

# imGHUM: Implicit Generative Models of 3D Human Shape and Articulated Pose

Thiemo Alldieck\*

Hongyi Xu\*

Cristian Sminchisescu

Google Research

{alldieck, hongyixu, sminchisescu}@google.com

## Abstract

We present imGHUM, the first holistic generative model of 3D human shape and articulated pose, represented as a signed distance function. In contrast to prior work, we model the full human body implicitly as a function zero-level-set and without the use of an explicit template mesh. We propose a novel network architecture and a learning paradigm, which make it possible to learn a detailed implicit generative model of human pose, shape, and semantics, on par with state-of-the-art mesh-based models. Our model features desired detail for human models, such as articulated pose including hand motion and facial expressions, a broad spectrum of shape variations, and can be queried at arbitrary resolutions and spatial locations. Additionally, our model has attached spatial semantics making it straightforward to establish correspondences between different shape instances, thus enabling applications that are difficult to tackle using classical implicit representations. In extensive experiments, we demonstrate the model accuracy and its applicability to current research problems.

## 1. Introduction

Mathematical models of the human body have been proven effective in a broad variety of tasks. In the last decades models of varying degrees of realism have been successfully deployed e.g. for 3D human motion analysis [48], 3D human pose and shape reconstruction [25, 54], personal avatar creation [3, 56], medical diagnosis and treatment [17], or image synthesis and video editing [55, 22]. Modern statistical body models are typically learnt from large collections of 3D scans of real people, which are used to capture the body shape variations among the human population. Dynamic scans, when available, can be used to further model how different poses affect the deformation of the muscles and the soft-tissue of the human body.

The recently released GHUM model [51] follows this methodology by describing the human body, its shape variation, articulated pose including fingers, and facial expressions as a moderate resolution mesh based on a low-dimensional, partly interpretable parameterization. In the



Figure 1. imGHUM is the first parametric full human body model represented as an implicit signed distance function. imGHUM successfully models broad variations in pose, shape, and facial expressions. The level sets of imGHUM are shown in blue-scale.

deep learning literature GHUM and similar models [29, 24] are typically used as fixed function layers. This means that the model is parameterized with the output of a neural network or some other non-linear function, and the resulting mesh is used to compute the final function value. While this approach works well for several tasks, including, more recently, 3D reconstruction, the question of how to best represent complex 3D deformable and articulated structures is open. Recent work dealing with the 3D visual reconstruction of general objects aimed to represent the output not as meshes but as implicit functions [30, 34, 7, 31]. Such approaches thus describe surfaces by the zero-level-set (decision boundary) of a function over points in 3D-space. This has clear benefits as the output is neither constrained by a template mesh topology, nor is it discretized and thus of fixed spatial resolution.

In this work, we investigate the possibility to learn a data-driven statistical body model as an implicit function. Given the maturity of state of the art explicit human models, it is crucial that an equivalent implicit representation maintains their key, attractive properties – representing comparable variation in shape and pose and similar level of detail. This is challenging since recently-proposed implicit function networks tend to produce overly smooth shapes and fail for articulated humans [8]. We propose a novel network architecture and a learning paradigm that enable, for the first time, constructing detailed generative models of human pose, shape, and semantics, represented as Signed Distance Functions (SDFs) (see fig. 1). Our multi-part architecture focuses on difficult to model body components like

\* The first two authors contributed equally.

generative pose	generative shape	gen. hands	gen. expression	interpolation	signed distances	semantics	continuous rep.	
✓	✓	✓	✓	✓	✗	✓	✗	GHUM [51]
✗	✗	✗	✗	✗	✗	✗	✓	IF-Net [8]
✗	✗	✗	✗	✓	✓	✗	✓	IGR [14]
✓	✗	✗	✗	✓	✗	✗	✓	NASA [11]
✓	✓	✓	✓	✓	✓	✓	✓	imGHUM

Table 1. Comparison of different approaches to model human bodies. GHUM is meshed-based and thus discretized. IGR only allows for shape interpolation. NASA lacks generative capabilities for shape, hands, and facial expressions and only returns occupancy values. Only imGHUM combines all favorable properties.

hands and faces. Moreover, imGHUM models its neighborhood through distance values, enabling e.g. collision tests. Our model is not bound to a specific resolution and thus can be easily queried at arbitrary locations. Being template-free further paves the way to our ultimate goal to fairly represent diversity of mankind, including disabilities which may not be always well covered by a generic template of standard topology. Finally, in contrast to recent implicit function networks, our model additionally carries on the explicit semantics of mesh-based models. Specifically, our implicit function also returns correspondences to a canonical representation near and on its zero-level-set, enabling e.g. texturing or body part labeling. This holistic approach is novel and significantly more difficult to produce, as can be noted in prior work which could only demonstrate individual properties, *c.f.* tab. 1. Our contribution – and the key to success – stems from the novel combination of adequate, generative latent representations, network architectures with fine grained encoding, implicit losses with attached semantics, and the consistent aggregation of multi-part components. Besides extensive evaluation of 3D deformable and articulated modeling capabilities, we also demonstrate surface completion using imGHUM and give an outlook to modeling varying topologies. Our models are available for research [1].

## 1.1. Related Work

We review developments in 3D human body modeling, variants of implicit function networks, and applications of implicit function networks for 3D human reconstruction.

**Human Body Models.** Parametric human body models based on geometric primitives have been proposed early on [50] and successfully applied e.g. for human reconstruction from video data [38, 48, 47]. SCAPE [37] was one of the first realistic large scale data-driven human body models. Later variants inspired by blend skinning [18] modeled correlations between body shape and pose [16], as well as soft-tissue dynamics [39]. SMPL variants [29, 24, 35, 33] are also popular parametric body models, with linear shape spaces, compatible with standard graphics pipelines and offering good full-body representation functionality.

GHUM is a recent parametric model [51] that represents the full body model using deep non-linear models – VAEs for shape and normalizing flows for pose, respectively – with various trainable parameters, learned end-to-end. In this work, we rely on GHUM to build our novel implicit model. Specifically, besides the static and dynamic 3D human scans in our dataset, we also rely on GHUM (1) to represent the latent pose and shape state of our implicit model, (2) to generate supervised training data in the form of latent pose and shape codes with associated 3D point clouds, sampled from the underlying, posed, GHUM mesh.

**Implicit Function Networks (IFNs)** have been proposed recently [30, 34, 7, 31]. Instead of representing shapes as meshes, voxels, or point clouds, IFNs learn a shape space as a function of a low-dimensional global shape code and a 3D point. The function either classifies the point as inside/outside [30, 7] (occupancy networks), or returns its distance to the closest surface [34] (distance functions). The global shape is then defined by the decision boundary or the zero-level-set of this function.

Despite advantages over mesh- and voxel-based representations in tasks like e.g. 3D shape reconstruction from partial views or given incomplete data, initial work has limitations. First, while the models can reliably encode rigid axis-aligned shape prototypes, they often fail for more complex shapes. Second, the reconstructions are often overly smooth, hence they lack detail. Different approaches have been presented to address these. Part-based models [13, 23, 12] assemble a global shape from smaller local models. Some methods do not rely on a global shape code but on features computed from convolving with an input observation [8, 10, 36, 9]. Others address such limitations by changing the learning methodology: tailored network initialization [4] and point sampling strategies [52], or second-order losses [14, 46] have been proposed towards this end. We found the latter to be extremely useful and rely on similar losses in this work.

**IFNs for Human Reconstruction.** Recently implicit functions have been explored to reconstruct humans. Huang et al. [19] learn an occupancy network that conditions on image features in a multi-view camera setup. Saito et al. [43] use features from a single image and an estimated normal image [44] together with depth values along camera rays as conditioning variables. ARCH [20] combines implicit function reconstruction and explicit mesh-based human models to represent dressed people. Karunratanakul et al. [26] propose to use SDFs to learn human grasps and augment their SDFs output with sparse regional labels. Similarly to us, Deng et al. [11] represent a pose-able human subject as a number of binary occupancy functions modeled in a kinematic structure. In contrast to our work, this framework is restricted to a single person and the body is only coarsely approximated, lacking facial features and hand detail. Also

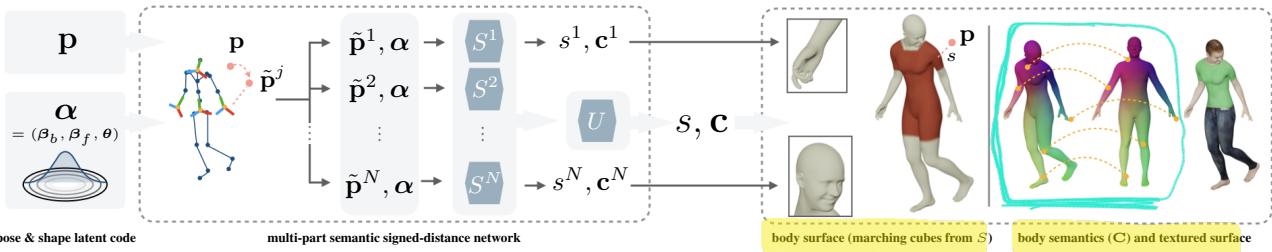


Figure 2. Overview of imGHUM. We compute the signed distance  $s = S(\mathbf{p}, \alpha)$  and the semantics  $\mathbf{c} = \mathbf{C}(\mathbf{p}, \alpha)$  of a spatial point  $\mathbf{p}$  to the surface of an articulated human shape defined by the generative latent code  $\alpha$ . Using an explicit skeleton, we transform the point  $\mathbf{p}$  into the normalized coordinate frames as  $\{\tilde{\mathbf{p}}^j\}$  for  $N = 4$  sub-part networks, modeling body, hands, and head. Each sub-model  $\{S^j\}$  represents a semantic signed-distance function. The sub-models are finally combined consistently using an MLP  $U$  to compute the outputs  $s$  and  $\mathbf{c}$  for the full body. Our multi-part pipeline builds a full body model as well as sub-part models for head and hands, jointly, in a consistent training loop. On the right, we visualize the zero-level-set body surface extracted with marching cubes and the implicit correspondences to a canonical instance given by the output semantics. The semantics allows e.g. for surface coloring or texturing.

related, SCANimate [45] builds personalized avatars from multiple scans of a single person. Concurrent to our work, LEAP [32] learns an occupancy model of human shape and pose also without hand poses, expressions, or semantics. In this work we aim for a full implicit body model, featuring a large range of body shapes corresponding to diverse humans and poses, with detailed hands, and facial expressions.

