGB image, they recover only the geometry of the clothed body without appearance information [Corona et al. 2021; Jiang et al. 2020; Zhu et al. 2020]. Similarly, Xiang et al. [2021] require multi-view video data and accurately registered 3D clothing meshes to build a subject-specific avatar; their method is not applicable to loose clothing like skirts or dresses.

Our goal is to go beyond existing work to capture realistic avatars from monocular videos that have detailed and animatable hands and faces as well as clothing that can be easily transferred between avatars. We observe that the body and clothing have different modeling requirements. Human bodies have similar shapes that can be modeled well by a statistical mesh model. In contrast, clothing shape and appearance are much more varied, thus require more flexible 3D representations that could handle changing topologies and transparent materials. With these observations, we propose SCARF (Segmented Clothed Avatar Radiance Field), a hybrid representation combining a mesh with a NeRF, to capture disentangled clothed human avatars from monocular videos. Specifically, we use SMPL-X to represent the human body and a NeRF on top of the body mesh to capture clothing of varied topology. There are four main challenges in building such a model from monocular video. First, SCARF must accurately capture human motion in monocular video and relate the body motion to the clothing. The NeRF is modeled in canonical space, and we use the skinning transformation from the SMPL-X body model to deform points in observation space to the canonical space. This requires accurate estimates of body shape and pose for every video frame. We estimate body pose and shape parameters with PIXIE [Feng et al. 2021a]. However, these estimates are not accurate enough, resulting in blurry reconstructions. Thus, we refine the body pose and shape during optimization. Second, the cloth deformations are not fully explained by the SMPL-X skinning, particularly in the presence of loose clothing. To overcome this, we learn a non-rigid deformation field to correct clothing deviations from the body. Third, SCARF’s hybrid representation, combining a NeRF and a mesh, requires customized volumetric rendering. Specifically, rendering the clothed body must account for the occlusions between the body mesh and the clothing layer. To integrate a mesh into volume rendering, we sample a ray from the camera’s optical center until it intersects the body mesh, and accumulate the colors along the ray up to the intersection point with the colored mesh

surface. Fourth, to disentangle the body and clothing, we must prevent the NeRF from capturing all image information including the body. To that end, we use clothing segmentation masks to penalize the NeRF outside of clothed regions.

In summary, SCARF automatically creates a 3D clothed human avatar from monocular video (Fig. 1) with disentangled clothing on top of the human body. SCARF offers the best of two worlds by combining different representations – a 3D parametric model for the body and a NeRF for the clothing. Based on SMPL-X, the reconstructed avatar offers animator control over body shape, pose, hand articulation, and facial expression. Since SCARF factors clothing from the body, the clothing can be extracted and transferred between avatars, enabling applications such as virtual try-on.

## 2 RELATED WORK

**3D Bodies from images.** The 3D surface of a human body is typically represented by a learned statistical 3D model [Alldieck et al. 2021; Anguelov et al. 2005; Joo et al. 2018; Loper et al. 2015; Osman et al. 2020; Pavlakos et al. 2019; Xu et al. 2020]. Numerous optimization and regression methods have been proposed to compute 3D shape and pose parameters from images, videos, and scans. See [Liu et al. 2021a; Tian et al. 2022] for recent surveys. We focus on methods that capture full-body pose and shape, including the hands and facial expressions [Choutas et al. 2020; Feng et al. 2021a; Pavlakos et al. 2019; Rong et al. 2021; Xiang et al. 2019; Xu et al. 2020; Zhou et al. 2021]. Such methods, however, do not capture hair, clothing, or anything that deviates the body. Also, they rarely recover texture information, due to the large geometric discrepancy between the clothed human in the image and captured minimal clothed body mesh. Unlike these prior works, we consider clothing as an important component and capture both the parametric body and non-parametric clothing from monocular videos.

**Capturing clothed humans from images.** Clothing is more complex than the body in terms of geometry, non-rigid deformation, and appearance, making the capture of clothing from images challenging. Mesh-based methods to capture clothing often use additional vertex offsets relative to the body mesh [Alldieck et al. 2019a, 2018a,b, 2019b; Jin et al. 2020; Lazova et al. 2019; Ma et al. 2020a,b]. While such an approach works well for clothing that is similar to the body, it does not capture clothing of varied topology like skirts and dresses. To handle clothing shape variations, recent methods exploit non-parametric models. For example, [He et al. 2021; Huang et al. 2020; Saito et al. 2019, 2020; Xiu et al. 2022; Zheng et al. 2021] extract pixel-aligned spatial features from images and map them to an implicit shape representation. To animate the captured non-parametric clothed humans, Yang et al. [2021] predict skeleton and skinning weights from images to drive the representation. Although such non-parametric models can capture various clothing styles much better than mesh-based approaches, faces and hands are usually poorly recovered due to the lack of a strong prior on how the human body should be. In addition, such approaches typically require a large set of manually cleaned 3D scans as training data. Recently, various methods recover 3D clothed humans directly from multi-view or monocular RGB videos [Chen et al. 2021b; Jiang et al. 2022; Liu et al. 2021b; Peng et al. 2021a, 2022, 2021b; Su et al. 2021;

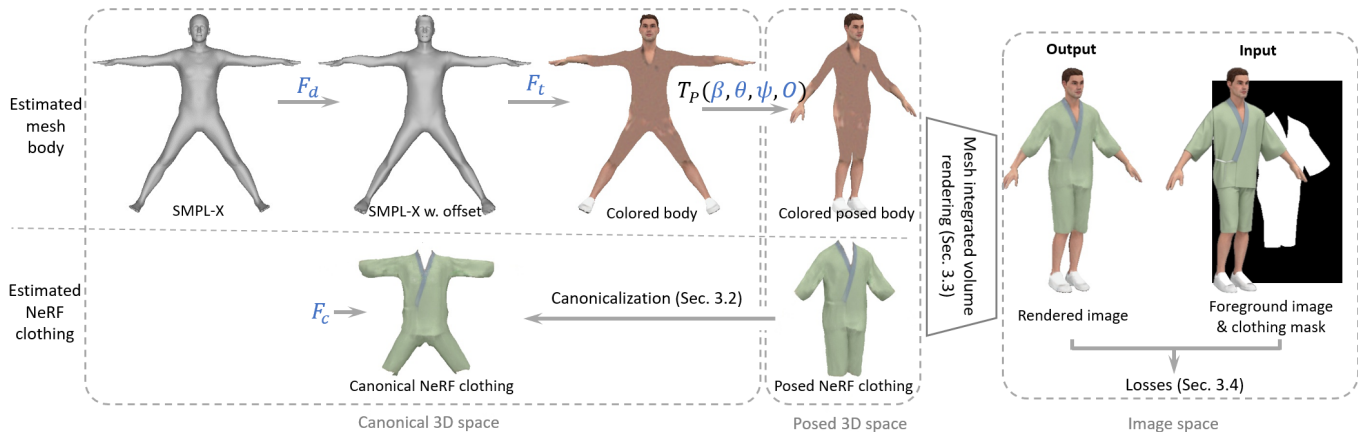


Fig. 2. SCARF takes monocular RGB video and clothing segmentation masks as input, and outputs a human avatar with separate body and clothing layers. Blue letters indicate optimizable modules or parameters.

Weng et al. 2022]. They optimize avatars from image information using implicit shape rendering [Liu et al. 2020; Niemeyer et al. 2020; Yariv et al. 2021, 2020] or volume rendering [Mildenhall et al. 2020], no 3D scans are needed. Although these approaches demonstrate impressive performance, hand gestures and facial expressions are difficult to capture and animate due to the lack of model expressivity and controllability. Unlike previous work, we capture clothing as a separate component on top of the body. With such a formulation, we use models tailored specifically to bodies and clothing, enabling applications such as virtual try-on and clothing transfer.

