

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

X-Avatar: Expressive Human Avatars

Anonymous CVPR submission

Paper ID 1050

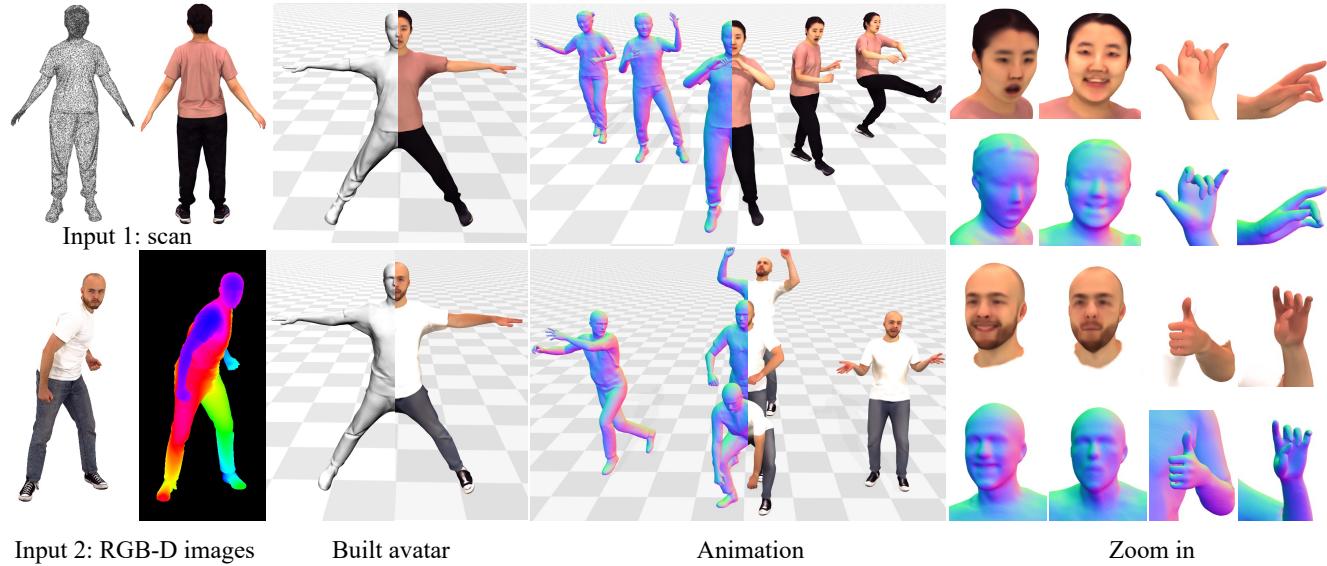


Figure 1. We propose **X-Avatar**, an animatable implicit human avatar model capable of capturing human body pose, hand pose, facial expressions, and appearance. X-Avatar can be created from input 3D scans (top row) or RGB-D images (bottom row) and displays high-quality geometry as well as appearance under animation. X-Avatar captures facial expressions and hand gestures (right), making it the first implicit human avatar model to capture the richness of the human state in a unified model.

Abstract

We present **X-Avatar**, a novel avatar model that captures the full expressiveness of digital humans to bring about life-like experiences in telepresence, AR/VR and beyond. Our method models bodies, hands, facial expressions and appearance in a holistic fashion and can be learned from either full 3D scans or RGB-D data. To achieve this, we propose a part-aware learned forward skinning module that can be driven by the parameter space of SMPL-X, allowing for expressive animation of X-Avatars. To efficiently learn the neural shape and deformation fields, we propose novel part-aware sampling and initialization strategies. This leads to higher fidelity results, especially for smaller body parts while maintaining efficient training despite increased number of articulated bones. To capture the appearance of the avatar with high-frequency details, we extend the geometry and deformation fields with a texture network that is conditioned on pose, facial expression, geometry and the normals of the deformed surface. We show

experimentally that our method outperforms strong baselines in both data domains both quantitatively and qualitatively on the animation task. To facilitate future research on expressive avatars we contribute a new dataset, called **X-Humans**, containing 234 sequences of high-quality textured scans from 20 participants, totalling 35,000 data frames.

1. Introduction

A significant part of human communication is non-verbal in which body pose, appearance, facial expressions, and hand gestures play an important role. Hence, it is clear that the quest towards immersive, life-like remote telepresence and other experiences in AR/VR, will require methods to capture the richness of human expressiveness in its entirety. Yet, it is not clear how to achieve this. Non-verbal communication involves an intricate interplay of several articulated body parts at different scales, which makes it difficult to capture and model algorithmically.

108 Parametric body models such as the SMPL family [32,
109 44, 46] have been instrumental in advancing the state-of-
110 the-art in modelling of digital humans in computer vision
111 and graphics. However, they rely on mesh-based represen-
112 tations and are limited to fixed topologies and in resolution
113 of the 3D mesh. These models are focused on minimally
114 clothed bodies and do not model garments or hair. Hence,
115 it is difficult to capture the full appearance of humans.
116

117 Neural implicit representations hold the potential to
118 overcome these limitations. Chen et al. [12] introduced a
119 method to articulate human avatars that are represented by
120 continuous implicit functions combined with learned for-
121 ward skinning. This approach has been shown to generalize
122 to arbitrary poses. While SNARF [12] only models the ma-
123 jor body bones, other works have focused on creating im-
124 plicit models of the face [21, 62], the hands [16], or how
125 to model humans that appear in garments [19] and how to
126 additionally capture appearance [47, 52]. Although neural
127 implicit avatars hold great promise, to date no model exists
128 that holistically captures the body and all the parts that are
129 important for human expressiveness jointly.
130

131 In this work, we introduce X-Avatar, an animatable, im-
132 plicit human avatar model that captures the shape, appear-
133 ance and deformations of complete humans and their hand
134 poses, facial expressions, and clothing. To this end we adopt
135 the full-body pose space of SMPL-X [44]. This causes two
136 key challenges for learning X-Avatars from data: (i) the
137 significantly increased number of involved articulated body
138 parts (9 used by [12] vs. 45 when including hands and face)
139 and (ii) the different scales at which they appear in obser-
140 vations. The hands and the face are much smaller in size
141 compared to the torso, arms and legs, yet they exhibit simi-
142 larly or even more complex articulations.
143

144 X-Avatar consists of a shape network that models the ge-
145 ometry in canonical space and a deformation network to es-
146 tablish correspondences between canonical and deformed
147 space via learned linear blend skinning (LBS). The par-
148 ameters of the shape and deformation fields must be learned
149 only from posed observations. SNARF [12] solves this via
150 iterative correspondence search. This optimization prob-
151 lem is initialized by transforming a large number of can-
152 didate points via the bone transformations. Directly adopt-
153 ing SNARF and initializing root-finding with only the body
154 bones leads to poor results for the hands and face. Hence,
155 to account for the articulation of these smaller body parts,
156 their bone transformations must also be considered. How-
157 ever, correspondence search scales poorly with the num-
158 ber of bones, so naïvely adding them makes training slow.
159 Therefore, we introduce a *part-aware initialization* strategy
160 which is almost 3 times faster than the naïve version while
161 outperforming it quantitatively. Furthermore, to counteract
the imbalance in scale between the body, hands, and face,
we propose a *part-aware sampling* strategy, which increases

162 the sampling rate for smaller body parts. This significantly
163 improves the fidelity of the final result. To model the ap-
164 pearance of X-Avatars, we extend the shape and deforma-
165 tion fields with an additional appearance network, condi-
166 tioned on pose, facial expression, geometry and the normals
167 in deformed space. All three neural fields are trained jointly.
168

169 X-Avatar can learn personalized avatars for multiple peo-
170 ple and from multiple input modalities. To demonstrate this,
171 we perform several experiments. First, we compare our
172 method to its most related work (SCANimate [47], SNARF
173 [12]) on the GRAB dataset [9, 51] on the animation task of
174 minimally clothed humans. Second, we contribute a novel
175 dataset consisting of 234 sequences of 20 clothed partici-
176 pants recorded in a high-quality volumetric capture stage
177 [15]. The dataset consists of subjects that perform diverse
178 body and hand poses (*e.g.*, counting, pointing, dancing) and
179 facial expressions (*e.g.*, laughing, screaming, frowning). On
180 this dataset we show that X-Avatar can learn from 3D scans
181 and (synthesized) RGB-D data. Our experiments show that
182 X-Avatar outperforms strong baselines both in quantitative
183 and qualitative measures in terms of animation quality. In
184 summary, we contribute:
185

- X-Avatar, the first expressive implicit human avatar
model that captures body pose, hand pose, facial ex-
pressions and appearance in a holistic fashion.
• Part-aware initialization and sampling strategies,
which together improve the quality of the results and
keep training efficient.
• X-Humans, a new dataset consisting of 234 sequences,
of high-quality textured scans showing 20 participants
with varied body and hand movements, and facial ex-
pressions, totalling 35,000 frames.

