

Supplementary Material

001 0.1. Addition related work

002 The literature [22, 24] on recovering 3D human representations from RGB images is vast. Techniques fall broadly
003 into two categories. Parametric methods [3, 16] characterize the human body in terms of a parametric model. Model
004 parameters defining body pose and shape are then estimated from images via direct optimization [20, 25, 26] or regression
005 with deep networks [5, 6, 11, 12]. Non-parametric methods directly regress a 3D body representation from images
006 using convolutional neural networks [13], transformers [15], intermediate representations [18] or implicit functions
007 [7, 21]. On the other hand, the development of human shape estimation erupt 3D digital human research applications
008 [4, 9] and the parametric models [16] from images and videos have attracted increasing attention. Optimization-
009 based methods [20] detect 2D features corresponding to the whole body and fit the SMPL-X model. However, they suffer
010 from slow speed and are ultimately limited by the quality of the 2D keypoint detectors. Hence, learning-based
011 models are proposed. Due to the highly complex multi-stage pipelines, the reconstructed results inevitably generated
012 unnatural articulation mesh and implausible 3D wrist rotations. [14] proposes the first ViT-based backbone [8]
013 to relieve the issues in previous approaches. This provides a promising and concise way to leverage the scaling-up
014 model for two-stage body measurements. However, there is a scarcity of benchmark datasets for comparing body
015 measurements, and few researchers are exploring the integration of additional data toward generalizable and accurate
016 body measurement results.

031 0.2. Detailed about datasets

032 We describe the datasets mentioned in the main paper. Note that all these are public academic datasets, each holding a
033 license. We follow the common practice of using them in our non-commercial research and refer readers to their policies
034 to ensure personal information protection.

037 **Lidar dataset:** We utilize a high-performance commercial Mojave Sensor, available at an affordable price, for
038 scanning subjects aged between 18 and 24 years old. This choice allows us to conveniently conduct tests on LiDAR
039 data. The output of this sensor is a distance and amplitude file. The distance image can be converted to a point cloud
040 (x,y,z values for each distance point) by leveraging the sensor lens parameters that are stored on the device. Further-
041 more, each subject was measured by a skilled anthropologist, providing chest, waist, hip, wrist, and shoulder width
042 as the GT measurement values for our experiments. RGB images were captured in a well-lit, indoor setup, with sub-

049 jects standing in A-pose. Note that we never require subjects to wear tight-fitting clothing. Capture distance varied
050 from 2.5-3.5 meters.

052 **AGORA** [19] is a synthetic dataset, rendered with high-quality human scans and realistic 3D scenes. It consists
053 of 4240 textured human scans with diverse poses and appearances, each fitted with accurate SMPL-X annotations.
054 There are 14K training images and 3K test images, and 173K instances.

058 **3DPW** [23] is the first in-the-wild dataset with a considerable amount of data, captured with a moving phone
059 camera and IMU sensors. It features accurate SMPL annotations and 60 video sequences captured in diverse environ-
060 ments. We follow the official definition of train, val, and test splits.

064 **Human3.6M** [10] is a studio-based 3D motion capture dataset including 3.6M human poses and corresponding im-
065 ages captured by a high-speed motion capture system. In this paper, we use the annotation generated by NeuralAnnot,
066 which fits the SMPL-X to the GT 2D joints and includes a total of 312.2K annotated data.

070 **MPI-INF-3DHP** [17] is captured with a multi-camera markerless motion capture system in constrained indoor and
071 complex outdoor scenes. It records 8 actors performing 8 activities from 14 camera views. We use the annotations
072 generated by NeuralAnnot, which fits the SMPL-X to the GT 2D joints and includes a total of 939,847 annotated data.

076 **MPII** [1] is a widely used in-the-wild dataset that offers a diverse collection of approximately 25K images. Each
077 image within the dataset contains one or more instances, resulting in a total of over 40K annotated people instances.
078 Among the 40K samples, 28K samples are used for training, while the remaining samples are reserved for testing. We
079 use the annotations generated by NeuralAnnot, which fits the SMPL-X to the GT 2D joints and includes a total of
080 ~28.9K annotated data.