**Capturing both clothing and body.** Several methods model clothing as a separate layer on top of the human body. They use training data produced by physics-based simulations [Bertiche et al. 2020; Patel et al. 2020; Santesteban et al. 2019; Vidaurre et al. 2020] or require template meshes fit to 3D scans [Chen et al. 2021a; Pons-Moll et al. 2017; Tiwari et al. 2020; Xiang et al. 2021]. It is a much harder problem to recover the body and clothing from images alone, where 3D data is not provided. Jiang et al. [2020] and Zhu et al. [2020] train a multi-clothing model on 3D datasets with various clothing styles. Then during inference, a trained network produces the 3D clothing as a separate layer by recognizing and predicting the clothing style from an image. Zhu et al. [2022] fit template meshes to non-parametric 3D reconstructions. While these methods recover the clothing and body from images, they are limited in visual fidelity, as they do not capture clothing appearance. Additionally, methods with such predefined clothing style templates can not easily handle the real clothing variations, limiting their applications. In contrast, Corona et al. [2021] represent clothing layers with deep unsigned distance functions [Chibane et al. 2020], and learn the clothing style and clothing cut space with an auto-decoder. Once trained, the clothing latent code can be optimized to match image observations, but it produces over-smooth results without detailed wrinkles. Instead, SCARF models the clothing layer with a neural radiance field, and optimizes the body and clothing layer from scratch instead of the latent space of a learned model. Therefore, SCARF produces avatars with higher visual fidelity (see Section 4).

### 3 METHOD

SCARF extracts a clothed 3D avatar from a monocular video. SCARF enables us to synthesize novel views of the reconstructed avatar, and to animate the avatar with SMPL-X identity shape and pose control. The disentanglement of body and clothing further enables us to transfer clothing between subjects for virtual try-on applications.

**Key idea.** SCARF is grounded in the observation that statistical mesh models can represent human bodies well, but are ill-suited for clothing due to the large variation in clothing shape and topology (e.g., open & closed jackets, shirt, trousers, and skirts cannot be modeled with meshes of the same topology). Instead, NeRF [Mildenhall et al. 2020] offers more flexibility for modeling clothing, but is less appropriate for bodies where good models already exist. In particular, body NeRFs often lack facial details, poorly reconstruct hands, and lack fine-grained control of hand articulation and facial expression [Chen et al. 2021b; Peng et al. 2022, 2021b; Su et al. 2021]. Motivated by the strengths and weaknesses of the different representations, we use a hybrid representation that combines the strengths of body mesh models (specifically SMPL-X) with the flexibility of NeRFs; see Figure 2 for an overview.

#### 3.1 Hybrid Representation

We define the clothed body model in a canonical space, where body and clothing are represented separately.

**Body representation.** We represent the body with the expressive body model, SMPL-X [Pavlakos et al. 2019], which captures whole-body shape and pose variations, including finger articulation, and facial expressions. Given parameters for identity body shape  $\beta \in \mathbb{R}^{|\beta|}$ , pose  $\theta \in \mathbb{R}^{3n_k+3}$ , and facial expression  $\psi \in \mathbb{R}^{|\psi|}$ , SMPL-X is defined as a differentiable function  $M(\beta, \theta, \psi) \rightarrow (V, F)$  that outputs a 3D human body mesh with  $n_v$  vertices  $V \in \mathbb{R}^{n_v \times 3}$ , and  $n_f$  faces  $F \in \mathbb{R}^{n_f \times 3}$ . To increase the flexibility of the model, we add an additional set of vertex offsets  $O \in \mathbb{R}^{n_v \times 3}$  to capture localized geometric details, and define the model as

$$M(\beta, \theta, \psi, O) = \text{LBS}(T_P(\beta, \theta, \psi, O), J(\beta), \theta, \mathcal{W}), \quad (1)$$

with  $n_k$  shape dependent joints  $\mathbf{J} \in \mathbb{R}^{n_k \times 3}$ , which are a function of body shape. The linear blend skinning function LBS uses blend skinning weights  $\mathcal{W} \in \mathbb{R}^{n_k \times n_v}$ , and

$$T_P(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}, \mathbf{O}) = \mathbf{T} + \mathbf{O} + B(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}), \quad (2)$$

where  $\mathbf{T} \in \mathbb{R}^{n_v \times 3}$  is a template in rest pose, and the blend shapes

$$B(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}) = B_S(\boldsymbol{\beta}; \mathcal{S}) + B_P(\boldsymbol{\theta}; \mathcal{P}) + B_E(\boldsymbol{\psi}; \mathcal{E}). \quad (3)$$

Here,  $B_S(\boldsymbol{\beta}; \mathcal{S}) : \mathbb{R}^{|\beta|} \rightarrow \mathbb{R}^{n_v \times 3}$  are the identity blend shapes,  $B_P(\boldsymbol{\theta}; \mathcal{P}) : \mathbb{R}^{3n_k+3} \rightarrow \mathbb{R}^{n_v \times 3}$  are the pose blend shapes, and  $B_E(\boldsymbol{\psi}; \mathcal{E}) : \mathbb{R}^{|\psi|} \rightarrow \mathbb{R}^{n_v \times 3}$  are the expression blend shapes with the learned identity  $\mathcal{S}$ , pose  $\mathcal{P}$ , and expression  $\mathcal{E}$  subspaces.

Specifically, given a template vertex  $t_i$ , the vertex  $v_i$  ( $t_i$  and  $v_i$  are column vectors in homogeneous coordinates) is computed as  $v_i = M_i(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}, \mathbf{O})t_i$ , with  $M_i(\cdot) \in \mathbb{R}^{4 \times 4}$  as

$$M_i(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}, \mathbf{O}) = \left( \sum_{k=1}^{n_k} w_{k,i} G_k(\boldsymbol{\theta}, \mathbf{J}) \right) \begin{bmatrix} \mathbf{I} & \mathbf{o}_i + B_i(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}) \\ 0^T & 1 \end{bmatrix}, \quad (4)$$

where  $w_{k,i}$  is a blend weight element of  $\mathcal{W}$ ,  $G_k(\boldsymbol{\theta}, \mathbf{J}) \in \mathbb{R}^{4 \times 4}$  is the world transformation of joint  $k$ ,  $\mathbf{I} \in \mathbb{R}^{3 \times 3}$  is the identity matrix, and  $\mathbf{o}_i$  and  $B_i(\cdot)$  are the elements of the  $i$ -th vertex of  $\mathbf{O}$  and  $B(\cdot)$ , respectively. For more details regarding the SMPL-X formulation, we refer to Pavlakos et al. [2019].

To capture more geometric details, we use an upsampled version of SMPL-X with  $n_v = 38,703$  vertices and  $n_t = 77,336$  faces. We obtain this by subdividing a quad version of the model’s template, and upsampling the blend shape bases and skinning weights using barycentric coordinates obtained from the upsampled template. As the upsampling does not increase the variability of the model, we add additional learnable vertex offsets  $\mathbf{O}$  for each subject. Similar to Grassal et al. [2022], we use implicit models  $F_d : \mathbf{t} \rightarrow \mathbf{o}$  to describe the offset from every vertex  $\mathbf{t}$  of  $\mathbf{T}$ , and  $F_t : \mathbf{t} \rightarrow \mathbf{c}$  to predict the RGB color of every vertex  $\mathbf{t}$ .

**Clothing representation.** Due to the large variety of clothing in in-the-wild videos, we represent clothing using NeRF [Mildenhall et al. 2020] due to its ability to handle diverse topologies and transparent cloth materials. Following previous work (e.g., [Chen et al. 2021b; Peng et al. 2021a]), we define the NeRF model in canonical space as  $F_c : x^c \rightarrow (c, \sigma)$  to predict RGB color  $c$  and density  $\sigma$  for each query point  $x^c \in \mathbb{R}^3$ . Note that unlike previous work that models entire clothed bodies with a NeRF (e.g., [Chen et al. 2021b; Liu et al. 2021b; Peng et al. 2021a,b; Weng et al. 2022]), we only represent clothing part with a NeRF. The whole clothed body then consists of an implicit representation NeRF for clothing and an explicit surface representation for the underlying human body.

The skinning articulation of a body model like SMPL-X is not sufficient to model pose-dependent clothing deformations. Following previous work ([Liu et al. 2021b; Peng et al. 2022; Weng et al. 2022]), to model pose-dependent effects, we learn a deformation function  $F_m : \mathbb{R}^6 \rightarrow \mathbb{R}^3$  in the canonical space to model the residual non-rigid deformation. Specifically, given a body mesh  $M(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}, \mathbf{O}) \rightarrow V$  and a point  $\mathbf{x}$  and  $x^c$  in observation space and canonical space, respectively, we optimize the weights of an MLP  $F_m : (x^c, v_{\text{nn}(x)}^p) \rightarrow d^c$ , where  $\text{nn}(x)$  is the index of the nearest neighbor vertex of  $\mathbf{x}$  in  $V$ . This MLP conditions  $x^c$  on a vertex

$v^p$  from the posed mesh  $M(0, \boldsymbol{\theta}, 0, 0) \rightarrow V^p$ . Instead of  $x^c$ , the displaced point  $x^c + d^c$  in canonical space is then input to  $F_c$ .