## 2. Methodology

In this section, we describe our models and the losses used for training. We introduce two variants: a single-part model that encodes the whole human in a single network and a multi-part model. The latter constructs the full body from the output superposition of four body part networks.

**Background.** We rely on neural networks and implicit functions to generate 3D human shapes and articulated poses. Given a latent representation  $\alpha$  of the human shape and pose, together with an underlying probability distribution, we model the posed body as the zero iso-surface decision boundaries of Signed Distance Functions (SDFs) given by deep feed-forward neural networks. A signed distance  $S(\mathbf{p}, \alpha) \in \mathbb{R}$  is a continuous function which, given an arbitrary spatial point  $\mathbf{p} \in \mathbb{R}^3$ , outputs the shortest distance to the surface defined by  $\alpha$ , where the sign indicates the inside (negative) or outside (positive) side w.r.t. the surface. The posed human body surface is implicitly given by  $S(\cdot, \alpha) = 0$ .

GHUM [51] represents the human model as an articulated mesh  $\mathbf{X}(\alpha)$ . GHUM has a minimally-parameterized skeleton with  $J = 63$  joints (124 Euler angle DOFs), and skinning deformations, explicitly sensitive to the pose kinematics  $\theta \in \mathbb{R}^{124}$ . A kinematic prior based on normalizing flows defines the distribution of valid poses [54]. Each kinematic pose  $\theta$  represents a set of joint transformations  $\mathbf{T}(\theta, \mathbf{j}) \in \mathbb{R}^{J \times 3 \times 4}$  from the neutral to a posed state, where  $\mathbf{j} \in \mathbb{R}^{J \times 3}$  are the joint centers that are dependent on the neutral body shape. The statistical body shapes are mod-

eled using a nonlinear embedding  $\beta_b \in \mathbb{R}^{16}$ . In addition to skeleton articulation, a nonlinear latent code  $\beta_f \in \mathbb{R}^{20}$  drives facial expressions. The implicit model we design here shares the same probabilistic latent representation as GHUM,  $\alpha = (\beta_b, \beta_f, \theta)$ , but in contrast to computing an articulated mesh, we estimate a signed distance value  $s = S(\mathbf{p}, \alpha)$  for each arbitrary spatial point  $\mathbf{p}$ .

### 2.1. Models and Training

Given a collection of full-body human meshes  $\mathbf{Y}$ , together with the corresponding GHUM encodings  $\alpha = (\beta_b, \beta_f, \theta)$ , our goal is to learn a MLP-based SDF representation  $S(\mathbf{p}, \alpha)$  so that it approximates the shortest signed distance to  $\mathbf{Y}$  for any query point  $\mathbf{p}$ . Note that  $\mathbf{Y}$  could be arbitrary meshes, such as raw human scans, mesh registrations, or samples drawn from the GHUM latent space. The zero iso-surface  $S(\cdot, \alpha) = 0$  is sought to preserve all geometric detail in  $\mathbf{Y}$ , including body shapes and poses, hand articulation, and facial expressions.

**Single-part Network.** We formulate one global neural network that decodes  $S(\mathbf{p}, \alpha)$  for a given latent code  $\alpha$  and a spatial point  $\mathbf{p}$ . Instead of pre-computing the continuous SDFs from point samples as in DeepSDF [34], we train a MLP network  $S(\mathbf{p}, \alpha; \omega)$  with weights  $\omega$ , similar in spirit to IGR [14], to output a solution to the Eikonal equation

$$\|\nabla_{\mathbf{p}} S(\mathbf{p}, \alpha; \omega)\| = 1, \quad (1)$$

where  $S$  is a signed distance function that vanishes at the surface  $\mathbf{Y}$  with gradients equal to surface normals. Mathematically, we formulate our total loss as a weighted combination of

$$L_o(\omega) = \frac{1}{|O|} \sum_{i \in O} (|S(\mathbf{p}_i, \alpha)| + \|\nabla_{\mathbf{p}_i} S(\mathbf{p}_i, \alpha) - \mathbf{n}_i\|) \quad (2)$$

$$L_e(\omega) = \frac{1}{|F|} \sum_{i \in F} (\|\nabla_{\mathbf{p}_i} S(\mathbf{p}_i, \alpha)\| - 1)^2 \quad (3)$$

$$L_l(\omega) = \frac{1}{|F|} \sum_{i \in F} \text{BCE}(l_i, \phi(kS(\mathbf{p}_i, \alpha))), \quad (4)$$

where  $\phi$  is the sigmoid function,  $O$  are surface samples from  $\mathbf{Y}$  with normals  $\mathbf{n}$ , and  $F$  are off surface samples with inside/outside labels  $l$ , consisting of both uniformly sampled points within a bounding box and sampled points near the surface. The first term  $L_o$  encourages the surface samples to be on the zero-level-set and the SDF gradient to be equal to the given surface normals  $\mathbf{n}_i$ . The Eikonal loss  $L_e$  is derived from (1) where the SDF is differentiable everywhere with gradient norm 1. We obtain the SDF gradient  $\nabla_{\mathbf{p}_i} S(\mathbf{p}_i, \alpha)$  analytically via network back-propagation. In practice, we also find it useful to include a binary cross-entropy error (BCE) loss  $L_l$  for off-the-surface samples, where  $k$  controls the sharpness of the decision boundary. We use  $k = 10$  in our experiments. Our training losses only require surface samples with normals and inside/outside labels for off-surface samples. Those are much easier and faster to obtain than pre-computing ground truth SDF values.

Recent work suggests that standard coordinate-based MLP networks encounter difficulties in learning high-frequency functions, a phenomenon referred to as *spectral bias* [41, 49]. To address this limitation, inspired by [49], we therefore encode our samples using the basic Fourier mapping  $\mathbf{e}_i = [\sin(2\pi\tilde{\mathbf{p}}_i), \cos(2\pi\tilde{\mathbf{p}}_i)]^\top$ , where we first unpose the samples with the root rigid transformation  $\mathbf{T}_0^{-1}$  and normalize them into  $[0, 1]^3$  with a shared bounding box  $\mathbf{B} = [\mathbf{b}_{min}, \mathbf{b}_{max}]$ , as

$$\tilde{\mathbf{p}}_i = \frac{\mathbf{T}_0^{-1}(\theta, \mathbf{j})[\mathbf{p}_i, 1]^\top - \mathbf{b}_{min}}{\mathbf{b}_{max} - \mathbf{b}_{min}}. \quad (5)$$

Note that our SDF is defined w.r.t. the original meshes  $\mathbf{Y}$  and therefore we do not unpose and scale the sample normals. Also, the loss gradients are derived w.r.t.  $\mathbf{p}_i$ .

**Multi-part Network.** Our single-part network represents well the global geometric features for various human body shapes and kinematic poses. However, despite its spatial encoding, the network still has difficulties capturing facial expressions and articulated hand poses, where the SDF has local high-frequency variations. To augment geometric detail on face and hands regions, we therefore propose a multi-part network that decomposes the human body into  $N = 4$  local regions, i.e. the head, left and right hand, and the remaining body, respectively. This significantly reduces spectral frequency variations within each local region allowing the specialized single-part networks to capture local geometric detail. A consistent full-body SDF  $S(\mathbf{p}, \alpha)$  is composed from the local single-part SDF network outputs  $s^j = S^j(\mathbf{p}, \alpha), j \in \{1, \dots, N\}$ .

We follow the training protocol described in §2.1 for each local sub-part network with surface and off-surface samples within a bounding box  $\mathbf{B}^j$  defined for each part. Note that we use the neck and wrist joints as the the root transformation for the head and hands respectively. In GHUM, the joint centers  $\mathbf{j}$  are obtained as a function given

the neutral body shapes  $\bar{\mathbf{X}}(\beta_b)$ . However,  $\bar{\mathbf{X}}$  is not explicitly presented in our implicit representation. Therefore, we build a nonlinear joint regressor from  $\beta_b$  to  $\mathbf{j}$ , which is trained, supervised, using GHUM’s latent space sampling.