186 Data, models and SMPL[-X] registrations, will be released
187 for scientific purposes.
188

2. Related Work

189 **Explicit Human Models** Explicit models use a triangu-
190 lated 3D mesh to represent the underlying shape and are
191 controlled by a lower-dimensional set of parameters. Some
192 models focus on capturing a specific part of the human, *e.g.*,
193 the body [3, 32, 41], the hands [46], or the face [6, 31], while
194 others treat the human more holistically like we do in this
195 work, *e.g.* [25, 42, 44, 55, 60]. Explicit models are popu-
196 lar because the 3D mesh neatly fits into existing computer
197 graphics pipelines and because the low-dimensional param-
198 eter space lends itself well for learning. Only naturally have
199 such models thus been applied to tasks such as RGB-based
200 pose estimation [14, 20, 26–29, 48–50, 59], RGB-D fitting
201 [7, 13, 58], fitting to body-worn sensor data [23, 53, 57], or
202 3D hand pose estimation [8, 22] with a resounding success.
203 Because the SMPL family does not natively model cloth-
204 ing, researchers have investigated ways to extend it, *e.g.*
205 via fixed additive 3D offsets [1, 2], also dubbed SMPL+D,
206

216 pose-dependent 3D offsets [34], by modelling 3D garments
 217 and draping them over the SMPL mesh [5, 17] or via lo-
 218 cal small surface patches [33]. Explicit models have seen
 219 a trend towards unification to model human expressiveness,
 220 e.g. SMPL-X [44] and Adam [25]. X-Avatar shares this
 221 goal, but for implicit models.
 222

223 **Implicit Human Models** Explicit body models are lim-
 224 ited by their fixed mesh topology and resolution, and thus
 225 the expressive power required to model clothing and appear-
 226 ance necessitates extending these models beyond their orig-
 227 inal design. In contrast, using implicit functions to represent
 228 3D geometry grants more flexibility. With implicit mod-
 229 els, the shape is defined by neural fields, typically par-
 230 ameterized by MLPs that predict signed distance fields [43],
 231 density [39], or occupancy [36] given a point in space. To
 232 extend this idea to articulated shapes like the human body,
 233 NASA [18] used per body-part occupancy networks [36].
 234 This per-part formulation creates artifacts, especially for
 235 unseen poses, which works such as [12, 37, 38] improve.
 236 SNARF [12] does so via a forward warping field which
 237 is compatible with the SMPL [32] skeleton, learns pose-
 238 independent skinning and generalizes well to unseen poses
 239 and people in clothing. Other works [47, 52] model ap-
 240 pearance and are learned from scans. [54] creates avatars
 241 from video by relying on normals extracted from a 3D body
 242 model fitted to images and [19] does so from RGB-D video.
 243

244 Moving beyond bodies, other work has investigated im-
 245 plicit models for faces [21, 45, 56, 62] and hands [16]. Yet,
 246 an implicit model that incorporates body, hands, face, and
 247 clothing in a single model is missing. X-Avatar fills this
 248 gap. We do so by adopting neural forward skinning [12]
 249 driven by SMPL-X [44]. This seemingly simple change ne-
 250 cessitates non-trivial improvements to the correspondence
 251 search as otherwise the iterative root finding is too slow and
 252 leads to poor results which we show empirically. We pro-
 253 pose to do so by introducing part-aware initialization and
 254 sampling strategies, which are incorporated into a single
 255 model. Similar to [47], we obtain color with an MLP that is
 256 fed with canonical points and conditioned on the predicted
 257 geometry. Thanks to the part-aware sampling strategy, our
 258 method produces higher quality results than [47] for the
 259 hands and faces. Furthermore, in contrast to [12, 47, 52],
 260 X-Avatars can be fit to 3D scans *and* RGB-D videos.
 261

262 **Human Datasets** Publicly available datasets that show
 263 the full range of human expressiveness and contain clothed
 264 and textured ground-truth are rare. GRAB [51], a subset of
 265 AMASS [35], contains minimally clothed SMPL-X regis-
 266 trations. BUFF [61] and CAPE [34] do not model appear-
 267 ance and facial expressions. The CMU Panoptic Studio [24]
 268 dataset was used to fit Adam [25] which does model hands
 269 and faces, but is neither textured nor clothed. Also, [24]

270 does not contain scans. To study X-Avatars on real clothed
 271 humans, we thus contribute our own dataset, X-Humans
 272 which contains 35,000 frames of high-quality, texturized
 273 scans of real clothed humans with corresponding SMPL-[
 274 X] registrations.
 275

276 3. Method

277 We introduce X-Avatar, a method for the modeling of
 278 implicit human avatars with full body control including
 279 body movements, hand gestures, and facial expressions. For
 280 an overview, please refer to Fig. 2 and Fig. 3. Our model can
 281 be learned from two types of inputs, *i.e.*, 3D posed scans
 282 and RGB-D images. We first recap the SMPL-X full body
 283 model. Then we describe the X-Avatar formulation, train-
 284 ing scheme, and our part-aware initialization and sampling
 285 strategies. For simplicity, we discuss the scan-based version
 286 without loss of generality and list the differences to depth-
 287 based acquisition in the Supp. Mat.
 288

289 3.1. Recap: SMPL-X Unified Human Body Model

290 Our goal is to create fully controllable human avatars.
 291 We use the parameter space of SMPL-X [44], which it-
 292 self extends SMPL to include fully articulated hands and
 293 an expressive face. SMPL-X is defined by a function
 294 $M(\theta, \beta, \psi) : \mathbb{R}^{|\theta|} \times \mathbb{R}^{|\beta|} \times \mathbb{R}^{|\psi|} \rightarrow \mathbb{R}^{3N}$, parameterized
 295 by the shape β , whole body pose θ and facial expressions
 296 ψ . The pose can be further divided into the global pose
 297 θ_g , head pose θ_f , articulated hand poses θ_h , and remain-
 298 ing body poses θ_b . Here, $|\theta_g| = 3$, $|\theta_b| = 63$, $|\theta_h| = 90$,
 299 $|\theta_f| = 9$, $|\beta| = 10$, $|\psi| = 10$, $N = 10,475$.
 300

301 3.2. Implicit Neural Avatar Representation

302 To deal with the varying topology of clothed humans
 303 and to achieve higher geometric resolution and increased
 304 fidelity of overall appearance, X-Avatar proposes a human
 305 model defined by articulated neural implicit surfaces. We
 306 define three neural fields: one to model the geometry via an
 307 implicit occupancy network, one to model deformation via
 308 learned forward linear blend skinning (LBS) with continu-
 309 ous skinning weights, and one to model appearance as an
 310 RGB color value.
 311

312 **Geometry** We model the geometry of the human avatar
 313 in the canonical space with an MLP that predicts the occu-
 314 pancy value f_{occ} for any 3D point \mathbf{x}_c in this space. To cap-
 315 ture local non-rigid deformations such as facial or garment
 316 wrinkles, we condition the geometry network on the body
 317 pose θ_b and facial expression coefficients ψ . We found em-
 318 pirically that high-frequency details are preserved better if
 319 positional encodings [40] are applied to the input. Hence,
 320 the shape model f_{occ} is denoted by:
 321

$$f_{\text{occ}} : \mathbb{R}^3 \times \mathbb{R}^{|\theta_b|} \times \mathbb{R}^{|\psi|} \rightarrow [0, 1]. \quad (1)$$