### 3.2 Canonicalization

To model the body and clothing in canonical space, we need to transfer points in observation space to the canonical space. Following Chen et al. [2021b], we use the inverse transformation of the underlying SMPL-X model to transform from the pose  $\boldsymbol{\theta}$  in observation space to the “star-like” body pose  $\boldsymbol{\theta}^c$  (Fig. 2) in canonical space.

As the transformation between canonical space and observation space (Eq. 4) is only defined for surface vertices of the body model, Zheng et al. [2021] and Chen et al. [2021b] generalize the model transformation to the entire space. Formally, given a body mesh  $M(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}, \mathbf{O}) \rightarrow V$  and a point  $\mathbf{x}$  (in homogeneous coordinates) in observation space,  $\mathbf{x}$  is transferred to canonical space with

$$\sum_{v_i \in \mathcal{N}(\mathbf{x})} \frac{\omega_i(\mathbf{x})}{\omega(\mathbf{x})} M_i(0, \boldsymbol{\theta}^c, 0, 0) (M_i(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}, \mathbf{O}))^{-1} \mathbf{x} \rightarrow \mathbf{x}^c, \quad (5)$$

where  $\mathcal{N}(\mathbf{x})$  is the set of nearest neighbor vertices of  $\mathbf{x}$  in  $V$ . Further, the transformations are weighted with

$$\omega_i(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - v_i\|_2 \|\mathbf{w}_{\text{nn}(\mathbf{x})} - \mathbf{w}_i\|_2}{2\sigma^2}\right), \text{ and} \quad (6)$$

$$\omega(\mathbf{x}) = \sum_{v_i \in \mathcal{N}(\mathbf{x})} \omega_i(\mathbf{x}),$$

where  $\text{nn}(x)$  is the index of the nearest neighbor vertex of  $\mathbf{x}$  in  $V$ ,  $\mathbf{w}_i \in \mathbb{R}^{n_k}$  are the blend weights of  $v_i$ , and  $\sigma$  is a constant weight.

### 3.3 Mesh Integrated Volume Rendering

**Camera.** To reconstruct SMPL-X from images, we use a scaled-orthographic camera model  $\mathbf{p} = [s, \mathbf{t}^T]^T$  with isotropic scale  $s \in \mathbb{R}$  and translation  $\mathbf{t} \in \mathbb{R}^2$ .

**Mesh rendering.** Given geometry parameters  $(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi})$ , vertex offsets  $\mathbf{O}$ , colors  $F_t : t_i \rightarrow c_i$  for every vertex in the upsampled SMPL-X template, and camera information  $\mathbf{p}$ , we render the colored mesh into an image as  $\mathcal{R}_m(M(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}, \mathbf{O}), c, \mathbf{p})$ , where  $\mathcal{R}_m$  denotes the differentiable rasterizer function.

**Volume rendering.** We follow Mildenhall et al. [2020] to use volumetric rendering. Given a camera ray  $R(t) = \mathbf{o} + t\mathbf{d}$  with center  $\mathbf{o} \in \mathbb{R}^3$  and direction  $\mathbf{d} \in \mathbb{R}^3$ , the rendering interval  $t \in [t_n, t_f] \subset \mathbb{R}$  (near and far bounds) is evenly split into  $n_s$  bins. A random sample  $t_i$  ( $1 \leq i \leq n_s$ ) from every bin is taken and the colors are aggregated across the ray samples  $R(t_i) \rightarrow r_i$ . Unlike previous work, we integrate the body model,  $M(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}, \mathbf{O})$ , into the volumetric rendering. Specifically, if  $R(t)$  intersects  $M$ , we set the  $t_f$  such that  $R(t_{n_s})$  is the intersection point with  $M$ . In this case, we use the mesh color instead of the NeRF color  $c_{n_s}$  (see Fig. 3). Formally, the aggregated



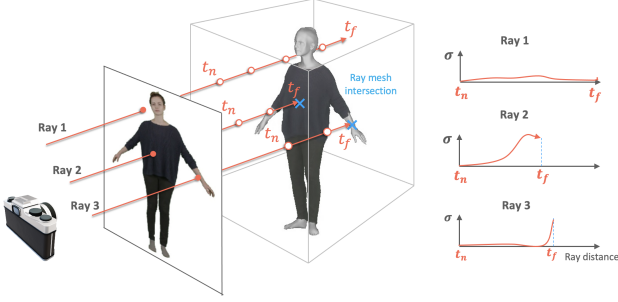


Fig. 3. Graphic illustration for mesh integrated volume rendering in Sec. 3.3.

color is

$$C(\mathbf{R}) = \sum_{i=1}^{n_s-1} \alpha_i c_i + \tau c, \text{ with } \alpha_i = \gamma_i (1 - \exp(-\sigma_i \delta_i)), \text{ where}$$

$$\gamma_i = \prod_{j=1}^{i-1} \exp(-\sigma_j \delta_j), \text{ and } \tau = 1 - \sum_{i=1}^{n_s-1} \alpha_i, \text{ and} \quad (7)$$

$$c = \begin{cases} F_t(r_{n_s}^c), & \text{if } \mathbf{R}(t) \text{ intersects } M \\ c_{n_s}, & \text{otherwise.} \end{cases}$$

Here,  $\delta_i = t_{i+1} - t_i$  is the distance between adjacent samples,  $\mathbf{R} = \{r_1, \dots, r_{n_s}\}$ ,  $F_c(r_i) \rightarrow (c_i, \sigma_i)$ , and  $r_{n_s}^c$  is the canonicalized  $r_{n_s}$ . For the scaled-orthographic camera, we use  $\mathbf{o} = [o_x, o_y, 0]$  and  $\mathbf{d} = [0, 0, 1]$  to compute the color of the pixel  $[o_x, o_y]$ . We denote the image rendered by sampling rays for all image pixels as  $\mathcal{R}_v$ .

### 3.4 Objectives

Given a sequence of  $n_f$  images,  $I_f$  ( $1 \leq f \leq n_f$ ), we optimize  $\beta$  and the weights of the MLPs  $F_d, F_c, F_t, F_m$  jointly across the entire sequence, and  $\theta_f$  and  $p_f$  per frame. The objective is

$$L = L_{\text{recon}} + L_{\text{clothing}} + L_{\text{body}}, \quad (8)$$

with reconstruction loss  $L_{\text{recon}}$ , clothing segmentation loss  $L_{\text{clothing}}$ , and body loss  $L_{\text{body}}$ . For simplicity, we omit the frame index  $f$  and the optimized parameters whenever possible. The sequence objective function is the sum over all frames.

**Reconstruction loss.** we minimize the difference between the rendered image and the input image as

$$L_{\text{recon}} = \lambda_{\text{vol}} L_{\delta}(\mathcal{R}_v - I) + \lambda_{\text{mrf}} L_{\text{mrf}}(\mathcal{R}_v - I), \quad (9)$$

where  $L_{\delta}$  is the Huber loss [Huber 1964], and  $L_{\text{mrf}}$  is an ID-MRF loss [Wang et al. 2018]. While the Huber loss focuses on the overall reconstruction, the ID-MRF loss allows us to reconstruct more details as previously shown by Feng et al. [2021b]. Solely minimizing  $L_{\text{recon}}$  results in a NeRF that models the entire clothed body including the non-clothing regions.

**Cloth segmentation loss.** Our goal is to only capture clothing with  $F_c$  instead of modeling the entire clothed body. This requires us to disentangle body and clothing. Given a clothing mask  $S_c$ , which is 1 for every clothing pixel and 0 elsewhere, we minimize the clothing segmentation loss as

$$L_{\text{clothing}} = \lambda_{\text{clothing}} \|S_v - S_c\|_{1,1}, \quad (10)$$

with the rendered NeRF mask  $S_v$ , which is obtained by sampling rays for all image pixels and computing per ray

$$S(\mathbf{R}) = \sum_{i=1}^{n_s-1} \prod_{j=1}^{i-1} \exp(-\sigma_j \delta_j) (1 - \exp(-\sigma_i \delta_i)). \quad (11)$$

Minimizing  $L_{\text{clothing}}$  ensures that the aggregated density across rays (excluding the far bound) outside of clothing is 0 and therefore nothing outside of the clothing mask is modeled by the NeRF.