In order to fuse the local SDFs into a consistent full-body SDF, while at the same time preserving local detail, we merge the last hidden layers of the local networks using an additional light-weight MLP  $U$ . To train the combined network, a sample point  $\mathbf{p}_i$ , defined for the full body, is transformed into the  $N$  local coordinate frames using  $\mathbf{T}_0^j$  and then passed to the single-part local networks, see fig. 2. The union SDF MLP then aggregates the shortest distance to the full body among the local distances. We apply our losses to the union full-body SDF as well, to ensure that the output for full body satisfies the SDF property (1). Our multi-part pipeline produces sub-part models and a full-body one, trained jointly and leveraging data correlations among different body components.

Our spatial point encoding  $\mathbf{e}_i$  requires all samples  $\mathbf{p}$  to be inside the bounding box  $\mathbf{B}$ , which otherwise might result in periodic SDFs due to sinusoidal encoding. However, a point sampled from the full body is likely to be outside of a sub-part’s local bounding box  $\mathbf{B}^j$ . Instead of clipping or projecting to the bounding box, we augment our encoding of sample  $\mathbf{p}_i$  for sub-part networks  $S^j$  as  $\mathbf{e}_i^j = [\sin(2\pi\tilde{\mathbf{p}}_i^j), \cos(2\pi\tilde{\mathbf{p}}_i^j), \tanh(\pi(\tilde{\mathbf{p}}_i^j - 0.5))]^\top$ , where the last value indicates the relative spatial location of the sample w.r.t. the bounding box. If a point  $\mathbf{p}_i$  is outside the bounding box  $\mathbf{B}^j$ , the union SDF MLP will learn to ignore  $S^j(\mathbf{p}_i^j, \alpha)$  for the final union output.

**Implicit Semantics.** In contrast to explicit models like GHUM, implicit functions do not naturally come with point correspondences between different shape instances. However, many applications, such as pose tracking, texture mapping, semantic segmentation, surface landmarks, or clothing modeling, largely benefit from such correspondences. Given an arbitrary spatial point, on or near the surface  $\mathbf{Y}$ , i.e.,  $|S(\mathbf{p}_i, \alpha)| < \sigma$ , we are therefore interested to interpret its *semantics*. We define the semantics as a 3D implicit function  $\mathbf{C}(\mathbf{p}, \alpha) \in \mathbf{R}^3$ . Given a query point  $\mathbf{p}_i$ , it returns a correspondence point on a canonical GHUM mesh  $\mathbf{X}(\alpha_0)$  as

$$\mathbf{C}(\mathbf{p}_i, \alpha) = \mathbf{w}_i \mathbf{v}_f(\alpha_0) = \mathbf{c}_i, \quad \mathbf{p}_i^* = \mathbf{w}_i \mathbf{v}_f(\alpha) \quad (6)$$

where  $\mathbf{p}_i^*$  is the closest point of  $\mathbf{p}_i$  in the GHUM mesh  $\mathbf{X}(\alpha)$  with  $f$  the nearest face and  $\mathbf{w}$  the barycentric weights of the vertex coordinates  $\mathbf{v}_f$ . In contrast to alternative semantic encodings, such as 2D texture coordinates, our semantic function  $\mathbf{C}(\mathbf{p}, \alpha)$  is smooth in the spatial domain without distortion and boundary discontinuities, which favors the learning process, c.f. [5].

By definition, implicit SDFs return the shortest distance to the underlying implicit surface for a spatial point whereas implicit semantics associate the query point to its closest

surface neighbor. Hence, we consider implicit semantics as highly correlated to SDF learning. We co-train both tasks with our augmented multi-part network (§2.1) computing both  $S(\mathbf{p}, \alpha)$  and  $C(\mathbf{p}, \alpha)$ . Semantics are trained fully supervised, using an  $L_1$  loss for a collection of training sample points near and on the surface  $\mathbf{Y}$ . Due to the correlation between tasks, our network is able to predict both signed distance and semantics, without expanding its capacity.

Using trained implicit semantics, we can e.g. apply textures to arbitrary iso-surfaces at level set  $|z| \leq \sigma$ , reconstructed from our implicit SDF. During inference, an iso-surface mesh  $S(\cdot, \alpha) = z$  can be extracted using Marching Cubes [28]. Then for every generated vertex  $\tilde{\mathbf{v}}_i$  we query its semantics  $C(\tilde{\mathbf{v}}_i, \alpha)$ . The queried correspondence point  $C(\tilde{\mathbf{v}}_i, \alpha)$  might not be exactly on the canonical surface and therefore we project it onto  $\mathbf{X}(\alpha_0)$ . Now, we can interpolate the UV texture coordinates and assign them to  $\tilde{\mathbf{v}}_i$ . Similarly, we can also assign segmentation labels or define on- or near-surface landmarks. In fig. 2 (right) we show an imGHUM reconstruction textured and with a binary ‘clothing’ segmentation. We use the latter throughout the paper demonstrating that our semantics allow the transfer of segmentation labels to different iso-surface reconstructions. Please refer to §3.3 for more applications of our implicit semantics e.g. landmarks or clothed human reconstruction.

**Architecture.** For the single-part network we use a similar feed-forward architecture as DeepSDF [34] or IGR [14] with eight 512-dimensional fully-connected layers. To enable higher-order derivatives, we use Swish nonlinear activation [42] instead of ReLU. IGR originally proposed SoftPlus, however, we found Swish superior (see tab. 3). The multi-part network is composed out of one 8-layer 256-dimensional MLP for the body and three 4-layer 256-dimensional MLPs for hands and head. Each sub-network has a skip connection to the middle layer. The last hidden layers of sub-networks are aggregated in a 128-dimensional fully-connected layer with Swish nonlinear activation, before the final network output. The final model features 2.49 million parameters and performs 4.99 million FLOPs per point query.

**Dataset.** Our training data consists of a collection of full-body human meshes  $\mathbf{Y}$  together with the corresponding GHUM latent code  $\alpha$ , where  $\mathbf{X}(\alpha)$  best approximates  $\mathbf{Y}$ . For each mesh, we perform Poisson disk sampling on the surface and obtain  $|O| = 32K$  surface samples, together with their surface normals. In addition, within a predefined  $2.2 \times 2.8 \times 2.2 \text{ m}^3$  bounding box centered at the origin, we sample  $|F|/2 = 16K$  points uniformly. Another 16K samples are generated by randomly displacing surface sample points with isotropic normal noise with  $\sigma = 0.05\text{m}$ . All off-surface samples are associated with inside/outside labels, computed by casting randomized rays and checking parity.

We also label semantics for on and near surface samples, which are drawn with random face indices and barycentric weights of the GHUM mesh and randomly displaced for near-surface samples. With the corresponding face and barycentric weights, semantic labels are generated using (6) in a light-weight computation with no need for projection or nearest neighbor search. Each mesh  $\mathbf{Y}$  is then decomposed into  $N = 4$  parts and we generate the same number of training samples per body part (we use  $\sigma = 0.02\text{m}$  for surface samples near the hands).

We use two types of human meshes for our imGHUM training. We first randomly sample 75K poses from H36M and the CMU mocap dataset, with Gaussian sampled body shapes, expressions and hand poses from the GHUM latent priors, where  $\mathbf{Y}$  are the posed GHUM meshes. In addition, we collect 35K human scans, on which we perform As-Conformal-As-Possible (ACAP) registrations [53] with the GHUM topology and fit GHUM parameters as well. Our human scans include the CAESAR dataset, full body pose scans, as well as close-up head and hand scans. Due to the noise and incompleteness in some of the raw scans we use the registrations for training. We fine-tune imGHUM – initially trained on GHUM sampling – using the registration dataset. In this way, imGHUM can capture geometric detail not well represented by GHUM (see tab. 2).

### 3. Experiments

We evaluate imGHUM qualitatively and quantitatively in multiple experiments. First, we compare imGHUM with its explicit counterpart GHUM (§3.1). Then, we perform an extensive baseline and ablation study, demonstrating the effect of imGHUM’s architecture and training scheme (§3.2). We also build a model to compare to the recent single-subject occupancy model NASA. Finally, we show the performance of imGHUM on three representative applications demonstrating its usefulness and versatility (§3.3).

We report three different metrics. Bi-directional Chamfer- $L_2$  distance measures the accuracy and completeness of the surface (lower is better). Normal Consistency (NC) evaluates estimated surface normals (higher is better). Volumetric Intersection over Union (IoU) compares the reconstructed volume with the ground truth shape (higher is better). The latter can only be reported for watertight shapes. Please note that metrics not always correlate with the perceived quality of the reconstructions. We therefore additionally include qualitative side-by-side comparisons.

For visualization and numerical evaluation we extract meshes from imGHUM using Marching Cubes [28]. To this end, we approximate the bounding box of the surface through probing and then run Marching Cubes with a resolution of  $256^3$  within the bounding box. Hereby, the signed distances support acceleration using Octree sampling: we use the highest grid density only near the surface and sam-



Figure 3. Bodies generated and reconstructed using imGHUM. *Left:* imGHUM with Gaussian sampling of the shape, expression and pose latent space. *Middle:* Reconstructed motion sequence from the CMU mocap dataset [2] (fixed body shape). *Right:* Body shape and facial expressions latent code interpolation (fixed pose). See supplementary material for more examples.

ple far less frequently away from it. However, we note that for most applications, such as human reconstruction and collision detection, Marching Cubes are not needed, except only once for the final mesh visualization.

### 3.1. Representational Power

In fig. 3, we show reconstructions of a motion capture sequence applied to imGHUM. Our model captures well the articulated full-body motion, with consistent body shape for various poses. By sharing the latent priors with GHUM, imGHUM supports realistic body shape and pose generation (fig. 3, left) as well as smooth interpolation within the shape and expression latent spaces (fig. 3, right). Our model generalizes well to novel body shapes, expressions, and poses, and has interpretable and decoupled latent representations.