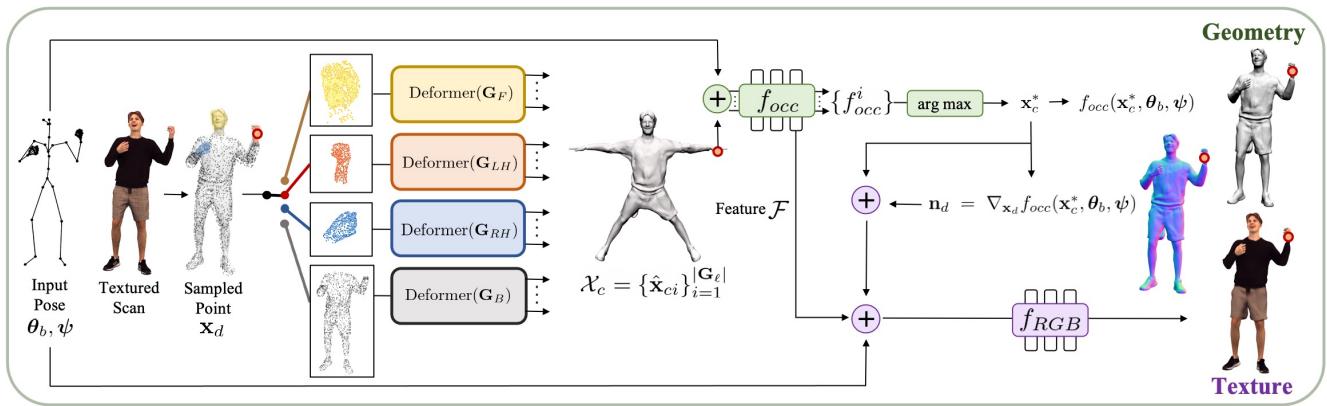
324
325
326
327
328
329
330
331
332
333
334
335
336
337378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398

Figure 2. **Method Overview.** Given a posed scan with an SMPL-X registration, we first adaptively sample points \mathbf{x}_d in deformed space per body part ℓ (face F , left hand LH , right hand RH , body B). A part-specific deformers network finds the corresponding candidate points $\hat{\mathbf{x}}_{ci}$ (for $1 \leq i \leq |\mathbf{G}_\ell|$) in canonical space via iterative root finding. The deformers share the parameters of the skinning network, but each deformer is initialized with only the bone transformations \mathbf{G}_ℓ (cf. Fig. 3). The final shape is obtained via an occupancy network f_{occ} . We further model appearance via a texture network that takes as input the body pose θ_b , facial expression ψ , the last layer \mathcal{F} of f_{occ} , the canonical point \mathbf{x}_c^* and the normals \mathbf{n}_d in deformed space. The normals correspond to the gradient $\nabla_{\mathbf{x}_d} f_{occ}(\mathbf{x}_c^*, \theta_b, \psi)$.

344

345
346
347
348
349
350
351
352

$\text{minimize } \mathbf{d}_w(\mathbf{x}_c, \mathbf{B}) - \mathbf{x}_d$ $\text{s.t. } \mathbf{d}_w(\mathbf{x}_c, \mathbf{B}) = \sum_{i=1}^{ \mathbf{G} } w_i(\mathbf{x}_c) \mathbf{B}_i \mathbf{x}_c$ $\mathbf{x}_c - f_w \rightarrow \mathbf{w}(\mathbf{x}_c)$ $\mathbf{x}_{ci}^0 = \mathbf{B}_i^{-1} \cdot \mathbf{x}_d \quad \mathbf{J}_{ci}^0 = \frac{\partial \mathbf{d}_w(\mathbf{x}, \mathbf{B}_i)}{\partial \mathbf{x}} \Big _{\mathbf{x}=\mathbf{x}_{ci}^0}, i \in \mathbf{G}_\ell$	Deformer(\mathbf{G}_ℓ) Linear Blend Skinning Skinning weights field Initialization of root finding
---	--

353
354
355
356
357

Figure 3. **Part-specific Deformer.** Each deformer shown in Fig. 2 is initialized with the bone transformations belonging to a specific part \mathbf{G}_ℓ , $\ell \in \{F, LH, RH, B\}$, but shares the parameters of f_w .

358
359
360
361

The canonical shape is defined as the 0.5 level set of f_{occ} :

$$\mathcal{S} = \{ \mathbf{x}_c \mid f_{occ}(\mathbf{x}_c, \theta_b, \psi) = 0.5 \}. \quad (2)$$

362
363
364
365

Deformation To model skeletal deformation, we follow previous work [12, 19, 30, 62] and represent the skinning weight field in the canonical space by an MLP:

$$f_w : \mathbb{R}^3 \rightarrow \mathbb{R}^{n_b} \times \mathbb{R}^{n_h} \times \mathbb{R}^{n_f}, \quad (3)$$

366
367
368
369
370
371
372
373
374

where n_b, n_h, n_f denotes the number of body, finger, and face bones respectively. Similar to [12], we assume a set of bones \mathbf{G} and require the weights $\mathbf{w} \in \mathbb{R}^{|\mathbf{G}|}$ to fulfill $w_i \geq 0$ and $\sum_i w_i = 1$. With the learned deformation field \mathbf{w} and given bone transformations $\mathbf{B} = \{\mathbf{B}_1, \dots, \mathbf{B}_{|\mathbf{G}|}\}$, for each point \mathbf{x}_c in the canonical space, its deformed counterpart is then uniquely determined:

$$\mathbf{x}_d = \mathbf{d}_w(\mathbf{x}_c, \mathbf{B}) = \sum_{i=1}^{|\mathbf{G}|} w_i(\mathbf{x}_c) \mathbf{B}_i \mathbf{x}_c. \quad (4)$$

Note that the canonical shape is a-priori unknown and learned during training. Since the relationship between deformed and canonical points is only implicitly defined, we follow [12] and employ correspondence search. We use Broyden’s method [10] to find canonical correspondences \mathbf{x}_c for each deformed query point \mathbf{x}_d iteratively as the roots of $\mathbf{d}_w(\mathbf{x}_c, \mathbf{B}) - \mathbf{x}_d = 0$. In cases of self-contact, multiple valid solutions exist. Therefore the optimization is initialized multiple times by transforming deformed points \mathbf{x}_d rigidly to the canonical space with each bone transformation. Finally, the set of valid correspondences \mathcal{X}_c is determined via analysis of the local convergence.

Part-Aware Initialization At the core of our method lies the problem of jointly learning the non-linear deformations introduced by body poses *and* dexterous hand articulation *and* facial expressions. The above method to attain multiple correspondences scales poorly with the number of bones. Therefore, naively adding finger and face bones of SMPL-X to the initialization procedure, causes prohibitively slow training. Yet our ablations show that these are required for good animation quality (*cf.* Tab. 1). Hence, we propose a part-aware initialization strategy, in which we first separate all SMPL-X bones \mathbf{G} into four groups \mathbf{G}_B , \mathbf{G}_{LH} , \mathbf{G}_{RH} , \mathbf{G}_F . For a given deformed point with part label ℓ , we then initialize the states $\{\mathbf{x}_{ci}^0\}$ and Jacobian matrices $\{\mathbf{J}_{ci}^0\}$ as:

$$\mathbf{x}_{ci}^0 = \mathbf{B}_i^{-1} \cdot \mathbf{x}_d, \quad \mathbf{J}_{ci}^0 = \frac{\partial \mathbf{d}_w(\mathbf{x}, \mathbf{B}_i)}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}_{ci}^0}, i \in \mathbf{G}_\ell. \quad (5)$$

We explain how we obtain the label ℓ for each point further below. The final occupancy prediction is determined

399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

432 via the maximum over all valid candidates $\mathcal{X}_c = \{\hat{\mathbf{x}}_{ci}\}_{i=1}^{|\mathbf{G}_\ell|}$:

$$433 \quad o(\mathbf{x}_d, \boldsymbol{\theta}_b, \boldsymbol{\psi}) = \max_{\hat{\mathbf{x}}_c \in \mathcal{X}_c} \{f_{\text{occ}}(\hat{\mathbf{x}}_c, \boldsymbol{\theta}_b, \boldsymbol{\psi})\}. \quad (6)$$

436 The correspondence in canonical space is given by:

$$437 \quad \mathbf{x}_c^* = \arg \max_{\hat{\mathbf{x}}_c \in \mathcal{X}_c} \{f_{\text{occ}}(\hat{\mathbf{x}}_c, \boldsymbol{\theta}_b, \boldsymbol{\psi})\}. \quad (7)$$