**Human body loss.** To further disentangle body and clothing, we must ensure that the body model does not capture clothing variations. For this purpose, we define different losses based on four observations.

First, the body mesh should match the masked image. Given a binary mask  $S$  of the clothed body (1 for inside, 0 elsewhere), we minimize the difference between the silhouette of the rendered body  $\mathcal{R}_m^s(M, \mathbf{p})$  and the given mask as

$$L_{\text{silhouette}} = \lambda_{\text{silhouette}} L_{\delta}(\mathcal{R}_m^s(M, \mathbf{p}) - S). \quad (12)$$

Second, the body mesh should match visible body parts. Optimizing  $L_{\text{silhouette}}$  only results in meshes that also fit the clothing, which is undesired especially for loose clothing (i.e., this leads to visible artifacts when transferring clothing between subjects). Instead, given a binary mask  $S_b$  of the visible body parts (1 for body parts, 0 elsewhere), we minimize a part-based silhouette loss

$$L_{\text{bodymask}} = \lambda_{\text{bodymask}} L_{\delta}(S_b \odot \mathcal{R}_m^s(M, \mathbf{p}) - S_b), \quad (13)$$

and a part-based photometric loss

$$L_{\text{skin}} = \lambda_{\text{skin}} L_{\delta}(S_b \odot (\mathcal{R}_m(M, \mathbf{c}, \mathbf{p}) - I)), \quad (14)$$

to put special emphasis on fitting visible body parts.

Third, the body mesh should stay within clothing regions, as

$$L_{\text{inside}} = \lambda_{\text{inside}} L_{\delta}(\text{ReLU}(\mathcal{R}_m^s(M, \mathbf{p}) - S_c)). \quad (15)$$

Fourth, the skin color of occluded body vertices should be similar to non-occluded regions. For this, we assume that hands are visible for some parts of the sequence, and minimize the difference between the body colors in occluded regions and the hand color as

$$L_{\text{skininside}} = \lambda_{\text{skininside}} L_{\delta}(S_c \odot (\mathcal{R}_m(M, \mathbf{c}, \mathbf{p}) - C_{\text{hand}})), \quad (16)$$

where  $C = [c_{\text{hand}}^T, \dots, c_{\text{hand}}^T]^T \in \mathbb{R}^{n_o \times 3}$  is the tiled average color  $c_{\text{hand}}$  of the hand vertices.

**Regularization.** We regularize the reconstructed mesh surface as

$$L_{\text{reg}} = \lambda_{\text{edge}} L_{\text{edge}}(M) + \lambda_{\text{offset}} \|O\|_{2,2},$$

where  $L_{\text{edge}}$  is the relative edge loss [Hirshberg et al. 2012] between the optimized body mesh w/ and w/o applied offsets. For the offset loss, we apply different weights on the body, hands and face region. For more details see the Sup. Mat.

Overall, the body loss is

$$L_{\text{body}} = L_{\text{silhouette}} + L_{\text{bodymask}} + L_{\text{skin}} + L_{\text{skininside}} + L_{\text{inside}} + L_{\text{reg}}. \quad (17)$$

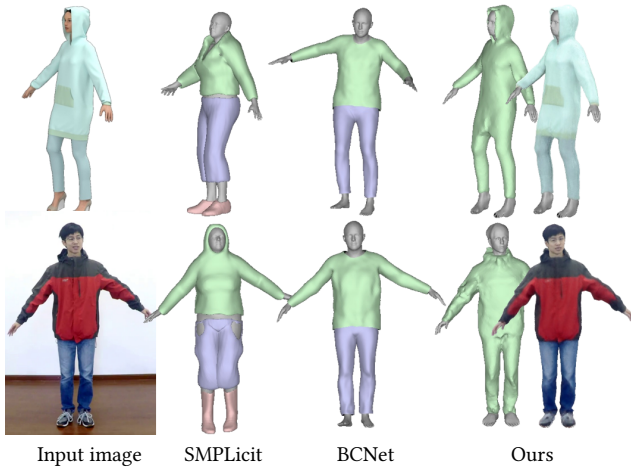


Fig. 4. Garment reconstruction comparison. SCARF reconstructs different clothing types more faithfully than SMPLicit [Corona et al. 2021] and BCNet [Jiang et al. 2020].

### 3.5 Implementation

SCARF is implemented in PyTorch, with a built-in PyTorch3D rasterizer [Ravi et al. 2020], and optimized with Adam [Kingma and Ba 2015]. For each frame, we run PIXIE [Feng et al. 2021a] to initialize  $(\beta, \theta, \psi)$ , and  $p$ . For datasets without provided silhouette masks, we compute  $S$  with [Lin et al. 2022], and [Dabhi 2022] for  $S_c$ . Following Mildenhall et al. [2020], we optimize both a coarse and a fine MLP to represent the NeRF. Our optimization pipeline has two stages. We first jointly optimize the canonical NeRF to estimate the entire clothed body (i.e., without clothing segmentation) and refine the SMPL-X pose for 100k iterations with a learning rate of  $5e-4$ . Then, we optimize the full model for another 50k iterations with learning rates of  $1e-4$  for the NeRF ( $F_c$  and  $F_m$ ) and pose refinement, and  $1e-5$  for the mesh color model ( $F_t$ ) and the offset ( $F_d$ ). For more details about the implementation, please refer to the Sup. Mat.

## 4 EXPERIMENTS

### 4.1 Datasets

We evaluate SCARF on sequences from People Snapshot [Alldieck et al. 2018b], iPER [Liu et al. 2019], SelfRecon [Jiang et al. 2022], and self-captured data. For People Snapshot, we use the provided SMPL pose as initialization instead of running PIXIE [Feng et al. 2021a]. For each subject, we use around 100-150 images for optimization. See the Sup. Mat. for more details.

### 4.2 Comparisons

Our method can capture the body and clothing from image sequences, enabling novel view synthesis. Previous works either model whole clothed body from video or reconstruct cloth geometry from a single image after training with plentiful 3D scan data. So we compare our method with others on two tasks: novel view synthesis and separate body and garment reconstruction from images.

**Body and garment reconstruction.** Similar to SCARF, SMPLicit [Corona et al. 2021] and BCNet [Jiang et al. 2020] separately model the body and clothing. Note that these methods and SCARF follow

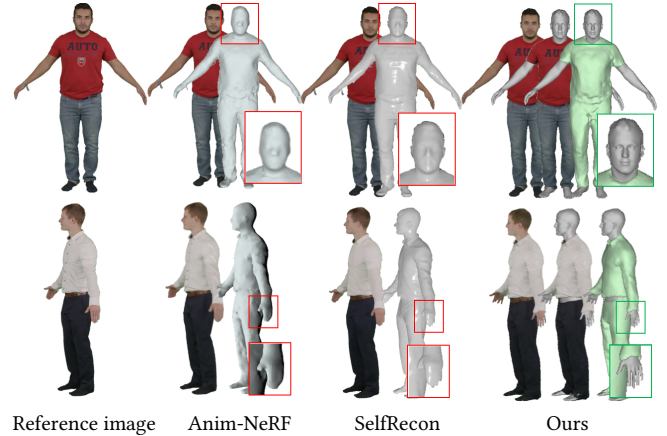


Fig. 5. Qualitative comparison with SelfRecon [Jiang et al. 2022] and Anim-NeRF [Chen et al. 2021b] for reconstruction. While all methods capture the clothing with comparable quality, our approach has much more detailed face and hands due to the disentangled representation of clothing and body.