In tab. 2, we compare the representation power of imGHUM with the explicit GHUM on our registration test set. imGHUM better captures present detail as numerically demonstrated. An imGHUM model trained only using GHUM samples captures the body deformation due to articulation less well, indicating that GHUM is a useful surrogate to ‘synthetically’ bootstrap the training of the implicit network, but that real data is important as well.

**Limitations** of imGHUM are sometimes apparent for very extreme pose configurations that have not been covered in the training set, such as anthropometrically invalid poses that are impossible for a human, e.g. resulting in self-intersection or by bending joints beyond their anatomical range of motion. imGHUM produces plausible results for inputs not too far from expected configurations, but the results occasionally feature some defects e.g. distorted or incomplete geometry or inaccurate semantics, see fig. 8 for examples.

### 3.2. Baseline Experiments

In the next section, we compare imGHUM to various baselines inspired by recent work. The first is an auto-encoder, where the encoder side is PointNet++ [40] and the decoder is our single-part network. The idea is to let the network find the best representation instead of pre-computing a low dimensional representation. In practice this means that latent codes are not interpretable. Further,

Model	IoU $\uparrow$	Chamfer $\times 10^{-3} \downarrow$	NC $\uparrow$
imGHUM $\ddagger$	0.900	0.071	0.977
GHUM	0.913	0.055	0.983
<b>imGHUM</b>	<b>0.932</b>	<b>0.040</b>	<b>0.984</b>

Table 2. GHUM comparisons on registration dataset. imGHUM marked with  $\ddagger$  is trained only based on GHUM sampling data.

Model	IoU $\uparrow$	Chamfer $\times 10^{-3} \downarrow$	NC $\uparrow$
Autoencoder	0.831	2.457	0.923
Single-part $\ddagger$	0.957	0.085	0.983
Single-part $\oplus$	0.958	0.070	0.983
Single-part	0.965	0.052	0.986
Single-part deeper $\ddagger$	0.961	0.070	0.984
Single-part deeper	0.967	0.058	0.986
imGHUM $\ddagger$	0.955	0.095	0.984
imGHUM w/o $L_l$ (4)	0.966	0.051	0.988
<b>imGHUM</b>	<b>0.969</b>	<b>0.036</b>	<b>0.989</b>

Table 3. Numerical comparison with baselines. Models marked with  $\ddagger$  don’t use the Fourier input mapping.  $\oplus$  marks Softplus activation as in [14].

Model	IoU $\uparrow$	Ch. $\times 10^{-3} \downarrow$	NC $\uparrow$
	Head / Hands	Head / Hands	Head / Hands
Single-part	0.967 / 0.818	0.010 / 0.201	0.937 / 0.790
Single-part deep.	0.968 / 0.832	0.011 / 0.271	0.938 / 0.811
<b>imGHUM</b>	<b>0.976 / 0.929</b>	<b>0.007 / 0.031</b>	<b>0.944 / 0.934</b>

Table 4. Unidirectional metrics (GT to generated mesh) for critical body parts. Our multi-part architecture significantly improves the head and hand reconstruction accuracy.

we experiment with our single part network without Fourier input mapping, largely following the training scheme proposed by IGR [14]. We also use input mapping and finally trained a deeper single-part network variant (10 layers) having roughly the same number of variables as imGHUM.

In tab. 3 we report the metrics for different variants on our test set containing 1 000 GHUM samples. In fig. 4, we show a side-by-side comparison. The Fourier input mapping consistently improves results for all variants. We have also tried higher-dimensional Fourier features but empirically found the basic encoding to work best in our setting. The auto-encoder produces large artifacts especially in the hand region. Similar problems, large blobs or missing pieces, can be observed in results from single-part variants, especially for the hands and, less severe, also for the facial region. These problems, however, are not well captured by

globally evaluating the whole shape. To this end, we evaluate imGHUM and our single-part models specifically for these critical regions, see tab. 4. Only imGHUM consistently produces high-quality results also for hands and the face, supporting the proposed architecture choices.

Next, we compare imGHUM to the recent single-subject multi-pose implicit human occupancy model NASA [11]. With a fixed body shape, we generate 22 500 random GHUM full-body training poses and 2 500 testing poses from Human3.6M [21] and the CMU mocap dataset [2], including head and hand poses. Using the original point sampling strategy in NASA, we have trained the network until convergence, based on the original source code. Please see the supplementary material for details on how we adapted NASA for the GHUM skeleton. For comparison, we have trained an imGHUM architecture with  $2 \times$  fewer layers than our full multi-subject model, each with half-dimensionality, using the same dataset. Even though GHUM-based NASA has  $3 \times$  more parameters, our smaller-size single-subject imGHUM still performs significantly better in representing both the global shape and local detail (see hand reconstructions in fig. 5). In contrast to NASA, which computes binary occupancy, imGHUM returns more informative signed distance values which produce smooth decision boundaries and preserve the detailed geometry much better. Further key differences to NASA are our considerably simpler architecture that requires far less computation to produce a reconstruction, our semantics, and the carefully chosen learning model (i.e. Fourier encoding, second-order losses) that pays particular attention to surface detail. Moreover, imGHUM additionally models body shape, fingers, and facial expressions using generative latent codes (tab. 1).

### 3.3. Applications

We apply imGHUM to three key tasks: body surface reconstruction, partial point cloud completion, and dressed and inclusive human reconstruction.

**Triangle Set Surface Reconstruction.** Given a triangle set ('soup') with  $n$  vertices  $\{\hat{v}\} \in \mathbb{R}^{3n}$  along with oriented normals  $\{\hat{n}\} \in \mathbb{R}^{3n}$ , we deploy our parametric implicit SDF for surface reconstruction with semantics. This task is necessary for triangle soups produced by 3D scanners. To extract the surface from an incomplete scan, we apply a BFGS optimizer to fit  $\alpha = (\beta_b, \beta_f, \theta)$  such that all vertices  $\hat{v}$  are close to the implicit surface  $S(\cdot, \alpha) = 0$ . Moreover, we enforce gradients at  $\hat{v}$  to be close to normals  $\hat{n}$ , and generated off-surface samples to have distances with the expected signs. In addition, we sample near surface points with a small distance  $\eta$  along surface normals, and enforce  $S(\hat{v} \pm \eta\hat{n}, \alpha) = \pm\eta$ , as in [34]. Note that all these operations can be easily implemented and are fully differential due to imGHUM being a SDF. When 3D landmarks are available on the target surface, e.g. as triangulated from

(Without Fourier input mapping)

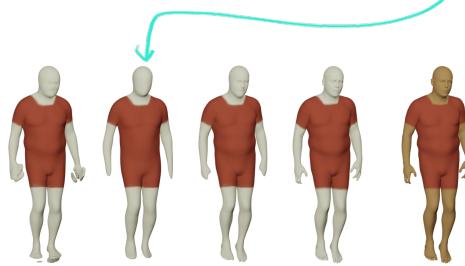


Figure 4. Qualitative comparison with baseline experiments. From left to right: autoencoder, single-part model without and with Fourier input mapping, our multi-part imGHUM, ground-truth GHUM. We use our semantics network to color baseline results.

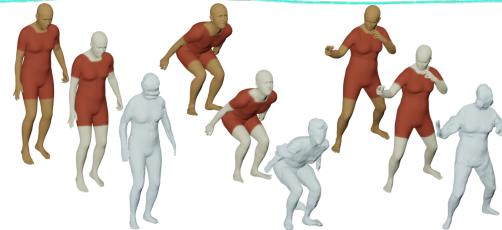


Figure 5. Comparison with NASA [11] on our single-subject multi-pose dataset. Top to bottom: GT, single-subject imGHUM, and NASA reconstructions. imGHUM better captures global and local geometry, despite using a significantly smaller network version in this experiment. Also numerically our results are superior: IoU ( $\uparrow$ ) 0.962 (ours) vs. 0.839 (theirs), Ch. ( $\downarrow$ )  $0.068 \times 10^{-3}$  (ours) vs.  $3.53 \times 10^{-3}$  (theirs), NC ( $\uparrow$ ) 0.985 (ours) vs. 0.903 (theirs).

2D detected landmarks of raw scanner images, we additionally augment the optimization with landmark losses based on the imGHUM semantics. Please see the supplementary material for details of the losses.

For reference, we also show results on IF-Net [8], a recent method for implicit surface extraction, completion, and voxel super-resolution. We trained IF-Net with the same pose and shape variation as used for imGHUM – presumably much more variation than the 2 183 scans in the original paper. In both training and testing we generate 15K random samples from the observed shape and pass them through IF-Net for surface reconstruction. Note that IF-Net is using less information compared to our method, but is also solving an easier task as it is not computing a global and semantically meaningful shape code. An entirely fair comparison is thus not possible. However, we believe that by comparing with IF-Net, we show that imGHUM is adequate for this task. Fig. 6 qualitatively shows examples of both imGHUM fits and of IF-Net inference results for 150 human scans containing 20 subjects. Our model not only fits well to the volume of the scans but also reconstructs the facial expressions and hand poses. Using landmarks and ICP losses, one could also fit GHUM to the triangle sets. However, our fully differential imGHUM losses show superior performance over ICP-based GHUM fitting (Chamfer ( $\downarrow$ )  $0.77 \times 10^{-3}$ , NC ( $\uparrow$ ) 0.921).