441 This part-aware initialization is based on the observation
442 that a point close to a certain body part is likely to be mostly
443 affected by the bones in that part. This scheme effectively
444 creates four deformor networks, as shown in Fig. 2. How-
445 ever, note that all deformers share the same skinning weight
446 network f_w as highlighted in Fig. 3. The only difference
447 between them is how the iterative root finding is initialized.
448

449 **Part-Aware Sampling** Because hands and faces are com-
450 paratively small, while still exhibiting complex deforma-
451 tions, we found that a uniform sampling strategy for points
452 \mathbf{x}_d leads to poor results (cf. Tab. 1, Fig. 4). Hence, we fur-
453 ther propose a part-aware sampling strategy, to over-sample
454 points per area for small body parts. Assuming part labels
455 $\mathcal{P} = \{F, LH, RH, B\}$, for each point \mathbf{p}_i in the 3D scan,
456 we first find the closest SMPL-X vertex \mathbf{v}_i and store its pre-
457 computed body part label $k_i \in \mathcal{P}$. Then, for each part $\ell \in \mathcal{P}$
458 we extract all points $\{\mathbf{p}_i \mid k_i = \ell\}$ and re-sample the result-
459 ing mesh with a sampling rate specific to part ℓ to obtain N_ℓ
460 many deformed points $\{\mathbf{x}_{di}\}_{i=1}^{N_\ell}$ for training.
461

462 **LBS regularization** To further account for the lower res-
463 olution and smaller scale of face and hands, we regularize
464 the LBS weights of these parts to be close to the weights
465 given by SMPL-X. A similar strategy has also been used
466 by [62]. Our ablations show that this greatly increases the
467 quality of the results (cf. Sec. 4.2).

468 **Texture** Similar to [47, 52] we introduce a third neu-
469 ral texture field to predict RGB values in canonical space.
470 Its output is the color value $c(\mathbf{x}_c, \mathbf{n}_d, \mathcal{F}, \boldsymbol{\theta}_b, \boldsymbol{\psi})$. This is,
471 in addition to pose and facial expression, the color de-
472 pends on the last layer \mathcal{F} of the geometry network and the
473 normals \mathbf{n}_d in deformed space. This conditions the
474 color prediction on the deformed geometry and local high-
475 frequency details, which has been shown to be helpful
476 [11, 62]. Following [62], the normals are obtained via
477 $\mathbf{n}_d = \nabla_{\mathbf{x}_d} f_{\text{occ}}(\mathbf{x}_c^*, \boldsymbol{\theta}_b, \boldsymbol{\psi})$. Therefore, the texture model
478 f_{RGB} is formulated as:
479

$$480 \quad f_{\text{RGB}} : \mathbb{R}^3 \times \mathbb{R}^3 \times \mathbb{R}^{512} \times \mathbb{R}^{|\boldsymbol{\theta}_b|} \times \mathbb{R}^{|\boldsymbol{\psi}|} \rightarrow \mathbb{R}^3. \quad (8)$$

484 We apply positional encoding to all inputs to obtain bet-
485 ter high-frequency details following best practices [39].

3.3. Training Process

Objective Function For each 3D scan, we minimize the following objective:

$$486 \quad \mathcal{L} = \mathcal{L}_{\text{occ}} + \mathcal{L}_{\text{RGB}} + \mathcal{L}_{\text{reg}}. \quad (9)$$

\mathcal{L}_{occ} supervises the geometry and consists of two losses: the binary cross entropy loss \mathcal{L}_{BCE} between the predicted occupancy $o(\mathbf{x}_d, \boldsymbol{\theta}_b, \boldsymbol{\psi})$ and the ground-truth value $o^{GT}(\mathbf{x}_d)$, and an L2 loss \mathcal{L}_n on the normals:

$$497 \quad \begin{aligned} \mathcal{L}_{\text{occ}} &= \lambda_{\text{BCE}} \mathcal{L}_{\text{BCE}} + \lambda_n \mathcal{L}_n \\ &= \lambda_{\text{BCE}} \sum_{\mathbf{x}_d \in \mathcal{P}_{\text{off}}} CE(o(\mathbf{x}_d, \boldsymbol{\theta}_b, \boldsymbol{\psi}), o^{GT}(\mathbf{x}_d)) \\ &\quad + \lambda_n \sum_{\mathbf{x}_d \in \mathcal{P}_{\text{on}}} \|\mathbf{n}_d - \mathbf{n}^{GT}(\mathbf{x}_d)\|_2, \end{aligned} \quad (10)$$

where \mathcal{P}_{on} , \mathcal{P}_{off} separately denote points on the scan surface and points within a thin shell surrounding the surface [19]. \mathcal{L}_{RGB} supervises the point color:

$$507 \quad \mathcal{L}_{\text{RGB}} = \lambda_{\text{RGB}} \sum_{\mathbf{x}_d \in \mathcal{P}_{\text{on}}} \|c(\mathbf{x}_c, \mathbf{n}_d, \mathcal{F}, \boldsymbol{\theta}_b, \boldsymbol{\psi}) - c^{GT}(\mathbf{x}_d)\|_1. \quad (11)$$

Finally, \mathcal{L}_{reg} represents the regularization term, consisting of the bone occupancy loss $\mathcal{L}_{\text{bone}}$, joint LBS weights loss $\mathcal{L}_{\text{joint}}$ and surface LBS weights loss $\mathcal{L}_{\text{surf}}$:

$$514 \quad \begin{aligned} \mathcal{L}_{\text{reg}} &= \lambda_{\text{bone}} \mathcal{L}_{\text{bone}} + \lambda_{\text{joint}} \mathcal{L}_{\text{joint}} + \lambda_{\text{surf}} \mathcal{L}_{\text{surf}} \\ &= \lambda_{\text{bone}} \sum_{\mathbf{x}_c \in \mathcal{P}_{\text{bone}}^c} CE(f_{\text{occ}}(\mathbf{x}_c, \boldsymbol{\theta}_b, \boldsymbol{\psi}), 1) \\ &\quad + \lambda_{\text{joint}} \sum_{\mathbf{x}_c \in \mathcal{P}_{\text{joint}}^c} \sum_{i \in \mathcal{N}(i)} (w_i(\mathbf{x}_c) - 0.5)^2 \\ &\quad + \lambda_{\text{surf}} \sum_{\mathbf{x}_c \in \mathcal{P}_{\text{surf}}^c} \sum_{i \in \mathbf{G} \setminus \mathbf{G}_B} (w_i(\mathbf{x}_c) - w_i^{GT}(\mathbf{x}_c))^2, \end{aligned} \quad (12)$$

where $\mathcal{N}(i)$ are the neighboring bones of joint i and w_i^{GT} are the skinning weights taken from SMPL-X. \mathcal{L}_{reg} makes use of the supervision from registered SMPL-X meshes. For more details on the registration, please refer to the Supp. Mat. $\mathcal{P}_{\text{bone}}^c$, $\mathcal{P}_{\text{joint}}^c$, $\mathcal{P}_{\text{surf}}^c$ refer to points sampled on the SMPL-X bones, the SMPL-X joints and from the SMPL-X mesh surface respectively. The first two terms follow the definition of [12]. We add the last term to regularize the LBS weights for fingers and face which have low resolution and are more difficult to learn.