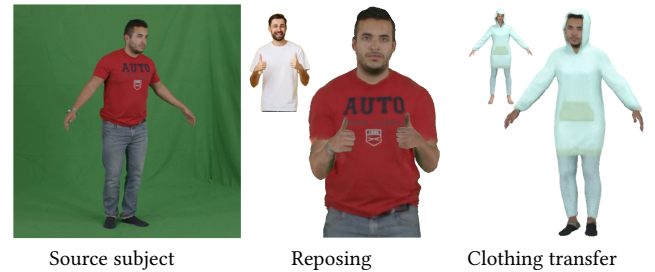


Fig. 6. Applications of SCARF. The hybrid representation enables (middle) reposing with detailed control over the body pose and (right) dressing up the source subject with target clothing. The target pose and clothing are shown in the inset images.

a different strategy. While they learn generative models from scans [Corona et al. 2021] or synthetic 3D data [Jiang et al. 2020] and then reconstruct the clothed body from a single image, SCARF extracts a clothed avatar from a video without 3D supervision. Figure 4 shows that SCARF reconstructs different clothing types more faithfully.

**Body and cloth modeling.** We quantitatively compare to NeRF [Omran et al. 2018], SMPLpix [Prokudin et al. 2021], Neural Body [Peng et al. 2021b] and Anim-NeRF [Chen et al. 2021b], following the evaluation protocol of [Chen et al. 2021b]. Table 1 shows that SCARF is more accurate than other methods under most metrics. Figure 5 provides qualitative comparisons demonstrating that SCARF better reconstructs hand and face geometry compared to SelfRecon [Jiang et al. 2022] and Anim-NeRF [Chen et al. 2021b].

### 4.3 Applications

**Animation.** Unlike previous methods that represent clothed bodies holistically, SCARF offers more fine grained control over body pose. Figure 6 shows reposing into novel poses.

**Cloth transfer.** Figures 1 and 6 and the Sup. Mat. show that our hybrid representation enables transfer of clothing between avatars.

Subject ID	PSNR $\uparrow$					SSIM $\uparrow$					LIPIS $\downarrow$				
	NeRF	SMPLpix	NB	Anim-NeRF	Ours	NeRF	SMPLpix	NB	Anim-NeRF	Ours	NeRF	SMPLpix	NB	Anim-NeRF	Ours
male-3-casual	20.64	23.74	24.94	29.37	<b>30.59</b>	.899	.923	.943	.970	<b>.977</b>	.101	.022	.033	<b>.017</b>	.024
male-4-casual	20.29	22.43	24.71	28.37	<b>28.99</b>	.880	.910	.947	.961	<b>.970</b>	.145	.031	.042	.027	<b>.025</b>
female-3-casual	17.43	22.33	23.87	28.91	<b>30.14</b>	.861	.929	.950	.974	<b>.977</b>	.170	.027	.035	<b>.022</b>	.028
female-4-casual	17.63	23.35	24.37	28.90	<b>29.96</b>	.858	.926	.945	.968	<b>.972</b>	.183	.024	.038	<b>.017</b>	.026

Table 1. Quantitative comparison of novel view synthesis on People-Snapshot [Alldieck et al. 2018b].

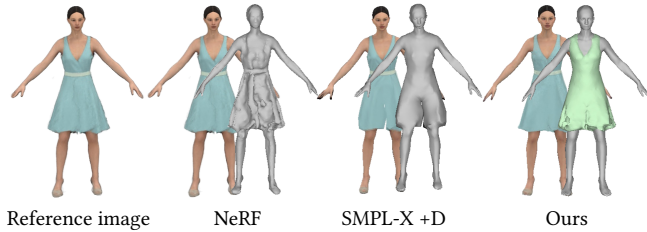


Fig. 7. Rendered images and extracted meshes from different components of SCARF. Our hybrid representation gives better estimated face, hand, and clothing geometry than vanilla NeRF or a mesh-based representation.

#### 4.4 Ablation Experiments

We run different ablation experiments to show the impact of different components of our hybrid representation (below), and to show the impact of the pose optimization in Sup. Mat.

**Effect of representations.** SCARF consists of a NeRF to represent clothing, and a mesh with vertex displacements. Figure 7 compares NeRF to holistically represent body and clothing (i.e., SCARF w/o body-clothing segmentation) and mesh-only based representation (i.e., SCARF w/o NeRF). Our hybrid representation is better able to estimate the face, hands, and complex clothing. Note that, unlike our hybrid representation, none of the existing body NeRF methods is able to transfer clothing between avatars.

## 5 DISCUSSION AND LIMITATIONS

**Segmentation.** SCARF requires body and cloth segmentation for training. Segmentation errors of the clothed body and background negatively impact the visual quality of the extracted avatar, and erroneous clothing segmentation results in poor separation of body and clothing. Enforcing temporal consistency by exploiting optical flow could improve the segmentation quality.

**Shoes & hair.** Modeling hair, shoes, or other accessories with NeRF would improve the visual quality of SCARF. We will explore alternative shoe and hair segmentation methods (e.g., [Yang et al. 2020]) to extend SCARF.

**Geometric quality.** The strength of NeRF is its visual quality and the ability to synthesize realistic images, even when the geometry is not perfect. In contrast, recent SDF-based methods have demonstrated good geometric reconstruction (e.g., [Jiang et al. 2022]). It may be possible to leverage their results to better represent the underlying clothed shape or to regularize NeRF.

**Novel poses.** While SCARF generalizes to unseen poses, extreme poses result in artifacts (see Sup. Mat.). Possible solutions include regularizing NeRF during optimization or learning a generative model from many training examples of different people and poses.

**Pose initialization.** SCARF refines the body pose during optimization. However, it may fail if the initial pose is far from the right pose. Handling difficult poses where PIXIE [Feng et al. 2021a] fails requires a more robust 3D body pose estimator.

**Dynamics.** SCARF handles non-rigid cloth deformation with the pose-conditioned deformation model. While the global pose accounts for some deformation, the modeling of clothing dynamics as a function of body movement is the subject of future work.

**Lighting.** As with other NeRF methods, we do not factor lighting and material properties. This results in baked-in shading and the averaging of specular reflections across frames. Factoring lighting from shape and material is a key next step to improve realism.

**Facial expressions.** SCARF uses the facial expressions estimated by PIXIE [Feng et al. 2021a] which is unable to capture the full spectrum of emotions (cf. [Danecek et al. 2022]). Also, we have not fully exploited neural radiance fields to capture complex changes in facial appearance, e.g. due to the mouth opening. We believe this is a promising future direction.

## 6 CONCLUSION

SCARF automatically extracts an animatable clothed 3D human avatar from a monocular video. Our key novelty is a hybrid representation that combines a mesh-based body model with a neural radiance field to separately model the body and clothing. This factored representation enables SCARF to transfer clothing between avatars, animate the body pose of the avatars including finger articulation, alter their body shape (see Sup. Mat.) and facial expression, and visualize them from unseen viewing directions. This property makes SCARF well suited to VR and virtual try-on applications. Finally, SCARF outperforms existing avatar extraction methods from videos in terms of visual quality and generality.

## ACKNOWLEDGMENTS

We thank Sergey Prokudin, Weiyang Liu, Yuliang Xiu, Songyou Peng, Qianli Ma for fruitful discussions, and Peter Kulits, Zhen Liu, Yandong Wen, Hongwei Yi, Xu Chen, Soubhik Sanyal, Omri Ben-Dov, Shashank Tripathi for proofreading. We also thank Betty Mohler, Sarah Danes, Natalia Marciniak, Tsvetelina Alexiadis, Claudia Gallatz, and Andres Camilo Mendoza Patino for their supports with data. This work was partially supported by the Max Planck ETH Center for Learning Systems.

**Disclosure.** MJB has received research gift funds from Adobe, Intel, Nvidia, Meta/Facebook, and Amazon. MJB has financial interests in Amazon, Datagen Technologies, and Meshcapade GmbH. While TB is part-time employee of Amazon, this research was performed solely at, and funded solely by, MPI.