**Partial Point Cloud Completion.** Another relevant task for many applications is shape completion. Here we show



Figure 6. Left: Triangle set surface reconstruction (input scan, imGHUM fit, and IF-Net inference from left to right). Numerically, imGHUM fits are better than IF-Net with Chamfer distance ( $\downarrow$ )  $0.156 \times 10^{-3}$  (ours) vs.  $0.844 \times 10^{-3}$  (IF-Net), and NC ( $\uparrow$ ) 0.954 (ours) vs. 0.914 (IF-Net). Right: Partial point cloud completion (input point cloud, imGHUM fit, IF-Net, and ground truth scan).

surface reconstruction and completion from partial point clouds as recorded e.g. using a depth sensor. We synthesize depth maps from A-posed scans of 10 subjects from the Faust dataset [6] using the intrinsics and the resolution of a Kinect V2 sensor. To complete the partial view, we search for the  $\alpha$  such that all points from the depth point cloud are close to imGHUM’s zero-level-set. We sample additional points along surface normals (estimated from depth image gradients) and enforce estimated distances by imGHUM to be close to true distances. We also sample points in front of the depth cloud and around it and enforce their  $L_1$  label loss. Finally, we also supervise the estimated normals. We do not rely on landmarks or other semantics in this experiment.

We show IF-Net [8] results for comparison. We trained IF-Net specifically for this task while we use the same imGHUM for all experiments. Our reconstructions are numerically better with Chamfer distance ( $\downarrow$ )  $0.103 \times 10^{-3}$  (ours) vs.  $0.315 \times 10^{-3}$  (theirs) and NC ( $\uparrow$ ) 0.962 (ours) vs. 0.936 (theirs). Qualitatively, our results contain much more of the desirable reconstruction detail, especially for hands and faces, see fig. 6, right. Note, again, that IF-Net only reconstructs a surface while we recover the parametrization of a body model, a considerably harder task.

**Dressed and Inclusive Human Modeling.** imGHUM is template-free which is a valuable property for future developments. While this work deals primarily with the methodology of learning a generative implicit human model – in itself a complex and novel task – we also give an outlook for possible future directions. Building a detailed model of the human body shape including hair and clothing, or learning inclusive models could be such directions. However, currently the data needed for building such models does not exist at large enough scale. To demonstrate that imGHUM is a valuable building block for such models, we leverage it as an inner layer for personalized human models. Concretely, we augment imGHUM with a light-weight residual SDF network, conditioned on the output of imGHUM, both the signed distances and semantics. We estimate the residual model using the same learning scheme as for imGHUM, but limit training to a single scan. The final output models the human with layers, including the inner body shape represented with imGHUM and the personalization (hair, clothing, non-standard body topology) as residuals, *c.f.* fig. 7. This layered representation can be reposed by changing the parameterization of the underlying imGHUM. Hereby,

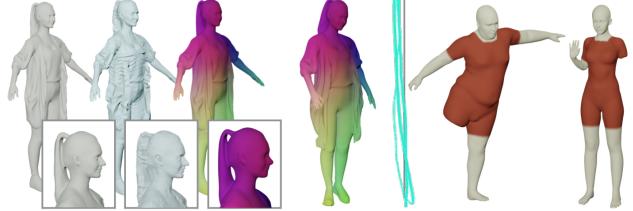


Figure 7. From left to right: scan, GHUM template mesh ACAP registration, imGHUM+residual fit (color-scale represents semantics), reposed imGHUM+residual, imGHUM+residual fits to people with limb differences. In contrast to the fitted template mesh, imGHUM+residual successfully models topologies different from the plain human body and captures more geometric detail.



Figure 8. Failure modes. Interpenetration can lead to unwanted shapes and semantics (leaked hand semantics to the cheek). Extreme poses may produce deformed body parts (thin arms).

residual model acts as a fitted layer around imGHUM and deforms according to the distance and semantic field defined by imGHUM. Please see the supplementary material for more examples, a numerical evaluation, and implementation details.

## 4. Discussion and Conclusion

We introduced imGHUM, the first 3D human body model, with controllable pose and shape, represented as an implicit signed distance function. imGHUM has comparable representation power to state-of-the-art mesh-based models and can represent significant variations in body pose, shape, and facial expressions, as well as underlying, precise, semantics. imGHUM has additional valuable properties, since its underlying implicit SDF represents not only the surface of the body but also its neighborhood, which e.g. enables collision tests with other objects or efficient distance losses. imGHUM can be used to build diverse, fair models of humans who may not match a standard template. This paves the way for transformative research and inclusive applications like modeling clothing, enabling immersive virtual apparel try-on, or free-viewpoint photorealistic visualization. Our models are available for research [1].

## References

- [1] <https://github.com/google-research/google-research/tree/master/imghum>. 2, 8
- [2] CMU graphics lab motion capture database. 2009. <http://mocap.cs.cmu.edu/>. 6, 7
- [3] Thiemo Alldieck, Marcus Magnor, Bharat Lal Bhatnagar, Christian Theobalt, and Gerard Pons-Moll. Learning to reconstruct people in clothing from a single RGB camera. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1175–1186. IEEE, 2019. 1
- [4] Matan Atzmon and Yaron Lipman. Sal: Sign agnostic learning of shapes from raw data. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2565–2574, 2020. 2
- [5] Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Loopreg: Self-supervised learning of implicit surface correspondences, pose and shape for 3d human mesh registration. In *Adv. Neural Inform. Process. Syst.*, 2020. 4
- [6] Federica Bogo, Javier Romero, Matthew Loper, and Michael J. Black. FAUST: Dataset and evaluation for 3D mesh registration. In *IEEE Conf. Comput. Vis. Pattern Recog.* IEEE, 2014. 8
- [7] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5939–5948, 2019. 1, 2
- [8] Julian Chibane, Thiemo Alldieck, and Gerard Pons-Moll. Implicit functions in feature space for 3d shape reconstruction and completion. In *IEEE Conf. Comput. Vis. Pattern Recog.* IEEE, jun 2020. 1, 2, 7, 8, 14
- [9] Julian Chibane, Aymen Mir, and Gerard Pons-Moll. Neural unsigned distance fields for implicit function learning. In *Adv. Neural Inform. Process. Syst.*, December 2020. 2
- [10] Julian Chibane and Gerard Pons-Moll. Implicit feature networks for texture completion from partial 3d data. In *Eur. Conf. Comput. Vis. Worksh.* Springer, August 2020. 2
- [11] Boyang Deng, JP Lewis, Timothy Jeruzalski, Gerard Pons-Moll, Geoffrey Hinton, Mohammad Norouzi, and Andrea Tagliasacchi. Neural articulated shape approximation. In *Eur. Conf. Comput. Vis.* Springer, August 2020. 2, 7, 14
- [12] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Local deep implicit functions for 3d shape. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4857–4866, 2020. 2
- [13] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T Freeman, and Thomas Funkhouser. Learning shape templates with structured implicit functions. In *Int. Conf. Comput. Vis.*, pages 7154–7164, 2019. 2
- [14] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. In *Int. Conf. on Mach. Learn.*, pages 3569–3579. 2020. 2, 3, 5, 6
- [15] John C Hart. Sphere tracing: A geometric method for the antialiased ray tracing of implicit surfaces. *The Visual Computer*, 12(10):527–545, 1996. 14
- [16] Nils Hasler, Carsten Stoll, Martin Sunkel, Bodo Rosenhahn, and H-P Seidel. A statistical model of human pose and body shape. *Comput. Graph. Forum*, 28(2):337–346, 2009. 2
- [17] N. Hesse, S. Pujades, M.J. Black, M. Arens, U. Hofmann, and S. Schroeder. Learning and tracking the 3D body shape of freely moving infants from RGB-D sequences. *IEEE Trans. Pattern Anal. Mach. Intell.*, 42(10):2540–2551, 2020. 1
- [18] David A Hirshberg, Matthew Loper, Eric Rachlin, and Michael J Black. Coregistration: Simultaneous alignment and modeling of articulated 3D shape. In *Eur. Conf. Comput. Vis.*, pages 242–255, 2012. 2
- [19] Zeng Huang, Tianye Li, Weikai Chen, Yajie Zhao, Jun Xing, Chloe LeGendre, Linjie Luo, Chongyang Ma, and Hao Li. Deep volumetric video from very sparse multi-view performance capture. In *Eur. Conf. Comput. Vis.*, pages 336–354, 2018. 2
- [20] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. Arch: Animatable reconstruction of clothed humans. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3093–3102, 2020. 2
- [21] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6M: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2014. 7
- [22] Arjun Jain, Thorsten Thormählen, Hans-Peter Seidel, and Christian Theobalt. Moviereshape: Tracking and reshaping of humans in videos. *ACM Trans. Graph.*, 29(5), 2010. 1
- [23] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, and Thomas Funkhouser. Local implicit grid representations for 3d scenes. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 6001–6010, 2020. 2
- [24] Hanbyul Joo, Tomas Simon, and Yaser Sheikh. Total capture: A 3D deformation model for tracking faces, hands, and bodies. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 8320–8329. IEEE, 2018. 1, 2
- [25] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018. 1
- [26] Korrawe Karunratanakul, Jinlong Yang, Yan Zhang, Michael J Black, Krikamol Muandet, and Siyu Tang. Grasping field: Learning implicit representations for human grasps. In *Int. Conf. on 3D Vis.*, pages 333–344. IEEE, 2020. 2
- [27] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *Int. Conf. Learn. Represent.*, 2015. 11
- [28] Thomas Lewiner, Hélio Lopes, Antônio Wilson Vieira, and Geovan Tavares. Efficient implementation of marching cubes’ cases with topological guarantees. *Journal of Graphics Tools*, 8(2):1–15, 2003. 5, 11
- [29] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graph.*, 34(6):248:1–248:16, 2015. 1, 2
- [30] Lars M. Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4460–4470, 2019. 1, 2