4. Experiments

We first introduce the datasets and metrics that we use for our experiments in Sec. 4.1. Sec. 4.2 ablates all important design choices. In Sec. 4.3 we briefly describe the state-of-the-art methods to which we compare our method. Finally

540 we show and discuss the results in Sec. 4.4-4.6. We focus on
 541 the challenging animation task, hence all the comparisons
 542 are conducted on entirely unseen poses. For completeness,
 543 we also report reconstruction results in the Supp. Mat.
 544

545 4.1. Datasets

546 **GRAB [51]** We use the GRAB subset of AMASS [35]
 547 for training and evaluate our model on SMPL-X meshes of
 548 minimally clothed humans. GRAB contains a diverse set of
 549 hand poses and facial expressions with several subjects. We
 550 pick the subject with the most pose variation and randomly
 551 select 9 sequences for training and one for validation. This
 552 results in 9,756 frames for training and 235 test frames.
 553

554 **X-Humans (Scans)** Currently, there exists no publicly
 555 available dataset containing textured 3D clothed scans of
 556 humans with a large variation of body poses, hand gestures
 557 and facial expressions. Therefore, we captured our own
 558 dataset, for which we leveraged a high-quality, multi-view
 559 volumetric capture stage [15]. We call the resulting dataset
 560 X-Humans. It consists of 20 subjects (11 males, 9 females)
 561 with various clothing types and hair style. The collection of
 562 this dataset has been approved by an internal ethics com-
 563 mittee. For each subject, we split the motion sequences
 564 into a training and test set. In total, there are 28,294 poses
 565 for training and 6,369 test poses. X-Humans also contains
 566 ground-truth SMPL-X parameters, obtained via a custom
 567 SMPL-X registration pipeline specifically designed to deal
 568 with low-resolution body parts. More details on the regis-
 569 tration process and contents of X-Humans are in Supp. Mat.
 570

571 **X-Humans (RGB-D)** We take the textured and posed
 572 scans from X-Humans and render them to obtain corre-
 573 sponding synthetic RGB-D images. For every time step,
 574 we render exactly one RGB-D image from a virtual camera,
 575 while the camera gradually rotates around the participant
 576 during the duration of the sequence. This is, the RGB-D
 577 version of X-Humans contains the same amount of frames
 578 as the scan version in both the training and test set.
 579

580 **Metrics** We evaluate the geometric accuracy via volu-
 581 metric IoU, Chamfer distance (CD) (mm) and normal con-
 582 sistency (NC) metrics, following the practice in PINA [19].
 583 Because these metrics are dominated by large surface areas,
 584 we always report the metrics for the entire body (*All*) and
 585 the hands separately (*Hands*).
 586

587 4.2. Ablation Study

588 **Part-Aware Initialization** The part-aware initialization
 589 for correspondence search is critical to accelerate training
 590 and to find good correspondences in small body parts. To
 591 verify this, we compare with two variations adapted from
 592

593 SNARF [12]. First, (A1) initiates the optimization states
 594 only via the body’s bone transformations, while (A2) initial-
 595 izes using all bones (body, hands, face). **Results:** A1 suffers
 596 from strong artifacts for hands and the jaw (*cf.* Fig. 4 and
 597 Tab. 1, A1). The final model is 3 times faster than A2 (0.7 it-
 598 erations per second vs. 0.25), yet it still retains high fidelity
 599 and even outperforms A2 by a small margin (*cf.* Tab. 1,
 600 A2 and the Supp. Mat. for qualitative results). Thus we
 601 conclude that part-based initialization of the deformer is an
 602 efficient way to find accurate correspondences.
 603

604 **Part-Aware Sampling** To verify the importance of part-
 605 aware sampling, we compare our model to a uniform sam-
 606 pling baseline (A3). **Results:** This component has two ef-
 607 fects: a) it strongly improves the hand shape (*cf.* second
 608 column of Fig. 4 and Tab. 1, A3) and b) it improves texture
 609 details in the eye and mouth region (*cf.* Fig. 5).
 610

611 **LBS Weights Regularization for Hands and Face** The
 612 first column in Fig. 4 shows that without regularizing the
 613 learned LBS weights with the SMPL-X weights, the learned
 614 hand shape is poor. This is further substantiated by a 75%
 615 increase in Chamfer distance for the hand region, compared
 616 to our final method (*cf.* Tab. 1, A4).
 617

618 4.3. Baselines

619 **Scan-based methods** We compare our 3D scan-based
 620 method variant on both GRAB and X-Humans to SM-
 621 PLX+D, SCANimate and SNARF baselines. We adapt SM-
 622 PLX+D from SMPL+D introduced in [4]. This baseline
 623 uses an explicit body model, SMPL-X, and models clothing
 624 with additive vertex offsets. To compare with SCANimate
 625 and SNARF, we use publicly available code. For details on
 626 the baselines, we refer to the Supp. Mat.
 627

628 **RGB-D Video-based methods** We compare our RGB-D
 629 method variant on the X-Humans (RGB-D) dataset to PINA
 630 [19], a SMPL-based implicit human avatar method learned
 631 from RGB-D inputs. We assume that the ground truth pose
 632 and shape are known. For a fair comparison we do not opti-
 633 mize these parameters in PINA.
 634

635 4.4. Results on GRAB Dataset

636 Tab. 2 summarizes results on the GRAB dataset. Over-
 637 all, our method beats all baselines, especially for the hands,
 638 where the margin is large. Fig. 6 visually shows that the
 639 quality of the hands and face learned by our method is much
 640 higher: SCANimate learns a mean hand and SNARF gen-
 641 eralizes badly to the unseen poses. Since GRAB meshes are
 642 minimally clothed, we omit SMPLX+D from comparison.
 643

ID	Method	CD↓		CD-MAX ↓		NC ↑		IoU ↑		702
		All	Hands	All	Hands	All	Hands	All	Hands	
A1	Ours (init w body bones)	5.42	5.05	57.54	25.10	0.940	0.824	0.964	0.812	704
A2	Ours (init w all bones)	4.55	4.35	44.86	20.71	0.945	0.845	0.974	0.811	705
A3	Ours (w/o part-aware sampling)	4.68	4.81	47.51	20.88	0.947	0.840	0.972	0.810	706
A4	Ours (w/o LBS reg.)	4.98	7.27	57.11	43.38	0.940	0.797	0.968	0.768	707
A	Ours (complete)	4.46	4.15	44.36	20.61	0.948	0.853	0.973	0.829	709

Table 1. **Ablation experiments for our major design choices.** We compute the metrics on the entire body (*All*) and separately on the hands (*Hands*) to better highlight the differences for the hands. All results are computed on a subset of X-Humans (Scans). Our final model (A) only marginally outperforms A2, but is roughly 3 times faster to train. For qualitative comparisons please refer to Fig. 4 and Fig. 5.

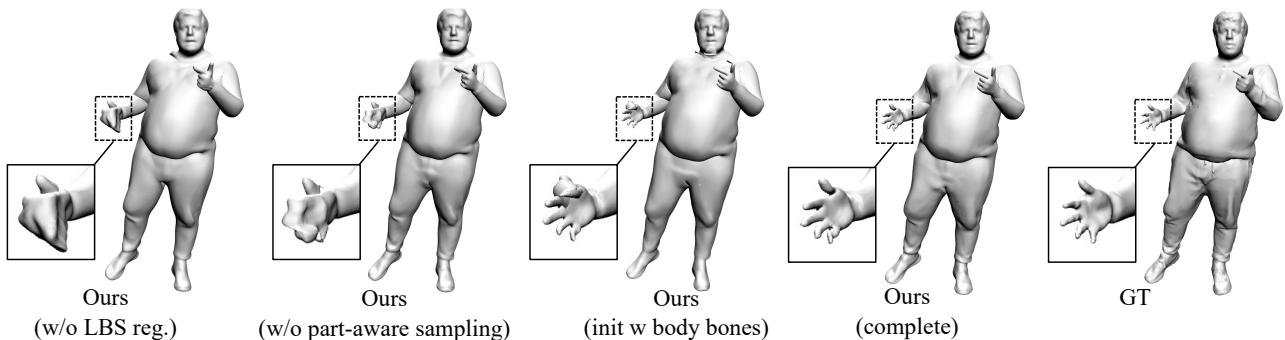


Figure 4. **Effect of our design decisions on the resulting geometry.** Notice how all baselines struggle to recover accurate hand geometry.

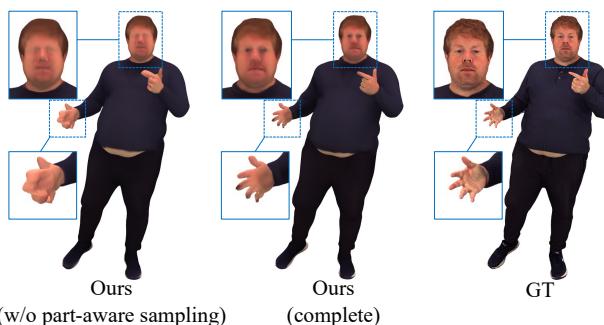


Figure 5. **Effect of our part-aware sampling strategy** on the hand geometry and texture prediction of the face.

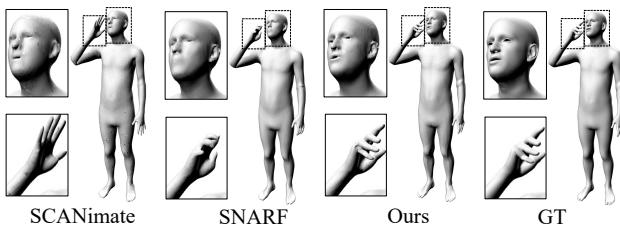


Figure 6. **Qualitative results on GRAB dataset.** Our method recovers hand articulation and facial expression most accurately.

