# HuMMan: Multi-Modal 4D Human Dataset for Versatile Sensing and Modeling

Zhongang Cai[*,1,2,3], Daxuan Ren[*,2], Ailing Zeng[*,4], Zhengyu Lin[*,3], Tao Yu[*,5], Wenjia Wang[*,3],
Xiangyu Fan[3], Yang Gao[3], Yifan Yu[3], Liang Pan[2], Fangzhou Hong[2], Mingyuan Zhang[2],
Chen Change Loy[2], Lei Yang[†,1,3], Ziwei Liu[†,2]

[1]Shanghai AI Laboratory, [2]S-Lab, Nanyang Technological University, [3]SenseTime Research
[4]The Chinese University of Hong Kong, [5]Tsinghua University
[*]co-first authors, [†]co-corresponding authors

arXiv:2204.13686v1 [cs.CV] 28 Apr 2022



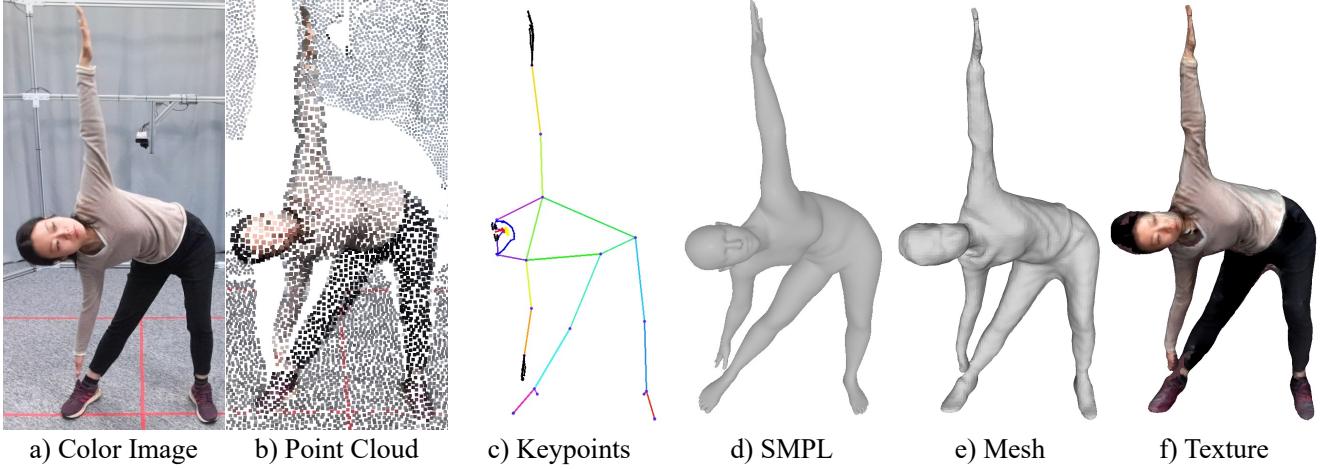a) Color Image  b) Point Cloud  c) Keypoints  d) SMPL  e) Mesh  f) Texture

Figure 1. HuMMan features multiple modalities of data format and annotations. We demonstrate a) color image, b) point cloud, c) keypoints, d) SMPL parameters and e) mesh geometry with f) texture. Each sequence is also annotated with an action label from 500 actions. Each subject has two additional high-resolution scans of naturally and minimally clothed body.

## Abstract

*4D human sensing and modeling are fundamental tasks in vision and graphics with numerous applications. With the advances of new sensors and algorithms, there is an increasing demand for more versatile datasets. In this work, we contribute **HuMMan**, a large-scale multi-modal 4D human dataset with 1000 human subjects, 400k sequences and 60M frames. HuMMan has several appealing properties: 1) multi-modal data and annotations including color images, point clouds, keypoints, SMPL parameters, and textured meshes; 2) popular mobile device is included in the sensor suite; 3) a set of 500 actions, designed to cover fundamental movements; 4) multiple tasks such as action recognition, pose estimation, parametric human recovery, and textured mesh reconstruction are supported and evaluated. Extensive experiments on HuMMan voice the need for further study on challenges such as fine-grained action recognition, dynamic human mesh reconstruction, point cloud-based parametric human recovery, and cross-device domain gaps.[1]*

## 1. Introduction

Sensing and modeling humans are longstanding problems for both computer vision and computer graphics research communities, which serve as the fundamental technology for a myriad of applications such as animation, gaming, augmented, and virtual reality. With the advent of deep learning, significant progress has been made alongside the introduction of large-scale datasets in human-centric sensing and modeling [32, 56, 63, 99, 109, 119]. In this work, we present **HuMMan**, a comprehensive human dataset consisting of 1000 human subjects, captured in total 400k se-

---

[1]https://caizhongang.github.io/projects/HuMMan

Table 1. Comparisons of HuMMan with published datasets. HuMMan has a competitive scale in terms of the number of subjects (#Subj), actions (#Act), sequences (#Seq) and frames (#Frame). Moreover, HuMMan features multiple modalities and supports multiple tasks. Video: sequential data, not limited to RGB sequences; Mobile: mobile device in the sensor suite; D/PC: depth image or point cloud, only genuine point cloud collected from depth sensors are considered; Act: action label; K2D: 2D keypoints; K3D: 3D keypoints; Param: statistical model (*e.g.* SMPL) parameters; Txtr: texture. -: not applicable or not reported.