- [31] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotagh, and Anders Eriksson. Deep level sets: Implicit surface representations for 3D shape inference. *arXiv preprint arXiv:1901.06802*, 2019. [1](#) [2](#)
- [32] Marko Mihajlovic, Yan Zhang, Michael J Black, and Siyu Tang. LEAP: Learning articulated occupancy of people. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. [3](#)
- [33] Ahmed A A Osman, Timo Bolkart, and Michael J. Black. STAR: A spare trained articulated human body regressor. In *Eur. Conf. Comput. Vis.*, 2020. [2](#)
- [34] Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 165–174, 2019. [1](#) [2](#) [3](#) [5](#) [7](#)
- [35] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3D hands, face, and body from a single image. In *IEEE Conf. Comput. Vis. Pattern Recog.* IEEE, 2019. [2](#)
- [36] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *Eur. Conf. Comput. Vis.*, Cham, 2020. Springer International Publishing. [2](#)
- [37] Leonid Pishchulin, Stefanie Wuhrer, Thomas Helten, Christian Theobalt, and Bernt Schiele. Building statistical shape spaces for 3d human modeling. *Pattern Recognition*, 67:276–286, 2017. [2](#)
- [38] Ralf Plankers and Pascal Fua. Articulated soft objects for video-based body modeling. In *Int. Conf. Comput. Vis.*, pages 394–401. IEEE, 2001. [2](#)
- [39] Gerard Pons-Moll, Javier Romero, Naureen Mahmood, and Michael J Black. Dyna: a model of dynamic human shape in motion. *ACM Trans. Graph.*, 34:120, 2015. [2](#)
- [40] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Adv. Neural Inform. Process. Syst.*, pages 5099–5108, 2017. [6](#) [11](#)
- [41] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *Int. Conf. on Mach. Learn.*, pages 5301–5310. PMLR, 2019. [4](#)
- [42] Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. *arXiv preprint arXiv:1710.05941*, 2017. [5](#)
- [43] 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 *Int. Conf. Comput. Vis.*, pages 2304–2314, 2019. [2](#)
- [44] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. [2](#)
- [45] Shunsuke Saito, Jinlong Yang, Qianli Ma, and Michael J. Black. SCANimate: Weakly supervised learning of skinned clothed avatar networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. [3](#)
- [46] Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *Adv. Neural Inform. Process. Syst.*, 2020. [2](#)
- [47] Cristian Sminchisescu, Atul Kanaujia, and Dimitris Metaxas. Learning joint top-down and bottom-up processes for 3D visual inference. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1743–1752. IEEE, 2006. [2](#)
- [48] Cristian Sminchisescu and Alexandru Telea. Human pose estimation from silhouettes. a consistent approach using distance level sets. In *10th Int. Conf. on Computer Graphics, Visualization and Computer Vision*, 2002. [1](#) [2](#)
- [49] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Adv. Neural Inform. Process. Syst.*, 2020. [4](#)
- [50] Daniel Thalmann, Jianhua Shen, and Eric Chauvineau. Fast realistic human body deformations for animation and VR applications. In *Proceedings of CG International'96*, pages 166–174, 1996. [2](#)
- [51] Hongyi Xu, Eduard Gabriel Bazavan, Andrei Zanfir, William T Freeman, Rahul Sukthankar, and Cristian Sminchisescu. GHUM & GHUML: Generative 3d human shape and articulated pose models. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 6184–6193, 2020. [1](#) [2](#) [3](#)
- [52] Yifan Xu, Tianqi Fan, Yi Yuan, and Gurprit Singh. Ladybird: Quasi-monte carlo sampling for deep implicit field based 3d reconstruction with symmetry. In *Eur. Conf. Comput. Vis.*, 2020. [2](#)
- [53] Yusuke Yoshiyasu, Wan-Chun Ma, Eiichi Yoshida, and Fumio Kanehiro. As-conformal-as-possible surface registration. *Comput. Graph. Forum*, 2014. [5](#)
- [54] 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 *Eur. Conf. Comput. Vis.*, 2020. [1](#) [3](#) [11](#)
- [55] Mihai Zanfir, Elisabeta Oneata, Alin-Ionut Popa, Andrei Zanfir, and Cristian Sminchisescu. Human synthesis and scene compositing. In *AAAI*, pages 12749–12756, 2020. [1](#)
- [56] Tiancheng Zhi, Christoph Lassner, Tony Tung, Carsten Stoll, Srinivasa G. Narasimhan, and Minh Vo. TexMesh: Reconstructing detailed human texture and geometry from RGB-D video. In *Eur. Conf. Comput. Vis.* Springer International Publishing, 2020. [1](#)
- [57] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5745–5753, 2019. [12](#)

## Supplementary Material

In this supplementary material, we detail our implementation and used architectures, show further results, discuss further ablations, and give more details on our experiments and comparisons. We also demonstrate how imGHUM can be used for differentiable rendering.

### A. Implementation Details

In the following, we detail the implementation of imGHUM. We specify the used hyper-parameters and the architectures used in the ablation experiments. Finally, we give running times for imGHUM mesh extraction via Marching Cubes [28].

**Hyper-parameters.** We train imGHUM with a batch-size of 32, each of which contains 32 instances of  $\alpha$  paired with 512 on surface, 256 near surface, and 256 uniform samples for each instance. Our loss is composed as

$$L = \lambda_{o_1} L_{o_1} + \lambda_{o_2} L_{o_2} + \lambda_e L_e + \lambda_l L_l + \lambda_s L_s, \quad (7)$$

where  $L_{o_1}$  refers to the first part of  $L_o$  (distance) and  $L_{o_2}$  to the second part (gradient direction), respectively, and  $L_s$  refers to the semantics loss. We choose  $\lambda_{o_1} = 1$ ,  $\lambda_{o_2} = 1$ ,  $\lambda_e = 0.1$ ,  $\lambda_l = 0.1$ , and  $\lambda_s = 0.5$ . Empirically we found that linearly increasing  $\lambda_{o_1}$  to 50 over 100K iterations leads to perceptually better results. We train imGHUM until convergence using the Adam optimizer [27] with a learning rate of  $0.2 \times 10^{-3}$  exponentially decaying by a factor of 0.9 over 100K iterations.

**Architectures.** The following architectures have been used for the baseline experiments: The single-part network has been used as described in the main paper totaling in 2.01M parameters. The deeper single-part network uses 10 instead of 8 layers, resulting in 2.53M parameters. The autoencoder is composed from a PointNet++ [40] encoder and our single-part decoder with a total number of parameters of 3.91M. The encoder consists of three PointNet++ set abstraction modules and two 512-dimensional fully-connected layers with ReLU activation.

**Running Times.** We extract meshes from imGHUM using Octree sampling. Reconstructing a mesh in its bounding box and with a maximum grid resolution of  $256^3$  takes on average 1.08s using a NVIDIA Tesla V100. Hereby, the network query time sums up to 0.44s and Marching Cubes [28] (on CPU) takes 0.34s. The rest of the time is used by identifying the bounding box through probing (0.17s), Octree logic (0.05s), and transforming the samples to the part reference frames (0.07s). We query imGHUM in batches with a maximum batch-size of  $64^3$  samples, where one full batch takes on average 0.13s to

Model	IoU $\uparrow$	Chamfer $\times 10^{-3} \downarrow$	NC $\uparrow$
Only scan registrations	0.901	0.091	0.975
imGHUM	<b>0.932</b>	<b>0.040</b>	<b>0.984</b>

Table 5. Numerical comparison of imGHUM trained with different data distributions evaluated on the registration test-set.

Model	IoU $\uparrow$	Chamfer $\times 10^{-3} \downarrow$	NC $\uparrow$
Only scan registrations	0.834	2.561	0.942
imGHUM	<b>0.969</b>	<b>0.036</b>	<b>0.989</b>

Table 6. Numerical comparison of imGHUM trained with different data distributions evaluated on the GHUM samples test-set.

compute. The resulting meshes feature approximately 100K vertices and 200K facets. Note that imGHUM allows creating meshes in arbitrary resolutions and can be queried and also rendered (*c.f.* §D.3) without generating an explicit mesh. For reference, we show imGHUM mesh reconstructions in different resolutions in fig. 9.

## B. Results

In this supplemental material we show additional results for our application experiments (fig. 11, 12, 14, 15). Additionally, fig. 10 displays a large number of imGHUM instances with great variety in poses, shapes, hand poses, and facial expressions sampled from imGHUM’s generative latent space. This demonstrates once more that imGHUM’s level of detail, expressiveness and generative power is on par with state-of-the-art mesh-based models. Moreover, imGHUM can additionally be queried at arbitrary resolutions and spatial locations and models not only the surface, but also the space around the person.

## C. Ablations

In this section, we report further results of our dataset ablation experiment and results of an additional ablation study on joint rotation parameterization.