4.5. Results on X-Humans (Scans)

Tab. 3 shows that our method also outperforms the baselines on X-Humans. Fig. 7 qualitatively shows differences.

Method	CD↓		CD-MAX ↓		NC ↑		IoU ↑		728
	All	Hands	All	Hands	All	Hands	All	Hands	
SCANimate [47]	2.32	6.25	54.86	55.18	0.970	0.804	0.938	0.656	729
SNARF [12]	1.11	3.75	29.98	29.72	0.980	0.851	0.975	0.724	730
Ours	0.88	0.75	16.76	4.87	0.984	0.962	0.994	0.901	731

Table 2. **Quantitative results on GRAB dataset.** Our method outperforms all baselines, especially for the hand part (*cf.* Fig. 6).

SMPLX+D, limited by its fixed topology and low resolution, cannot model details like hair and wrinkles in clothing. SCANimate and SNARF are SMPL-driven, so they either learn a static or incomplete hand. Our method balances the different body parts so that hands are well-structured, but also the details on the face and body are maintained. Fig. 1 and Fig. 9 show more animation results.

Method	CD↓		CD-MAX ↓		NC ↑		IoU ↑		742
	All	Hands	All	Hands	All	Hands	All	Hands	
SMPLX+D	5.75	5.19	48.41	23.48	0.921	0.790	0.957	0.774	743
SCANimate [47]	6.54	9.78	59.71	48.32	0.925	0.726	0.919	0.557	744
SNARF [12]	5.05	7.23	55.06	37.15	0.934	0.788	0.937	0.608	745
Ours	4.43	5.14	47.56	22.15	0.939	0.793	0.965	0.776	746

Table 3. **Quantitative results on X-Humans (Scans).** We beat all baselines both for the entire body (*All*) and hands only (*Hands*).

4.6. Results on X-Humans (RGB-D)

Tab. 4 shows our model's performance compared to PINA [19]. Our method outperforms PINA on all metrics. Fig. 8 further qualitatively shows that without utilizing the hand and face information in the modelling process, the face and hands produced by PINA are not consistent with the in-

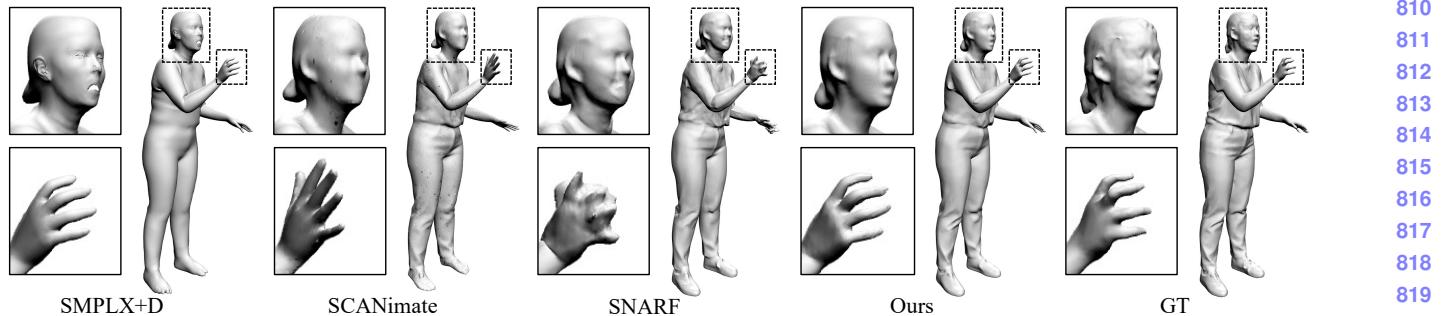


Figure 7. **Qualitative animation comparison on X-Humans (Scans).** SMPLX+D fails to model face and garment details. SCANimate and SNARF generate poor hands (static or incomplete). Our method produces the most plausible face and hands, and keeps the clothing details comparable to strong baselines.

Method	CD \downarrow		CD-MAX \downarrow		NC \uparrow		IoU \uparrow	
	All	Hands	All	Hands	All	Hands	All	Hands
PINA [19]	5.41	9.51	66.05	48.07	0.928	0.771	0.910	0.566
Ours	5.33	5.27	51.73	22.86	0.936	0.797	0.947	0.768

Table 4. **Quantitative results on X-Humans (RGB-D).** Our method outperforms PINA in all metrics. Improvements are more pronounced for hands (*cf.* Fig. 8 for visual comparison).

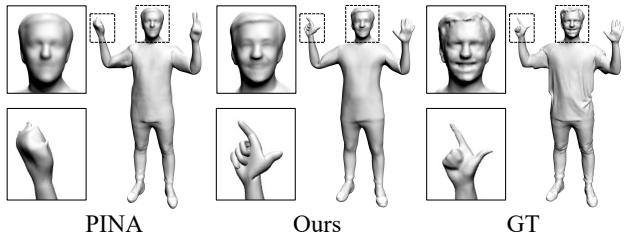


Figure 8. **X-Avatars created from RGB-D input compared to PINA.** Notice how we obtain better hand and face geometry.

put pose. Our model generates a) more realistic faces as the shape network is conditioned on facial expression and b) better hand poses because we initialize the root finding with hand bone transformations.

5. Conclusion

Limitations X-Avatar struggles to model loose clothing that is far away from the body (*e.g.* skirts). Furthermore, generalization capability beyond a single person is still limited, *i.e.* we train one model for each subject.

Conclusion We propose X-Avatar, the first expressive implicit human avatar model that captures body pose, hand pose, facial expressions and appearance in a holistic fashion. We have demonstrated our method’s expressive power, the benefit of our proposed part-aware initialization and sampling strategy, and the capability of creating it from multiple input modalities with the aid of our newly introduced X-Humans dataset. We believe that our method along with X-Humans will promote further scientific research in creating expressive digital avatars.

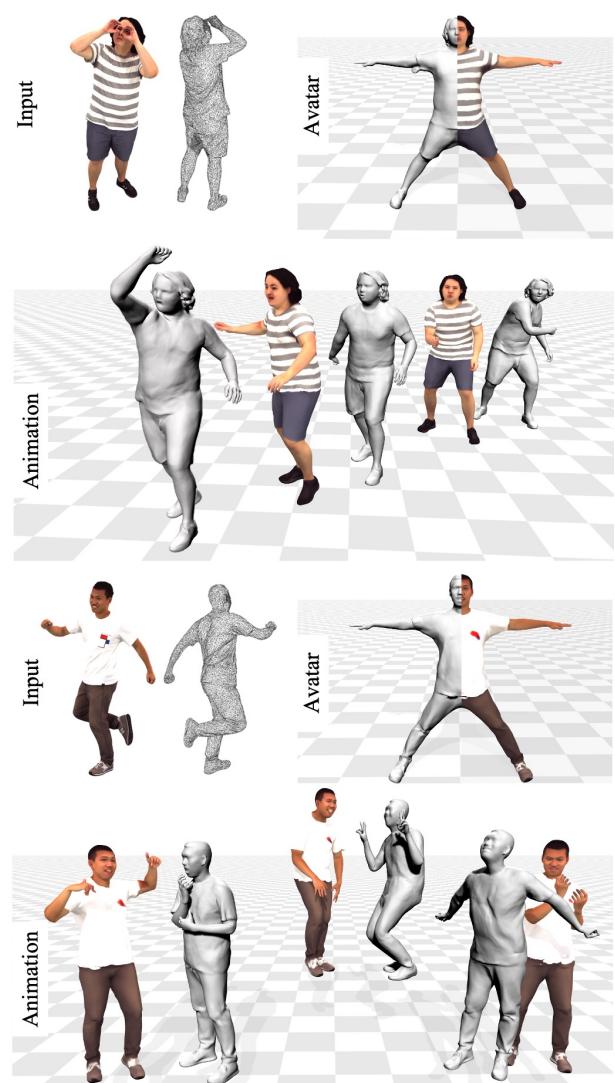


Figure 9. **Animation demonstration on X-Humans (Scans).** Our method can handle relatively complex clothing patterns, hair styles, and varied facial expressions, hand, and body poses.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