## REFERENCES

- Thiemo Alldieck, Marcus Magnor, Bharat Lal Bhatnagar, Christian Theobalt, and Gerard Pons-Moll. 2019a. Learning to reconstruct people in clothing from a single RGB camera. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1175–1186.
- Thiemo Alldieck, Marcus Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. 2018a. Detailed Human Avatars from Monocular Video. In *International Conference on 3D Vision (3DV)*.
- Thiemo Alldieck, Marcus Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. 2018b. Video based reconstruction of 3d people models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8387–8397.
- Thiemo Alldieck, Gerard Pons-Moll, Christian Theobalt, and Marcus Magnor. 2019b. Tex2shape: Detailed full human body geometry from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2293–2303.
- Thiemo Alldieck, Hongyi Xu, and Cristian Sminchisescu. 2021. imGHUM: Implicit Generative Models of 3D Human Shape and Articulated Pose. In *International Conference on Computer Vision (ICCV)*. IEEE, 5441–5450.
- Dragomir Anguelov, Praveen Srinivasan, Daphne Koller, Sebastian Thrun, Jim Rodgers, and James Davis. 2005. SCAPE: Shape completion and animation of people. In *ACM SIGGRAPH 2005 Papers*. 408–416.
- Hugo Bertiche, Meysam Madadi, and Sergio Escalera. 2020. CLOTH3D: clothed 3d humans. In *European Conference on Computer Vision (ECCV)*. Springer, 344–359.
- Jianchuan Chen, Ying Zhang, Di Kang, Xuefei Zhe, Linchao Bao, Xu Jia, and Huchuan Lu. 2021b. Animatable Neural Radiance Fields from Monocular RGB Videos. arXiv:2106.13629 [cs.CV]
- Xin Chen, Anqi Pang, Wei Yang, Peihao Wang, Lan Xu, and Jingyi Yu. 2021a. TightCap: 3D Human Shape Capture with Clothing Tightness Field. *Transactions on Graphics (TOG)* 41, 1 (2021), 1–17.
- Julian Chibane, Aymen Mir, and Gerard Pons-Moll. 2020. Neural Unsigned Distance Fields for Implicit Function Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Vasileios Choutas, Georgios Pavlakos, Timo Bolkart, Dimitrios Tzionas, and Michael J. Black. 2020. Monocular Expressive Body Regression through Body-Driven Attention. In *European Conference on Computer Vision (ECCV)*. 20–40.
- Enric Corona, Albert Pumarola, Guillem Alenyà, Gerard Pons-Moll, and Francesco Moreno-Noguer. 2021. SMPlicit: Topology-aware generative model for clothed people. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 11875–11885.
- Levin Dabhi. 2022. Clothes Segmentation using U2NET. <https://github.com/levindabhi/cloth-segmentation>
- Radek Danecek, Michael J. Black, and Timo Bolkart. 2022. EMOCA: Emotion Driven Monocular Face Capture and Animation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Yao Feng, Vasileios Choutas, Timo Bolkart, Dimitrios Tzionas, and Michael Black. 2021a. Collaborative Regression of Expressive Bodies using Moderation. In *International Conference on 3D Vision (3DV)*. 792–804.
- Yao Feng, Haiwen Feng, Michael J. Black, and Timo Bolkart. 2021b. Learning an Animatable Detailed 3D Face Model from In-the-Wild Images. *Transactions on Graphics, (Proc. SIGGRAPH)* 40, 4 (2021), 88:1–88:13.
- Philip-William Grassal, Malte Prinzler, Titus Leistner, Carsten Rother, Matthias Nießner, and Justus Thies. 2022. Neural Head Avatars from Monocular RGB Videos. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Tong He, Yuanlu Xu, Shunsuke Saito, Stefano Soatto, and Tony Tung. 2021. ARCH++: Animation-ready clothed human reconstruction revisited. In *International Conference on Computer Vision (ICCV)*. 11046–11056.
- David A. Hirshberg, Matthew Loper, Eric Rachtlin, and Michael J. Black. 2012. Coregistration: Simultaneous Alignment and Modeling of Articulated 3D Shape. In *European Conference on Computer Vision (ECCV) (Lecture Notes in Computer Science, Vol. 7577)*. Springer, 242–255.
- Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. 2022. HeadNeRF: A real-time nerf-based parametric head model. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 20374–20384.
- Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. 2020. ARCH: Animatable reconstruction of clothed humans. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 3093–3102.
- Peter J. Huber. 1964. Robust Estimation of a Location Parameter. *The Annals of Mathematical Statistics* 35, 1 (1964), 73–101.
- Boyi Jiang, Yang Hong, Hujun Bao, and Juyong Zhang. 2022. SelfRecon: Self Reconstruction Your Digital Avatar from Monocular Video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5605–5615.
- Boyi Jiang, Juyong Zhang, Yang Hong, Jinhao Luo, Ligang Liu, and Hujun Bao. 2020. BCNet: Learning body and cloth shape from a single image. In *European Conference on Computer Vision (ECCV)*. Springer, 18–35.
- Ning Jin, Yilin Zhu, Zhenglin Geng, and Ronald Fedkiw. 2020. A Pixel-Based Framework for Data-Driven Clothing. In *Computer Graphics Forum*, Vol. 39. Wiley Online Library, 135–144.
- Hanbyul Joo, Tomas Simon, and Yaser Sheikh. 2018. Total Capture: A 3D Deformation Model for Tracking Faces, Hands, and Bodies. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 8320–8329.
- Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. 2018. End-to-end Recovery of Human Shape and Pose. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 7122–7131.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *International Conference on Learning Representations (ICLR)*.
- Nikos Kolotouros, Georgios Pavlakos, Michael J. Black, and Kostas Daniilidis. 2019. Learning to Reconstruct 3D Human Pose and Shape via Model-Fitting in the Loop. In *International Conference on Computer Vision (ICCV)*. IEEE, 2252–2261.
- Verica Lazova, Eldar Insafutdinov, and Gerard Pons-Moll. 2019. 360-Degree Textures of People in Clothing from a Single Image. In *International Conference on 3D Vision (3DV)*.
- Shanchuan Lin, Linjie Yang, Imran Saleemi, and Soumyadip Sengupta. 2022. Robust Video Matting (RVM). <https://github.com/PeterLin/RobustVideoMatting>
- Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. 2021b. Neural actor: Neural free-view synthesis of human actors with pose control. *Transactions on Graphics (TOG)* 40, 6 (2021), 1–16.
- Shaohui Liu, Yinda Zhang, Songyou Peng, Boxin Shi, Marc Pollefeys, and Zhaopeng Cui. 2020. Dist: Rendering deep implicit signed distance function with differentiable sphere tracing. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019–2028.
- Wu Liu, Qian Bao, Yu Sun, and Tao Mei. 2021a. Recent Advances in Monocular 2D and 3D Human Pose Estimation: A Deep Learning Perspective. *CoRR* abs/2104.11536 (2021).
- Wen Liu, Zhixin Piao, Min Jie, Wenhan Luo, Lin Ma, and Shenghua Gao. 2019. Liquid Warping GAN: A Unified Framework for Human Motion Imitation, Appearance Transfer and Novel View Synthesis. In *International Conference on Computer Vision (ICCV)*. 5903–5912.
- Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. 2015. SMPL: A Skinned Multi-Person Linear Model. *Transactions on Graphics, (Proc. SIGGRAPH Asia)* 34, 6 (2015), 248:1–248:16.
- Qianli Ma, Jinlong Yang, Anurag Ranjan, Sergi Pujades, Gerard Pons-Moll, Siyu Tang, and Michael J. Black. 2020a. Learning to dress 3D people in generative clothing. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 6469–6478.
- Qianli Ma, Jinlong Yang, Anurag Ranjan, Sergi Pujades, Gerard Pons-Moll, Siyu Tang, and Michael J. Black. 2020b. Learning to Dress 3D People in Generative Clothing. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 6468–6477.
- Marko Mihajlovic, Yan Zhang, Michael J. Black, and Siyu Tang. 2021. LEAP: Learning Articulated Occupancy of People. In *Proceedings IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*. 10461–1047.
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision (ECCV)*. Springer, 405–421.
- Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. 2020. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 3504–3515.
- Mohamed Omran, Christoph Lassner, Gerard Pons-Moll, Peter Gehler, and Bernt Schiele. 2018. Neural body fitting: Unifying deep learning and model based human pose and shape estimation. In *International Conference on Computer Vision (ICCV)*. IEEE, 484–494.
- Ahmed A. A. Osman, Timo Bolkart, and Michael J. Black. 2020. STAR: Sparse trained articulated human body regressor. In *European Conference on Computer Vision (ECCV)*. 598–613.
- Chaitanya Patel, Zhouyingcheng Liao, and Gerard Pons-Moll. 2020. TailorNet: Predicting clothing in 3d as a function of human pose, shape and garment style. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 7365–7375.
- Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. 2019. Expressive Body Capture: 3D Hands, Face, and Body From a Single Image. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 10975–10985.
- Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. 2021a. Animatable Neural Radiance Fields for Modeling Dynamic Human Bodies. In *ICCV*.
- Sida Peng, Shangzhan Zhang, Zhen Xu, Chen Geng, Boyi Jiang, Hujun Bao, and Xiaowei Zhou. 2022. Animatable Neural Implicit Surfaces for Creating Avatars from Videos. *arXiv preprint arXiv:2203.08133* (2022).
- Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. 2021b. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 9054–9063.
- Gerard Pons-Moll, Sergi Pujades, Sonny Hu, and Michael J. Black. 2017. ClothCap: Seamless 4D clothing capture and retargeting. *Transactions on Graphics (TOG)* 36, 4 (2017), 1–15.