085 0.3. Motivation about focusing network

086 Traditional fine-tuning methods require modifying the top layer of the network to adapt differences in label spaces and
087 losses, which can disrupt the pretrained features and diminish the network’s reusability. In contrast, our focusing net-
088 work employs bypass network to extract various guidance features. This modification preserves the pretrained features
089 for consistent performance and facilitates efficient model sharing. On the other hand, the traditional two-stage meth-
090 ods of reconstructing before measuring have limitations in generalizing to different scene categories. Additionally, the
091 reconstruction quality lacks reliability under extreme viewing angles, making accurate measurement across various
092
093
094
095
096
097

098 scenarios challenging. We are the first to solve the above
099 problem by utilizing large-scale models for estimating anthropometric
100 measurements. Moreover, we introduce a bypass network to fine-tune the output of the large model, an
101 innovation not present in previous methods.
102

103 0.4. Details process about using Mojave Sensor

104 The output of this sensor for a subject is a distance and amplitude file. The distance image can be converted to a point
105 cloud (x,y,z values for each distance point) by leveraging the sensor lens parameters that are stored on the device. By
106 using the focal length and optical center of the lens a coefficient for x, y, and z can be calculated at each pixel location.
107 These arrays of coefficients are stored on the device referred to as the pixel rays. After querying this information
108 from the sensor, each value in the distance image can be multiplied by the three coefficients to obtain its x,y, and
109 z coordinates in 3D space.
110
111
112
113
114

115 0.5. Training details

116 Following previous work, we use the following typical datasets for training focusing network, i.e., Human3.6M [10],
117 MPI-INF-3DHP [17], MPII [2] and so on. Additionally, we provided 30 sets of samples, each containing frontal and
118 profile images along with additional information. During the training of the bypass network, we utilized both our
119 collected dataset and publicly available datasets, as described in Sec.5.1 of the main paper. Additionally, our
120 Fashion-body dataset can be continually enhanced with diverse scenarios and individuals to meet the needs of a
121 broader range of human reconstruction and measurement tasks. We considered the fairness of these baselines for
122 comparison. The datasets used for HMR-BMViT and 4D-BMViT training are consistent with the bypass network,
123 but NeuralAnthro could not support such a large amount of data for training.
124
125
126
127
128
129
130
131

132 0.6. Fixed hyper-parameter for loss function

133 In this part, we evaluate the effect of the loss design parameters on the measurement performance. We re-trained
134 the proposed method on our Fashion-body dataset using fixed values for the parameters in the loss function
135 introduced in main paper. We chose four different fixed values: 0.1,0.2,0.4,0.8. The validation error is presented in Tab. 1.
136 The result shows that the larger value of the α parameter multiplies the measurement parts the smaller error of
137 measurement results.
138
139
140
141

142 0.7. Discussion

143 The generalization ability of the proposed approach, such as its robustness to background changes and variations in
144 lighting conditions. Actually, our Fashion-body dataset
145

α	Chest	Waist	Hip	Wrist	Shoulder width
0.1	6.02	9.57	10.99	1.07	4.34
0.2	5.87	9.04	10.73	1.00	4.24
0.4	4.31	5.21	7.21	0.56	2.73
0.8	3.07	3.91	6.95	0.41	2.03

Table 1. The MAE error of five body part in hyper-parameter tuning.

146 consists of examples collected from diverse angles, lighting conditions, and backgrounds. Our experiments present
147 the body measurement outcomes and MAE of our approach across multiple datasets, including our Fashion-body
148 dataset, which validate the robustness of our approach.
149
150

151 **Applications.** We propose tailored measurement pipelines and scanner selections for diverse anthropometric
152 applications: medicine, fashion, fitness, and entertainment. Our model can estimate key parameters of the human
153 body from a simple RGB image captured by a camera or smartphone, making it applicable across various domains
154 such as health, fashion, and entertainment. In the health domain, our model enables users to monitor changes
155 in body measurements over fixed intervals (e.g., weekly or monthly) by capturing images, aiding in the
156 formulation of fitness and diet plans. In the fashion domain, our model’s outputs can assist users in
157 selecting the most suitable clothing sizes in online shopping scenarios. In the entertainment domain, our
158 model serves as a valuable tool for virtual character creation and clothing rendering, ensuring the physical
159 realism and coherence of virtual scenes.
160
161
162
163
164
165
166