References

- [1] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1175–1186, 2019. [2](#)
- [2] Thiemo Alldieck, Marcus Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. Video based reconstruction of 3d people models. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8387–8397, 2018. [2](#)
- [3] Dragomir Anguelov, Praveen Srinivasan, Daphne Koller, Sebastian Thrun, Jim Rodgers, and James Davis. Scape: Shape completion and animation of people. *ACM Trans. Graph.*, 24(3):408–416, jul 2005. [2](#)
- [4] Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Combining implicit function learning and parametric models for 3d human reconstruction. In *European Conference on Computer Vision (ECCV)*. Springer, aug 2020. [6](#)
- [5] Bharat Lal Bhatnagar, Garvita Tiwari, Christian Theobalt, and Gerard Pons-Moll. Multi-garment net: Learning to dress 3d people from images. In *IEEE International Conference on Computer Vision (ICCV)*. IEEE, Oct 2019. [3](#)
- [6] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH ’99*, page 187–194, USA, 1999. ACM Press/Addison-Wesley Publishing Co. [2](#)
- [7] Federica Bogo, Michael J. Black, Matthew Loper, and Javier Romero. Detailed full-body reconstructions of moving people from monocular rgb-d sequences. *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 2300–2308, 2015. [2](#)
- [8] Adnane Boukhayma, Rodrigo de Bem, and Philip HS Torr. 3d hand shape and pose from images in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10843–10852, 2019. [2](#)
- [9] Samarth Brahmbhatt, Cusuh Ham, Charles C. Kemp, and James Hays. ContactDB: Analyzing and predicting grasp contact via thermal imaging. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. [2](#)
- [10] Charles G Broyden. A class of methods for solving nonlinear simultaneous equations. *Mathematics of computation*, 19(92):577–593, 1965. [4](#)
- [11] Xu Chen, Tianjian Jiang, Jie Song, Jinlong Yang, Michael J. Black, Andreas Geiger, and Otmar Hilliges. gdna: Towards generative detailed neural avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20427–20437, June 2022. [5](#)
- [12] Xu Chen, Yufeng Zheng, Michael J Black, Otmar Hilliges, and Andreas Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. In *International Conference on Computer Vision (ICCV)*, 2021. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#)
- [13] Yin Chen, Z. Cheng, Chao Lai, Ralph Robert Martin, and Gang Dang. Realtime reconstruction of an animating human body from a single depth camera. *IEEE Transactions on Visualization and Computer Graphics*, 22:2000–2011, 2016. [2](#)
- [14] Vasileios Choutas, Georgios Pavlakos, Timo Bolkart, Dimitrios Tzionas, and Michael J. Black. Monocular expressive body regression through body-driven attention. In *European Conference on Computer Vision (ECCV)*, 2020. [2](#)
- [15] Alvaro Collet, Ming Chuang, Pat Sweeney, Don Gillett, Dennis Evseev, David Calabrese, Hugues Hoppe, Adam Kirk, and Steve Sullivan. High-quality streamable free-viewpoint video. *ACM Trans. Graph.*, 34(4), jul 2015. [2](#), [6](#)
- [16] Enric Corona, Tomas Hodan, Minh Vo, Francesc Moreno-Noguer, Chris Sweeney, Richard Newcombe, and Lingni Ma. Lisa: Learning implicit shape and appearance of hands. In *CVPR*, 2022. [2](#), [3](#)
- [17] Enric Corona, Albert Pumarola, Guillem Alenyà, Gerard Pons-Moll, and Francesc Moreno-Noguer. Smplicit: Topology-aware generative model for clothed people. In *CVPR*, 2021. [3](#)
- [18] Boyang Deng, JP Lewis, Timothy Jeruzalski, Gerard Pons-Moll, Geoffrey Hinton, Mohammad Norouzi, and Andrea Tagliasacchi. Neural articulated shape approximation. In *The European Conference on Computer Vision (ECCV)*. Springer, August 2020. [3](#)
- [19] Zijian Dong, Chen Guo, Jie Song, Xu Chen, Andreas Geiger, and Otmar Hilliges. Pina: Learning a personalized implicit neural avatar from a single rgb-d video sequence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20470–20480, June 2022. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [20] Yao Feng, Vasileios Choutas, Timo Bolkart, Dimitrios Tzionas, and Michael J. Black. Collaborative regression of expressive bodies using moderation. In *International Conference on 3D Vision (3DV)*, 2021. [2](#)
- [21] Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8649–8658, June 2021. [2](#), [3](#)
- [22] Yana Hasson, Gülcin Varol, Dimitrios Tzionas, Igor Kalevatykh, Michael J. Black, Ivan Laptev, and Cordelia Schmid. Learning joint reconstruction of hands and manipulated objects. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11799–11808, 2019. [2](#)
- [23] Yinghao Huang, Manuel Kaufmann, Emre Aksan, Michael J. Black, Otmar Hilliges, and Gerard Pons-Moll. Deep inertial poser: Learning to reconstruct human pose from sparse inertial measurements in real time. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 37:185:1–185:15, Nov. 2018. [2](#)
- [24] Hanbyul Joo, Hao Liu, Lei Tan, Lin Gui, Bart Nabbe, Iain Matthews, Takeo Kanade, Shohei Nobuhara, and Yaser Sheikh. Panoptic studio: A massively multiview system for social motion capture. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 3334–3342, 2015. [3](#)
- [25] Hanbyul Joo, Tomas Simon, and Yaser Sheikh. Total capture: A 3d deformation model for tracking faces, hands, and bodies. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. [2](#), [3](#)