**Dataset.** In the main paper we have shown that imGHUM benefits from being trained on both samples of GHUM and additionally on As-Conformal-As-Possible (ACAP) registrations of a corpus of human scans. While training only on scan data can represent the distribution of the scans well (tab. 5), it does not generalize sufficiently to poses that are not covered in this limited training set, as we show in tab. 6.

In fig. 13, we qualitatively show the effect of fine-tuning with scan data. Please note the increased level of detail in the faces and the enhanced soft-tissue deformation.

**Rotation Representations.** In tab. 7, we report metrics for imGHUM using different rotation representations for joint rotations  $\theta$ . We have experimented with Euler angles, basic sin, cos Fourier mapping [54], and the recently



Figure 9. imGHUM mesh reconstructions in different resolutions. Left to right: ground-truth shape,  $512^3$ ,  $256^3$ ,  $128^3$ ,  $64^3$ .

Model	IoU $\uparrow$	Chamfer $\times 10^{-3} \downarrow$	NC $\uparrow$
6D	<b>0.969</b>	0.044	<b>0.989</b>
sin, cos	0.967	0.046	0.988
Euler	<b>0.969</b>	<b>0.036</b>	<b>0.989</b>

Table 7. Numerical comparison of imGHUM models using different representations for joint angles evaluated on the GHUM samples test-set.

proposed 6D representation [57]. Perhaps surprisingly, we found only minor differences in imGHUM’s representational power using different rotation representations, both qualitatively and quantitatively. We, therefore, use Euler angles in this work as it is the most compact representation.

## D. Applications

In the following, we explain the losses used in our triangle set surface reconstruction experiment, detail the residual model of the dressed and inclusive modeling experiment, and finally introduce another application namely pose estimation from silhouettes using differentiable rendering.

### D.1. Triangle Set Surface Reconstruction

We describe our triangle set surface reconstruction experiment in the main paper (§3.3) and show more examples here in fig. 11. Our imGHUM reconstructions are performed under a weighted combination of losses as

$$\min_{\alpha} L_o(\alpha) + L_l^+(\alpha) + L_l^-(\alpha) \quad (8)$$

$$L_o(\alpha) = \frac{1}{n} \sum_i |S(\hat{\mathbf{v}}_i, \alpha)| + \|\nabla_{\hat{\mathbf{v}}_i} S(\hat{\mathbf{v}}_i, \alpha) - \hat{\mathbf{n}}_i\| \quad (9)$$

$$L_l^+(\alpha) = \frac{1}{n} \sum_i (\phi(kS(\hat{\mathbf{v}}_i + \gamma_i \hat{\mathbf{n}}_i, \alpha)) - 1)^2 \quad (10)$$

$$L_l^-(\alpha) = \frac{1}{n} \sum_i (\phi(kS(\hat{\mathbf{v}}_i - \gamma_i \hat{\mathbf{n}}_i, \alpha)))^2, \quad (11)$$

where  $L_o$  is a surface sample loss (similar to eq. 2 in the main paper), and  $L_l^+, L_l^-$  are sign classification losses defined for points sampled along and opposite to the normals respectively ( $\gamma_i \in [0, 0.05]$  is a Gaussian sampled distance).

Enabled by the implicit semantics of imGHUM, we can additionally exploit landmark losses as,

$$L_j(\alpha) = \frac{1}{|M_j|} \sum_{i \in M_j} \|\mathbf{T}_i(\alpha) \mathbf{j}_i(\alpha) - \mathbf{m}_{j,i}\|^2 \quad (12)$$

$$L_s(\alpha) = \frac{1}{|M_s|} \sum_{i \in M_s} \|\mathbf{C}(\mathbf{m}_{s,i}, \alpha) - \bar{\mathbf{m}}_{s,i}\|^2, \quad (13)$$

where  $M_j = \{\mathbf{m}_j\}$ ,  $M_s = \{\mathbf{m}_s\}$  are a collection of 3D landmarks defined over the joints and the surface, respectively.  $\bar{\mathbf{m}}_s$  are the corresponding surface landmarks defined on the canonical mesh  $\mathbf{X}(\alpha_0)$ .  $L_j$  aligns the transformed joint centers with the joint landmarks. The surface landmarks loss  $L_s$  queries the semantics for the ground-truth surface landmarks  $\mathbf{m}_s$  conditioned on  $\alpha$ . The semantics describe the position of the landmarks  $\mathbf{m}_s$  w.r.t. the canonical mesh, and thus should match their correspondences  $\bar{\mathbf{m}}_s$ .

Given a triangle set mesh, one could also fit GHUM with landmarks and ICP losses. However, we note that imGHUM is not only able to perform equivalently on the landmark losses to mesh-based representations, but also exploits more information of the triangle set with its differential losses (eq. (8)) compared to ICP. The process of finding the nearest point for ICP at each optimization iteration is non-differentiable and the accuracy of the nearest point correspondences are highly sensitive to the initialization. In contrast, our imGHUM losses are fully differential everywhere and also exploit additional information encoded in the surface normals and the sign labels. Numerical comparisons are reported in §3.3 of the main paper.

### D.2. Dressed and Inclusive Human Modeling

In the following, we detail our dressed and inclusive modeling experiment from the main paper. We also show more results in fig. 14. In order to learn a personalized shape of a given scan, we augment imGHUM with an MLP  $\hat{S}$  consisting of four 256-dimensional layers. Each layer is followed by Swish nonlinear activation, and a skip connection is added to the middle layer.  $\hat{S}$  modulates the signed distance field of the body to match the scan. These distance residuals could come from clothing, hair, other ap-



Figure 10. Random imGHUM full-body and part instances sampled from imGHUM’s generative latent codes. On the right, we show textured examples. Texturing and binary coloring is enabled by imGHUM’s semantics.



Figure 11. More examples for the triangle set surface reconstruction experiment. Each pair shows the ground truth scan (left) and our reconstruction (right). Notice the reconstructed facial expressions and hand poses.

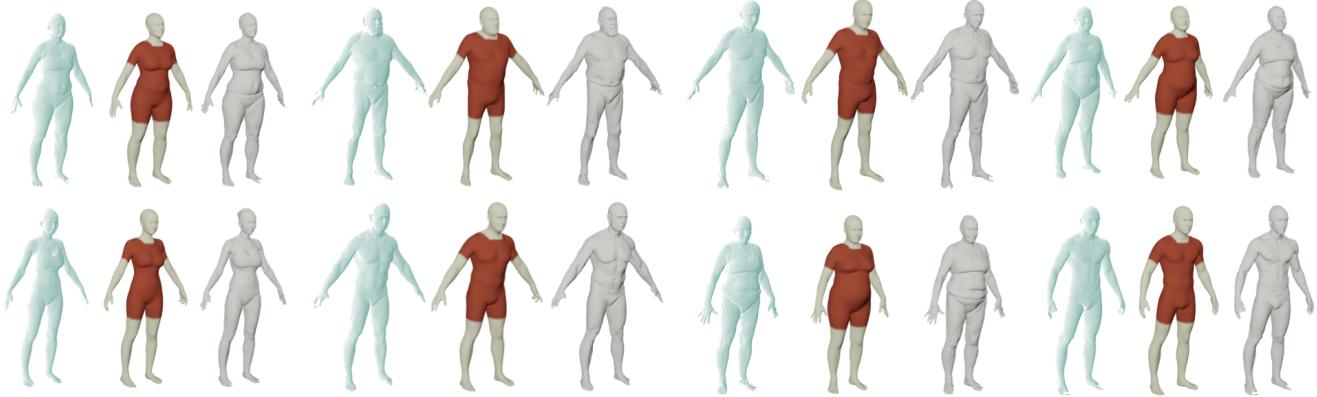


Figure 12. Remaining partial point cloud completion results. Left to right: input point cloud, imGHUM fit, and ground truth scan.

parel items, or any divergence from the standard human template. We condition the output signed distance of the scan with both the distance and semantics fields of the body defined by imGHUM:

$$\hat{s} = \hat{S}(S(\mathbf{p}, \boldsymbol{\alpha})) = \hat{S}(s, \mathbf{c}). \quad (14)$$

We first fit imGHUM to the scan, similar to the trinagle set surface reconstruction experiment. Next, we train  $\hat{S}$  on top of it. The training process is similar to imGHUM with the difference that we sample points from a single scan containing the desired personalizations. We only train the residual while keeping imGHUM fixed. We, therefore, have both the underlying human body and the personalized shape modeled separately as layers. We train a separate instance of  $\hat{S}$  for each scan observation. Learning a combined model using an auto-decoder style learning scheme is possible but beyond the scope of this work.

We show two categories of personalizations: dressed humans and humans with limb differences. We compare imGHUM+residual with mesh-based GHUM ACAP registrations. In contrast to our template-free imGHUM+residual model, GHUM ACAP registrations have difficulties in explaining complex and layered structure and unsurprisingly fail entirely for large structural changes. We fit to scans of ten subjects with limb differences and 30 dressed human scans. Numerically, imGHUM+residual performs better than GHUM ACAP registrations with Chamfer distance  $0.014 \times 10^{-3}$  (ours, limb differences) /  $0.018 \times 10^{-3}$  (ours, dressed) versus  $1.393 \times 10^{-3}$  (GHUM ACAP, limb differences) /  $0.021 \times 10^{-3}$  (GHUM ACAP, dressed) and Normal Consistency 0.993 (ours, limb differences) / 0.990 (ours, dressed) versus 0.984 (GHUM ACAP, limb differences) / 0.976 (GHUM ACAP, dressed). imGHUM+residual is especially superior in explaining the scans of people with limb differences, due to large structural differences compared to the GHUM template mesh. Also qualitatively imGHUM+residual explains much more of the detail present in the input scans, see fig. 14 and fig. 15.