167 **Limitations and Future works.** Due to constraints related to the specialized equipment and cost required for
168 dataset creation, our model is trained only on single-view human reconstruction, disregarding the potential
169 benefits of information redundancy from multiple views. Utilizing multiple views to assist in human
170 reconstruction can mitigate the impact of certain viewpoints or inappropriate poses on the model. In the
171 future, we anticipate generating multi-view human reconstruction datasets. Additionally, considering the
172 inference speed of the model is crucial for practical scenarios. Therefore, while ensuring model performance,
173 we aim to enhance the inference speed.
174
175
176
177
178

179 References

- 180 [1] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark
181 and state of the art analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014. 1
182
183
184 [2] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark
185 and state of the art analysis. In *Proceedings of the IEEE Conference on computer Vision and Pattern Recognition*, pages
186 3686–3693, 2014. 2
187
188

- 189 [3] Dragomir Anguelov, Praveen Srinivasan, Daphne Koller, Se- 247
190 bastian Thrun, Jim Rodgers, and James Davis. Scape: shape 248
191 completion and animation of people. In *ACM SIGGRAPH* 249
192 *2005 Papers*, pages 408–416. 2005. 1
- 193 [4] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter 250
194 Gehler, Javier Romero, and Michael J Black. Keep it smpl: 251
195 Automatic estimation of 3d human pose and shape from a 252
196 single image. In *Computer Vision–ECCV 2016: 14th Euro- 253
197 pean Conference, Amsterdam, The Netherlands, October 11- 254
198 14, 2016, Proceedings, Part V 14*, pages 561–578. Springer, 255
199 2016. 1
- 200 [5] Vasileios Choutas, Georgios Pavlakos, Timo Bolkart, Dim- 256
201 itrios Tzionas, and Michael J Black. Monocular expressive 257
202 body regression through body-driven attention. In *Computer 258
203 Vision–ECCV 2020: 16th European Conference, Glasgow, 259
204 UK, August 23–28, 2020, Proceedings, Part X 16*, pages 20– 260
205 40. Springer, 2020. 1
- 206 [6] Vasileios Choutas, Lea Müller, Chun-Hao P Huang, Siyu 261
207 Tang, Dimitrios Tzionas, and Michael J Black. Accurate 3d 262
208 body shape regression using metric and semantic attributes. 263
209 In *Proceedings of the IEEE/CVF Conference on Computer 264
210 Vision and Pattern Recognition*, pages 2718–2728, 2022. 1
- 211 [7] Eric Corona, Albert Pumarola, Guillem Alenya, Gerard 265
212 Pons-Moll, and Francesc Moreno-Noguer. Smplicit: 266
213 Topology-aware generative model for clothed people. In 267
214 *Proceedings of the IEEE/CVF conference on computer vi- 268
215 sion and pattern recognition*, pages 11875–11885, 2021. 1
- 216 [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, 269
217 Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, 270
218 Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl- 271
219 vain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is 272
220 worth 16x16 words: Transformers for image recognition at 273
221 scale. *CoRR*, abs/2010.11929, 2020. 1
- 222 [9] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang 274
223 Cai, Lei Yang, and Ziwei Liu. Avatarclip: Zero-shot text- 275
224 driven generation and animation of 3d avatars. *arXiv preprint 276
225 arXiv:2205.08535*, 2022. 1
- 226 [10] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian 277
227 Sminchisescu. Human3. 6m: Large scale datasets and pre- 278
228 dictive methods for 3d human sensing in natural environ- 279
229 ments. *IEEE transactions on pattern analysis and machine 280
230 intelligence*, 36(7):1325–1339, 2013. 1, 2
- 231 [11] Angjoo Kanazawa, Michael J Black, David W Jacobs, and 281
232 Jitendra Malik. End-to-end recovery of human shape and 282
233 pose. In *Proceedings of the IEEE conference on computer 283
234 vision and pattern recognition*, pages 7122–7131, 2018. 1
- 235 [12] Muhammed Kocabas, Chun-Hao P Huang, Otmar Hilliges, 284
236 and Michael J Black. Pare: Part attention regressor for 3d 285