- Sergey Prokudin, Michael J. Black, and Javier Romero. 2021. SMPLpix: Neural Avatars from 3D Human Models. In *Winter Conference on Applications of Computer Vision (WACV)*. 1810–1819.
- Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. 2020. Accelerating 3D Deep Learning with Py-Torch3D. *arXiv:2007.08501* (2020).
- Yu Rong, Takaaki Shiratori, and Hanbyul Joo. 2021. FrankMocap: A Monocular 3D Whole-Body Pose Estimation System via Regression and Integration. In *International Conference on Computer Vision Workshops (ICCV-W)*.
- Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. 2019. PIFu: Pixel-Aligned Implicit Function for High-Resolution Clothed Human Digitization. In *International Conference on Computer Vision (ICCV)*.
- Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. 2020. PIFuHD: Multi-Level Pixel-Aligned Implicit Function for High-Resolution 3D Human Digitization. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Igor Santesteban, Miguel A Otaduy, and Dan Casas. 2019. Learning-based animation of clothing for virtual try-on. In *Computer Graphics Forum*, Vol. 38. Wiley Online Library, 355–366.
- Shih-Yang Su, Frank Yu, Michael Zollhöfer, and Helge Rhodin. 2021. A-NeRF: Articulated neural radiance fields for learning human shape, appearance, and pose. *Advances in Neural Information Processing Systems (NeurIPS)* 34 (2021).
- Yating Tian, Hongwen Zhang, Yebin Liu, and Limin Wang. 2022. Recovering 3D Human Mesh from Monocular Images: A Survey. *arXiv preprint arXiv:2203.01923* (2022).
- Garvita Tiwari, Bharat Lal Bhatnagar, Tony Tung, and Gerard Pons-Moll. 2020. SIZER: A dataset and model for parsing 3d clothing and learning size sensitive 3d clothing. In *European Conference on Computer Vision (ECCV)*. Springer, 1–18.
- Raquel Vidas, Igor Santesteban, Elena Garces, and Dan Casas. 2020. Fully Convolutional Graph Neural Networks for Parametric Virtual Try-On. In *Computer Graphics Forum*, Vol. 39. Wiley Online Library, 145–156.
- Yi Wang, Xin Tao, Xiaojuan Qi, Xiaoyong Shen, and Jiaya Jia. 2018. Image inpainting via generative multi-column convolutional neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*. 331–340.
- Chung-Yi Weng, Brian Curless, Pratul P. Srinivasan, Jonathan T. Barron, and Ira Kemelmacher-Shlizerman. 2022. HumanNeRF: Free-Viewpoint Rendering of Moving People From Monocular Video. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 16210–16220.
- Donglai Xiang, Hanbyul Joo, and Yaser Sheikh. 2019. Monocular Total Capture: Posing Face, Body, and Hands in the Wild. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 10965–10974.
- Donglai Xiang, Fabian Prada, Timur Bagautdinov, Weipeng Xu, Yuan Dong, He Wen, Jessica Hodgins, and Chenglei Wu. 2021. Modeling clothing as a separate layer for an animatable human avatar. *Transactions on Graphics (TOG)* 40, 6 (2021), 1–15.
- Yuliang Xiu, Jinlong Yang, Dimitrios Tzionas, and Michael J. Black. 2022. ICON: Implicit Clothed humans Obtained from Normals. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Hongyi Xu, Thiemo Alldieck, and Cristian Sminchisescu. 2021. H-NeRF: Neural radiance fields for rendering and temporal reconstruction of humans in motion. *Advances in Neural Information Processing Systems (NeurIPS)* 34 (2021).
- Hongyi Xu, Eduard Gabriel Bazavan, Andrei Zanfir, William T Freeman, Rahul Sukthankar, and Cristian Sminchisescu. 2020. GHUM & GHUML: Generative 3d human shape and articulated pose models. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 6184–6193.
- Lu Yang, Qing Song, Zhihui Wang, Mengjie Hu, Chun Liu, Xueshi Xin, Wenhe Jia, and Songcen Xu. 2020. Renovating Parsing R-CNN for Accurate Multiple Human Parsing. In *European Conference on Computer Vision (ECCV) (Lecture Notes in Computer Science, Vol. 12357)*. Springer, 421–437.
- Ze Yang, Shenlong Wang, Sivabalan Manivasagam, Zeng Huang, Wei-Chiu Ma, Xinchun Yan, Ersin Yumer, and Raquel Urtasun. 2021. S3: Neural shape, skeleton, and skinning fields for 3D human modeling. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 13284–13293.
- Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. 2021. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems (NeurIPS)* 34 (2021).
- Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. 2020. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems (NeurIPS)* 33 (2020), 2492–2502.
- Andrei Zanfir, Eduard Gabriel Bazavan, Mihai Zanfir, William T. Freeman, Rahul Sukthankar, and Cristian Sminchisescu. 2021. Neural Descent for Visual 3D Human Pose and Shape. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. Computer Vision Foundation / IEEE, 14484–14493.
- Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai. 2021. PaMIR: Parametric model-conditioned implicit representation for image-based human reconstruction. *Transactions on Pattern Analysis and Machine Intelligence (PAMI)* (2021).
- Yuxiao Zhou, Marc Habermann, Ikhsanul Habibie, Ayush Tewari, Christian Theobalt, and Feng Xu. 2021. Monocular Real-Time Full Body Capture With Inter-Part Correlations. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 4811–4822.
- Heming Zhu, Yu Cao, Hang Jin, Weikai Chen, Dong Du, Zhangye Wang, Shuguang Cui, and Xiaoguang Han. 2020. Deep Fashion3D: A dataset and benchmark for 3D garment reconstruction from single images. In *European Conference on Computer Vision (ECCV)*. Springer, 512–530.
- Heming Zhu, Lingteng Qiu, Yuda Qiu, and Xiaoguang Han. 2022. Registering Explicit to Implicit: Towards High-Fidelity Garment mesh Reconstruction from Single Images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3845–3854.



## A APPENDIX

The supplementary material includes this document and an additional video. Here, we provide more details about the datasets and the implementation, and present further results.

### A.1 Implementation Details

We choose  $\sigma = 0.1$ ,  $|\mathcal{N}(x)| = 6$ ,  $t_n = -0.6$ , and  $t_f = 0.6$  and weight the individual losses with  $\lambda_{\text{vol}} = 1.0$ ,  $\lambda_{\text{mrf}} = 0.0005$ ,  $\lambda_{\text{clothing}} = 0.5$ ,  $\lambda_{\text{silhouette}} = 0.001$ ,  $\lambda_{\text{bodymask}} = 30$ ,  $\lambda_{\text{skin}} = 1.0$ ,  $\lambda_{\text{inside}} = 40$ ,  $\lambda_{\text{skininside}} = 0.01$ ,  $\lambda_{\text{jap}} = 500$ ,  $\lambda_{\text{offset}} = 400$ . For  $\lambda_{\text{offset}}$ , the weight ratio of body, face and hands region is 2 : 3 : 12. Note that it is important to perform the first stage NeRF training without optimizing the non-rigid deformation model. In this stage, we also set  $\lambda_{\text{mrf}} = 0$ . In the second stage, the non-rigid deformation model then explains clothing deformations that cannot be explained by the body transformation. And  $L_{\text{mrf}}$  helps capture more details that can not be modelled by the non-rigid deformation. The overall optimization time is around 40 hours with NVIDIA V100.