### D.3. Differentiable Rendering

A benefit of imGHUM’s SDF representation is the potential for rendering using sphere tracing [15]. During ray tracing the surface is located by stepping from the camera along a ray until a surface is passed. In sphere tracing the save step length is given by the current minimal distance to any point on the surface, i.e. the SDF value at the current location. For inexact SDFs, one can take a damped step to reduce the likelihood of over-shooting. Using this technique we can render among other things: imGHUM depth maps, normal maps, and semantics. Hereby, each pixel contains the last queried value of its corresponding camera ray. In the following, we compute differentiable binary silhouettes via sphere tracing and fit imGHUM to images using a silhouette alignment loss.

ettes via sphere tracing and fit imGHUM to images using a silhouette alignment loss.

We implement differentiable approximate sphere tracing by taking a fixed number of steps. Concretely, we step  $T = 15$  save steps into the SDF in the direction of each camera ray. At each final point  $\mathbf{p}_T$  of each camera ray, we query the signed distance value and generate the binarized pixel as:

$$b = \frac{1}{\eta S(\mathbf{p}_T, \boldsymbol{\alpha})^2 + 1}, \quad (15)$$

with  $\eta = 5000$  in our experiment.  $b$  is differentiable w.r.t.  $\boldsymbol{\alpha}$  and thus can be used in optimization losses. We formulate a standard silhouette overlap loss and a sparse 2D joint landmark loss and use both to fit imGHUM to image evidence. Fig. 16 shows results of fitting imGHUM to image silhouettes.

## E. Details on Compared Methods

As reported in the main paper, we change NASA [11] in contrast to their original version. Firstly, we train NASA based on the GHUM skeleton containing 63 parts. Originally, NASA was trained on SMPL containing only 24 parts. Another difference is the topology of GHUM. In contrast to SMPL, GHUM features an oral cavity that is also represented in our training data. Summarizing, we deploy NASA for a higher-dimensional model and thus a harder task. For a fair comparison, we therefore use a larger and deeper architecture with eight 64-dimensional fully-connected layers for each part instead of the original four 40-dimensional layers. The new architecture features 1.92M parameters (original version has 0.38M) and has shown significantly better representation power. In contrast, we use a much smaller imGHUM architecture in this experiment. imGHUM has been originally designed to also explain shape variation and facial expressions. Since this experiment only features variation in pose, we can use a much smaller version. We use 2× fewer layers in each part, each with half-dimensionality, resulting in only 0.64M parameters. This smaller-size imGHUM still performs significantly better than NASA in our experiments.

We have trained IF-Net [8] based on their original source code. Specifically, we use IF-Net for point clouds with  $128^3$  resolution featuring 2.6M parameters. We also follow their sampling and resizing strategy, such that the input point cloud always has a maximum side length of one unit. Finally, we train IF-Net task-specific (for full and partial point clouds), while we use the same imGHUM in all our comparisons.

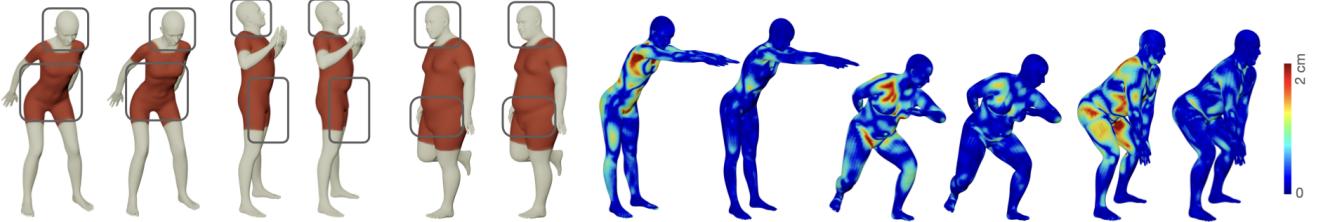


Figure 13. Added detail after fine-tuning on the registration dataset. We show imGHUM reconstructions before fine-tuning (left) and after fine-tuning (right) qualitatively and using error heat-maps (red means  $\geq 2\text{cm}$ ). Please pay attention to the faces, body shapes, and soft-tissue deformations (digital zoom in recommended).

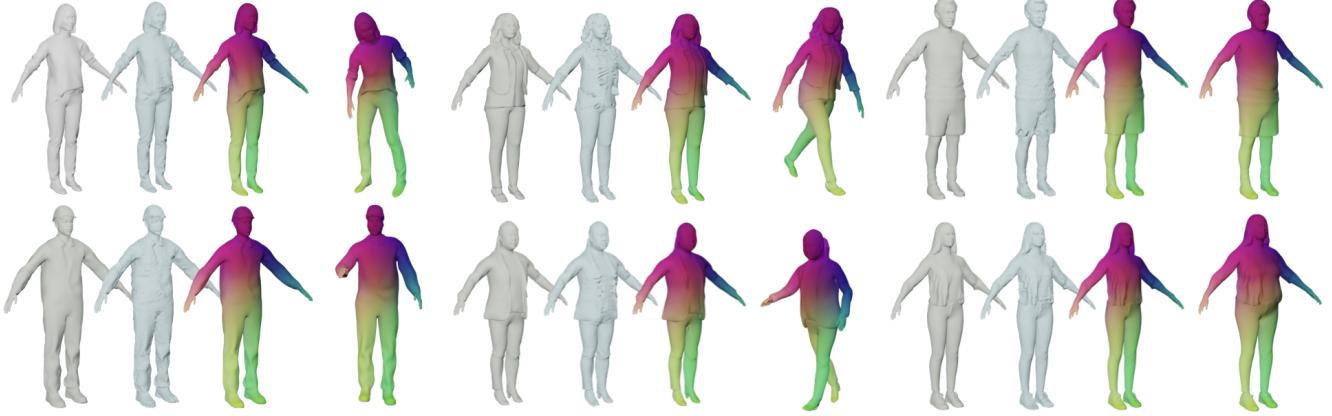


Figure 14. Dressed human modeling. From left to right: Scan, GHUM ACAP mesh registration, imGHUM+residual fit, reposed or reshaped imGHUM+residual. imGHUM+residual accurately explains all detail present in the input scan. GHUM ACAP mesh registrations have difficulties with complicated and layered structures. By changing the parameterization of the underlying imGHUM, we can repose and reshape the personalized models. The color-scale represents imGHUM semantics and thus correspondences between different instances.

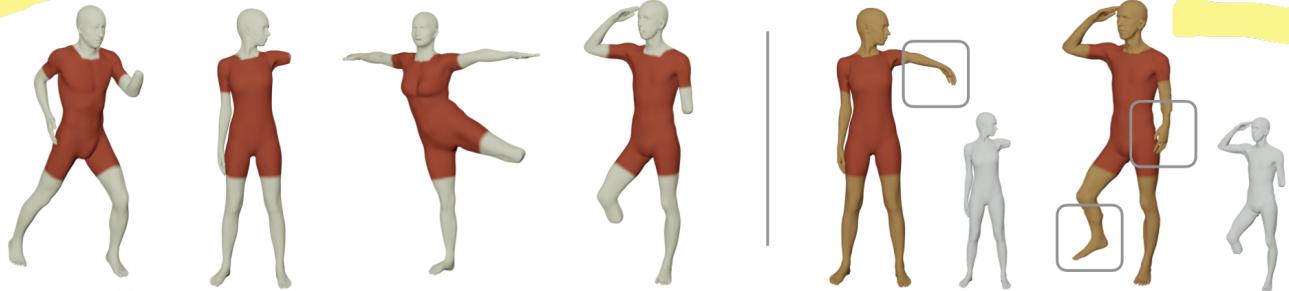


Figure 15. Inclusive human modeling. *Left:* imGHUM+residual can explain body shapes that do not match the standard template. *Right:* GHUM ACAP mesh registrations fail to explain these body shapes. For reference, we show ground truth scans in small. Missing limbs are deformed but still present.

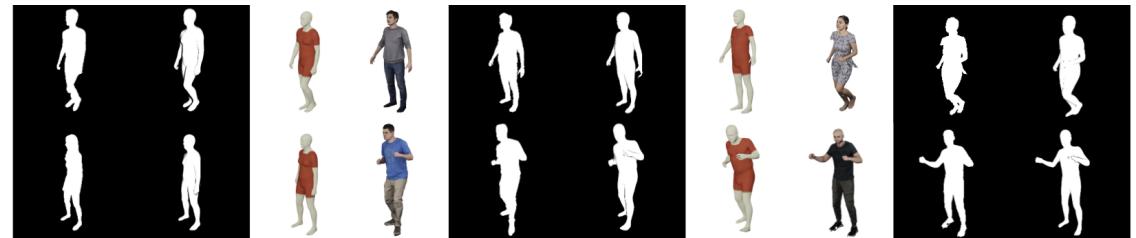


Figure 16. Visual 3D reconstruction of imGHUM using differentiable silhouette and landmark losses. Left to right: image, observed silhouette, estimated silhouette, imGHUM reconstruction. By using a silhouette loss, we are able to accurately reconstruct body shapes.