- 972 [26] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and
973 Jitendra Malik. End-to-end recovery of human shape and
974 pose. In *Computer Vision and Pattern Recognition (CVPR)*,
975 2018. 2
- 976 [27] Muhammed Kocabas, Chun-Hao P. Huang, Otmar Hilliges,
977 and Michael J. Black. PARE: Part attention regressor for
978 3D human body estimation. In *Proc. International Conference
979 on Computer Vision (ICCV)*, pages 11127–11137, Oct.
980 2021.
- 981 [28] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and
982 Kostas Daniilidis. Learning to reconstruct 3d human pose
983 and shape via model-fitting in the loop. In *ICCV*, 2019.
- 984 [29] Jiefeng Li, Chao Xu, Zhicun Chen, Siyuan Bian, Lixin Yang,
985 and Cewu Lu. Hybrik: A hybrid analytical-neural inverse
986 kinematics solution for 3d human pose and shape estimation.
987 In *Proceedings of the IEEE/CVF Conference on Computer
988 Vision and Pattern Recognition*, pages 3383–3393, 2021. 2
- 989 [30] Rui long Li, Julian Tanke, Minh Vo, Michael Zollhofer, Jur-
990 gen Gall, Angjoo Kanazawa, and Christoph Lassner. Tava:
Template-free animatable volumetric actors. 2022. 4
- 991 [31] Tianye Li, Timo Bolkart, Michael. J. Black, Hao Li, and
992 Javier Romero. Learning a model of facial shape and ex-
993 pression from 4D scans. *ACM Transactions on Graphics*,
(*Proc. SIGGRAPH Asia*), 36(6), 2017. 2
- 994 [32] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard
995 Pons-Moll, and Michael J. Black. SMPL: A skinned multi-
996 person linear model. *ACM Transactions on Graphics*, (*Proc.
997 SIGGRAPH Asia*), 34(6):248:1–248:16, Oct. 2015. 2, 3
- 998 [33] Qianli Ma, Shunsuke Saito, Jinlong Yang, Siyu Tang, and
999 Michael J. Black. SCALE: Modeling clothed humans with a
1000 surface codec of articulated local elements. In *Proceedings
1001 IEEE/CVF Conf. on Computer Vision and Pattern Recog-
1002 nition (CVPR)*, pages 16082–16093, 2021. 3
- 1003 [34] Qianli Ma, Jinlong Yang, Anurag Ranjan, Sergi Pujades,
1004 Gerard Pons-Moll, Siyu Tang, and Michael J. Black. Learn-
1005 ing to dress 3d people in generative clothing. In *Computer
1006 Vision and Pattern Recognition (CVPR)*, June 2020. 3
- 1007 [35] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Ger-
1008 ard Pons-Moll, and Michael J. Black. AMASS: Archive of
1009 motion capture as surface shapes. In *International Confer-
1010 ence on Computer Vision*, pages 5442–5451, Oct. 2019. 3,
1011 6
- 1012 [36] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Se-
1013 bastian Nowozin, and Andreas Geiger. Occupancy networks:
1014 Learning 3d reconstruction in function space. In *Proceedings
1015 of the IEEE/CVF Conference on Computer Vision and Pat-
1016 tern Recognition*, pages 4460–4470, 2019. 3
- 1017 [37] Marko Mihajlovic, Shunsuke Saito, Aayush Bansal, Michael
1018 Zollhoefer, and Siyu Tang. COAP: Compositional articu-
1019 lated occupancy of people. In *Proceedings IEEE Conf. on
1020 Computer Vision and Pattern Recognition (CVPR)*, June
1021 2022. 3
- 1022 [38] Marko Mihajlovic, Yan Zhang, Michael J Black, and Siyu
1023 Tang. LEAP: Learning articulated occupancy of people. In
1024 *Proceedings IEEE Conf. on Computer Vision and Pattern
1025 Recognition (CVPR)*, June 2021. 3
- 1026 [39] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik,
1027 Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf:
1028 Representing scenes as neural radiance fields for view syn-
1029 thesis. In *ECCV*, 2020. 3, 5
- 1030 [40] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and
1031 Andreas Geiger. Occupancy flow: 4d reconstruction by
1032 learning particle dynamics. In *Proceedings of the IEEE/CVF
1033 international conference on computer vision*, pages 5379–
1034 5389, 2019. 3
- 1035 [41] Ahmed A A Osman, Timo Bolkart, and Michael J. Black.
1036 STAR: A sparse trained articulated human body regressor.
1037 In *European Conference on Computer Vision (ECCV)*, pages
1038 598–613, 2020. 2
- 1039 [42] Ahmed A A Osman, Timo Bolkart, Dimitrios Tzionas, and
1040 Michael J. Black. SUPR: A sparse unified part-based human
1041 body model. In *European Conference on Computer Vision
1042 (ECCV)*, 2022. 2
- 1043 [43] Jeong Joon Park, Peter Florence, Julian Straub, Richard
1044 Newcombe, and Steven Lovegrove. DeepSDF: Learning con-
1045 tinuous signed distance functions for shape representation.
1046 In *Proceedings of the IEEE/CVF Conference on Computer
1047 Vision and Pattern Recognition*, pages 165–174, 2019. 3
- 1048 [44] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani,
1049 Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and
1050 Michael J. Black. Expressive body capture: 3d hands, face,
1051 and body from a single image. In *Proceedings IEEE Conf.
1052 on Computer Vision and Pattern Recognition (CVPR)*, 2019.
1053 2, 3
- 1054 [45] Eduard Ramon, Gil Triginer, Janna Escur, Albert Pumarola,
1055 Jaime Garcia, Xavier Giro-i Nieto, and Francesc Moreno-
1056 Noguer. H3d-net: Few-shot high-fidelity 3d head reconstruc-
1057 tion. In *Proceedings of the IEEE/CVF International Confer-
1058 ence on Computer Vision*, pages 5620–5629, 2021. 3
- 1059 [46] Javier Romero, Dimitrios Tzionas, and Michael J. Black.
1060 Embodied hands: Modeling and capturing hands and bodies
1061 together. *ACM Transactions on Graphics*, (*Proc. SIG-
1062 GRAPH Asia*), 36(6):245:1–245:17, Nov. 2017. 2
- 1063 [47] Shunsuke Saito, Jinlong Yang, Qianli Ma, and Michael J.
1064 Black. SCANimate: Weakly supervised learning of skinned
1065 clothed avatar networks. In *Proceedings IEEE/CVF Conf. on
1066 Computer Vision and Pattern Recognition (CVPR)*, June
1067 2021. 2, 3, 5, 7
- 1068 [48] Jie Song, Xu Chen, and Otmar Hilliges. Human body model
1069 fitting by learned gradient descent. 2020. 2
- 1070 [49] Yu Sun, Qian Bao, Wu Liu, Yili Fu, Black Michael J., and
1071 Tao Mei. Monocular, one-stage, regression of multiple 3d
1072 people. In *ICCV*, 2021.
- 1073 [50] Yu Sun, Wu Liu, Qian Bao, Yili Fu, Tao Mei, and Michael J
1074 Black. Putting people in their place: Monocular regression
1075 of 3d people in depth. In *CVPR*, 2022. 2
- 1076 [51] Omid Taheri, Nima Ghorbani, Michael J. Black, and Dim-
1077 itrios Tzionas. GRAB: A dataset of whole-body human
1078 grasping of objects. In *European Conference on Computer
1079 Vision (ECCV)*, 2020. 2, 3, 6
- 1080 [52] Garvita Tiwari, Nikolaos Sarafianos, Tony Tung, and Ger-
1081 ard Pons-Moll. Neural-gif: Neural generalized implicit func-
1082 tions for animating people in clothing. In *International Con-
1083 ference on Computer Vision (ICCV)*, October 2021. 2, 3, 5
- 1084 [53] Timo von Marcard, Roberto Henschel, Michael Black, Bodo
1085 Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d
1086 human pose in the wild using imus and a moving camera.
1087 In *European Conference on Computer Vision (ECCV)*, sep-
1088

- 1080 2018. 2 1134
1081 [54] Yuliang Xiu, Jinlong Yang, Dimitrios Tzionas, and Michael J. Black. ICON: Implicit Clothed humans Obtained from Normals. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13296–13306, June 2022. 3 1135
1082 [55] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6184–6193, 2020. 2 1136
1083 [56] T Yenamandra, A Tewari, F Bernard, HP Seidel, M Elgarib, D Cremers, and C Theobalt. i3dmm: Deep implicit 3d morphable model of human heads. In *Proceedings of the IEEE / CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021. 3 1137
1084 [57] Xinyu Yi, Yuxiao Zhou, Marc Habermann, Soshi Shimada, Vladislav Golyanik, Christian Theobalt, and Feng Xu. Physical inertial poser (pip): Physics-aware real-time human motion tracking from sparse inertial sensors. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022. 2 1138
1085 [58] Tao Yu, Zerong Zheng, Kaiwen Guo, Jianhui Zhao, Qionghai Dai, Hao Li, Gerard Pons-Moll, and Yebin Liu. Doublefusion: Real-time capture of human performances with inner body shapes from a single depth sensor. In *The IEEE International Conference on Computer Vision and Pattern Recognition(CVPR)*. IEEE, June 2018. 2 1139
1086 [59] Ye Yuan, Umar Iqbal, Pavlo Molchanov, Kris Kitani, and Jan Kautz. Glamr: Global occlusion-aware human mesh recovery with dynamic cameras. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2 1140
1087 [60] 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 *Computer Vision – ECCV 2020*, pages 465–481, 2020. 2 1141
1088 [61] Chao Zhang, Sergi Pujades, Michael J. Black, and Gerard Pons-Moll. Detailed, accurate, human shape estimation from clothed 3d scan sequences. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 3 1142
1089 [62] Yufeng Zheng, Victoria Fernández Abrevaya, Marcel C. Bühler, Xu Chen, Michael J. Black, and Otmar Hilliges. I M Avatar: Implicit morphable head avatars from videos. In *Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3, 4, 5 1143
1090 [63] 1144
1091 [64] 1145
1092 [65] 1146
1093 [66] 1147
1094 [67] 1148
1095 [68] 1149
1096 [69] 1150
1097 [70] 1151
1098 [71] 1152
1099 [72] 1153
1100 [73] 1154
1101 [74] 1155
1102 [75] 1156
1103 [76] 1157
1104 [77] 1158
1105 [78] 1159
1106 [79] 1160
1107 [80] 1161
1108 [81] 1162
1109 [82] 1163
1110 [83] 1164
1111 [84] 1165
1112 [85] 1166
1113 [86] 1167
1114 [87] 1168
1115 [88] 1169
1116 [89] 1170
1117 [90] 1171
1118 [91] 1172
1119 [92] 1173
1120 [93] 1174
1121 [94] 1175
1122 [95] 1176
1123 [96] 1177
1124 [97] 1178
1125 [98] 1179
1126 [99] 1180
1127 [100] 1181
1128 [101] 1182
1129 [102] 1183
1130 [103] 1184
1131 [104] 1185
1132 [105] 1186
1133 [106] 1187