### A.2 Datasets

We use 4 subjects ('male-3-casual', 'female-3-casual', 'male-4-casual', 'female-4-casual') from People Snapshot [Alldieck et al. 2018b] for qualitative and quantitative evaluation. We follow the settings of Anim-NeRF [Chen et al. 2021b], namely

- 'male-3-casual': frames 1-456 with step size of 4 for training, and frames 456-676 with step size 4 for test.
- 'male-4-casual': frame 1-660 with step size 6 for training, and frames 661-873 with step size 4 for test.
- 'female-3-casual': frame 1-446 with step size 4 for training, and frames 447-648 with step size 4 for test.
- 'female-4-casual': frame 1-336 with step size 4 for training, and frames 336-524 with step size 4 for test.

We further use 4 subjects ('subject003', 'subject016', 'subject022', 'subject023') with outfit 1 and motion 1 from iPER [Liu et al. 2019] for qualitative evaluation. For all subjects, we use frames 1-490 with step 4 for optimization. We use 4 synthetic video data ('female outfit1', 'female outfit2', 'female outfit3', 'male outfit1') and 1 self-captured video ('CHH female') from SelfRecon [Jiang et al. 2022]. For each subject, we use 100 frames for optimization. For self-captured data, we record videos of each subject wearing different clothing types. The subject wears different clothes and performs an A-pose video and a video with random actions. In our experiments, we use A-pose videos of subject 'Yao' with six types of clothing for qualitative evaluation, those videos include loose dressing and short skirts. For each video, we use frames 0-400 with step 2 for optimization.

### A.3 Ablation Experiments

**Effect of pose refinement.** Since the pose estimation for each frame is not accurate, the pose refinement is important to gain details. We try learning our method without pose refinement. Fig. 8 shows that pose refinement improves the image quality a lot.

### A.4 More Qualitative results

We show additional comparisons on Garment reconstruction with SMPLicit [Corona et al. 2021] and BCNet [Jiang et al. 2020] in Fig 13.

SCARF gives better visual quality than SMPLicit and BCNet. Note that the training/optimization settings are different, they reconstruct the body and garment from a single image, while our results are learned from video. However, they require a large set of 3D scans and manually designed cloth template for training, while we do not need any 3D supervision, and capture the garment appearance as well.

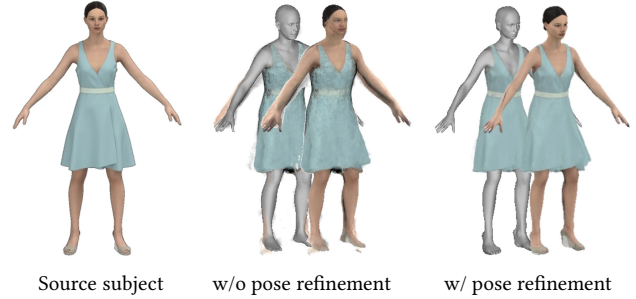


Fig. 8. Novel view synthesis w/o and w/ pose refinement. The pose refinement improves the visual quality of the reconstruction, as more texture details are reconstructed.

In Figure 9, we also show that SCARF can alter body shape and the clothing will adapt to the shape accordingly.

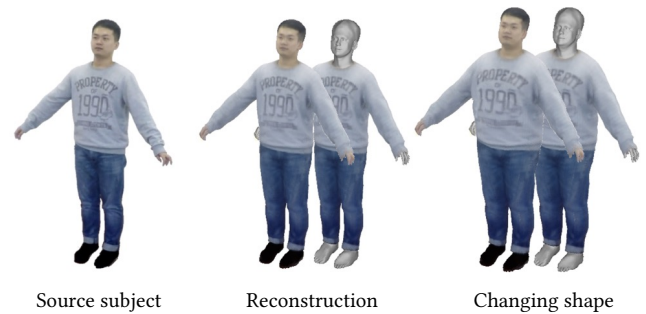


Fig. 9. SCARF can change underlying body shapes by altering SMPL-X shape parameters, the NeRF clothing will adapt to the body accordingly.

### A.5 Limitations

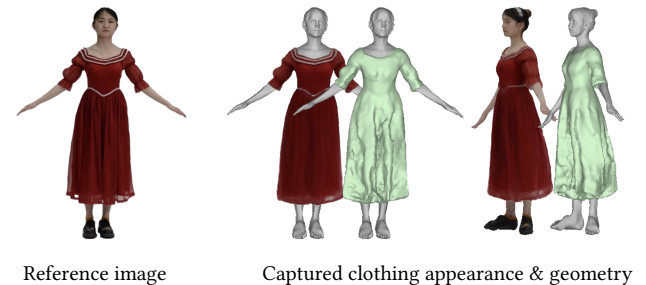


Fig. 10. While SCARF gives good visual quality for clothing renderings, the underlying geometry of the NeRF clothing is sometimes noisy.

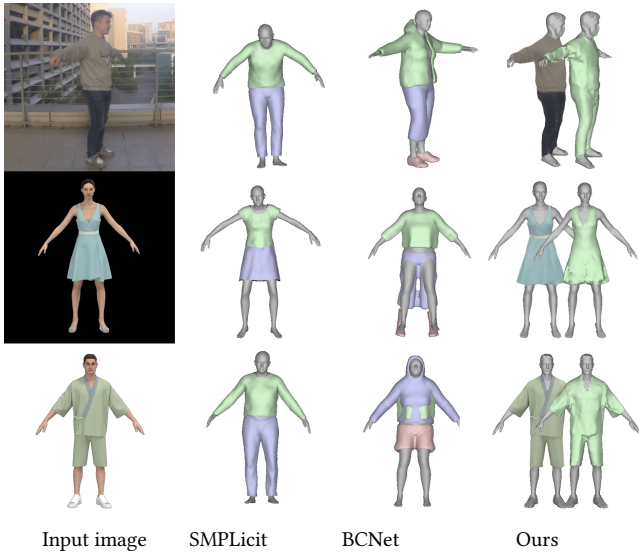


Fig. 13. Additional examples for qualitative comparison of garment reconstruction. SCARF reconstructs different clothing types more faithfully than SMPLicit [Corona et al. 2021] and BCNet [Jiang et al. 2020].

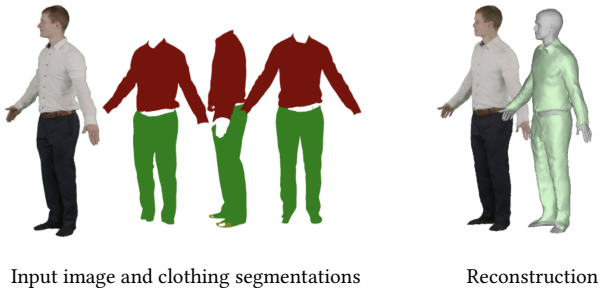


Fig. 11. The wrong clothing segmentation results in a visible gap within the reconstructed clothing.

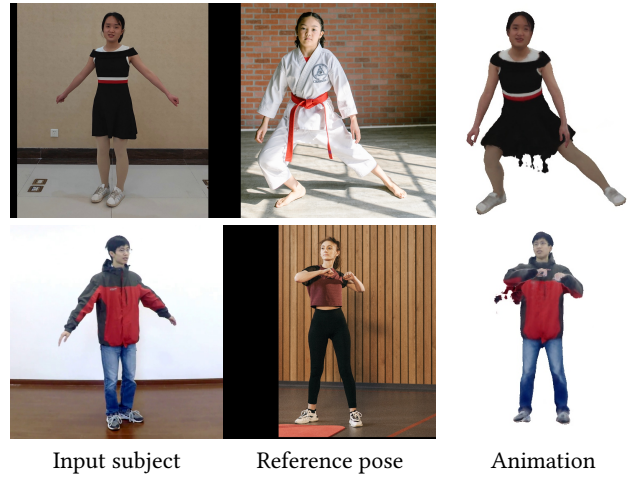


Fig. 12. Reposing can result in visual artifacts for unseen poses.

The main paper discusses limitations regarding the image segmentation, the quality of the reconstructed geometry, and the generalization to unseen poses. Figure 11 shows the wrong reconstruction due to consistent clothing segmentation errors, e.g. the belt is not recognized as part of clothing in segmentation, this results in wrong disentanglement between human body and clothing. Fig. 10 shows an example of noisy geometry despite good visual quality, and Figure 12 shows some reposing artifacts for unseen poses.