237 human body estimation. In *Proceedings of the IEEE/CVF 286
238 International Conference on Computer Vision*, pages 11127– 287
239 11137, 2021. 1
- 240 [13] Nikos Kolotouros, Georgios Pavlakos, and Kostas Dani- 288
241 ilidis. Convolutional mesh regression for single-image hu- 289
242 man shape reconstruction. In *Proceedings of the IEEE/CVF 290
243 Conference on Computer Vision and Pattern Recognition*, 291
244 pages 4501–4510, 2019. 1
- 245 [14] Jing Lin, Ailing Zeng, Haoqian Wang, Lei Zhang, and Yu Li. 292
246 One-stage 3d whole-body mesh recovery with component 293
aware transformer. In *Proceedings of the IEEE/CVF Con- 294
ference on Computer Vision and Pattern Recognition*, pages 295
21159–21168, 2023. 1
- [15] Kevin Lin, Lijuan Wang, and Zicheng Liu. End-to-end hu- 296
man pose and mesh reconstruction with transformers. In *Pro- 297
ceedings of the IEEE/CVF conference on computer vision 298
and pattern recognition*, pages 1954–1963, 2021. 1
- [16] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard 299
Pons-Moll, and Michael J Black. Smpl: A skinned multi- 300
person linear model. In *Seminal Graphics Papers: Pushing 301
the Boundaries, Volume 2*, pages 851–866. 2023. 1
- [17] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, 302
Oleksandr Sotnychenko, Weipeng Xu, and Christian 303
Theobalt. Monocular 3d human pose estimation in the wild 304
using improved cnn supervision. In *2017 international con- 305
ference on 3D vision (3DV)*, pages 506–516. IEEE, 2017. 1, 306
2
- [18] Gyeongsik Moon and Kyoung Mu Lee. I2l-meshnet: Image- 307
to-lixel prediction network for accurate 3d human pose and 308
mesh estimation from a single rgb image. In *Computer 309
Vision–ECCV 2020: 16th European Conference, Glasgow, 310
UK, August 23–28, 2020, Proceedings, Part VII 16*, pages 311
752–768. Springer, 2020. 1
- [19] Priyanka Patel, Chun-Hao P Huang, Joachim Tesch, David T 312
Hoffmann, Shashank Tripathi, and Michael J Black. 313
AGORA: Avatars in geography optimized for regression 314
analysis. In *Proceedings of the IEEE/CVF Conference on 315
Computer Vision and Pattern Recognition*, pages 13468– 316
13478, 2021. 1
- [20] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, 317
Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and 318
Michael J Black. Expressive body capture: 3d hands, 319
face, and body from a single image. In *Proceedings of 320
the IEEE/CVF conference on computer vision and pattern 321
recognition*, pages 10975–10985, 2019. 1
- [21] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul 322
Joo. Pifuhd: Multi-level pixel-aligned implicit function for 323
high-resolution 3d human digitization. In *Proceedings of 324
the IEEE/CVF Conference on Computer Vision and Pattern 325
Recognition*, pages 84–93, 2020. 1
- [22] Yating Tian, Hongwen Zhang, Yebin Liu, and Limin Wang. 326
Recovering 3d human mesh from monocular images: A sur- 327
vey. *IEEE transactions on pattern analysis and machine in- 328
telligence*, 2023. 1
- [23] Timo Von Marcard, Roberto Henschel, Michael J Black, 329
Bodo Rosenhahn, and Gerard Pons-Moll. Recovering ac- 330
curate 3d human pose in the wild using imus and a moving 331
camera. In *Proceedings of the European conference on com- 332
puter vision (ECCV)*, pages 601–617, 2018. 1
- [24] Jinbao Wang, Shujie Tan, Xiantong Zhen, Shuo Xu, Feng 333
Zheng, Zhenyu He, and Ling Shao. Deep 3d human pose 334
estimation: A review. *Computer Vision and Image Under- 335
standing*, 210:103225, 2021. 1
- [25] Donglai Xiang, Hanbyul Joo, and Yaser Sheikh. Monocu- 336
lar total capture: Posing face, body, and hands in the wild. 337
In *Proceedings of the IEEE/CVF conference on computer vi- 338
sion and pattern recognition*, pages 10965–10974, 2019. 1

- 304 [26] Andrei Zanfir, Eduard Gabriel Bazavan, Mihai Zanfir,
305 William T Freeman, Rahul Sukthankar, and Cristian Smin-
306 chisescu. Neural descent for visual 3d human pose and
307 shape. In *Proceedings of the IEEE/CVF Conference on Com-
308 puter Vision and Pattern Recognition*, pages 14484–14493,
309 2021. 1