| Dataset | #Subj | #Act | #Seq | #Frame | Video | Mobile | Modalities | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | | | | RGB | D/PC | Act | K2D | K3D | Param | Mesh | Txtr |
| UCF101 [91] | - | 101 | 13k | - | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| AVA [22] | - | 80 | 437 | - | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| FineGym [88] | - | 530 | 32k | - | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| HAA500 [15] | - | 500 | 10k | 591k | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| SYSU 3DHOI [30] | 40 | 12 | 480 | - | ✓ | - | ✓ | ✓ | ✓ | - | ✓ | - | - | - |
| NTU RGB+D [87] | 40 | 60 | 56k | - | ✓ | - | ✓ | ✓ | ✓ | - | ✓ | - | - | - |
| NTU RGB+D 120 [58] | 106 | 120 | 114k | - | ✓ | - | ✓ | ✓ | ✓ | - | ✓ | - | - | - |
| NTU RGB+D X [97] | 106 | 120 | 113k | - | ✓ | - | ✓ | ✓ | ✓ | - | ✓ | ✓ | - | - |
| MPII [3] | - | 410 | - | 24k | - | - | ✓ | - | ✓ | ✓ | - | - | - | - |
| COCO [56] | - | - | - | 104k | - | - | ✓ | - | - | ✓ | - | - | - | - |
| PoseTrack [2] | - | - | >1.35k | >46k | ✓ | - | ✓ | - | - | ✓ | - | - | - | - |
| Human3.6M [32] | 11 | 17 | 839 | 3.6M | ✓ | - | ✓ | ✓ | ✓ | ✓ | ✓ | - | - | - |
| CMU Panoptic [38] | 8 | 5 | 65 | 154M | ✓ | - | ✓ | ✓ | - | ✓ | ✓ | - | - | - |
| MPI-INF-3DHP [68] | 8 | 8 | 16 | 1.3M | ✓ | - | ✓ | - | - | ✓ | ✓ | - | - | - |
| 3DPW [99] | 7 | - | 60 | 51k | ✓ | ✓ | ✓ | - | - | - | - | ✓ | - | - |
| AMASS [65] | 344 | - | >11k | >16.88M | ✓ | - | - | - | - | - | ✓ | ✓ | - | - |
| AIST++ [52] | 30 | - | 1.40k | 10.1M | ✓ | - | ✓ | - | - | ✓ | ✓ | ✓ | - | - |
| CAPE [63] | 15 | - | >600 | >140k | ✓ | - | - | - | ✓ | - | ✓ | ✓ | ✓ | - |
| BUFF [113] | 6 | 3 | >30 | >13.6k | ✓ | - | ✓ | ✓ | ✓ | - | ✓ | ✓ | ✓ | ✓ |
| DFAUST [6] | 10 | >10 | >100 | >40k | ✓ | - | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| HUMBI [109] | 772 | - | - | ~26M | ✓ | - | ✓ | - | - | ✓ | ✓ | ✓ | ✓ | ✓ |
| ZJU LightStage [82] | 6 | 6 | 9 | >1k | ✓ | - | ✓ | - | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| THuman2.0 [107] | 200 | - | - | >500 | - | - | - | - | - | - | - | ✓ | ✓ | ✓ |
| **HuMMan (ours)** | 1000 | 500 | 400k | 60M | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

quences and 60M frames. More importantly, HuMMan features four main properties listed below.

**- Multiple Modalities**. HuMMan provides a basket of data formats and annotations in the hope to assist exploration in their potential complementary nature [29]. We build HuMMan with a set of 10 synchronized RGB-D cameras to capture both video and depth sequences. Our toolchain then post-process the raw data into sequences of colored point clouds, 2D/3D keypoints, statistical model (SMPL) parameters, and model-free textured mesh. Note that all data and annotations are temporally synchronized, while 3D data and annotations are spatially aligned. In addition, we provide a high-resolution scan for each of the subjects in a canonical pose.

**- Mobile Device**. With the development of 3D sensors, it is common to find depth cameras or low-power LiDARs on a mobile device in recent years. In view of the surprising gap between emerging real-life applications and the insufficiency of data collected with mobile devices, we add a mobile phone with built-in LiDAR in the data collection to facilitate the relevant research.

**- Action Set**. We design HuMMan to empower comprehensive studies on human actions. Instead of empirically selecting daily activities, we propose to take an anatomical point of view and systematically divide body movements by their driving muscles. Specifically, we design 500 movements by categorizing major muscle groups to achieve a more complete and fundamental representation of human actions.

**- Multiple Tasks**. To facilitate research on HuMMan, we provide a whole suite of baselines and benchmarks for action recognition, 2D and 3D pose estimation, 3D parametric human recovery, and textured mesh reconstruction. Popular methods are implemented and evaluated using standard metrics. Our experiments demonstrate that HuMMan would be useful for multiple fields of study, such as fine-grained action recognition, point cloud-based parametric human recovery, dynamic mesh sequence reconstruction, and transferring knowledge across devices.

In summary, HuMMan is a large-scale multi-modal dataset for 4D (spatio-temporal) human sensing and modeling, with four main features: **1**) multi-modal data and annotations; **2**) mobile device included in the sensor suite; **3**) action set with atomic motions; **4**) standard benchmarks for multiple vision tasks. We hope HuMMan would pave the way towards more comprehensive sensing and modeling of humans.

## 2. Related Works

**Action Recognition.** As an important step towards understanding human activities, action recognition is the task to categorize human motions into predefined classes. RGB videos [17, 18, 95, 96] with additional information such as optical flow and estimated poses and 3D skeletons typically obtained from RGB-D sequences [89, 90, 105, 111] are the common input to existing methods. Datasets for RGB video-based action recognition are often collected from the Internet. Some have a human-centric action design [15, 22, 42, 49, 88, 91] whereas others introduce interaction and diversity in the setup [11, 71, 117]. Recently, fine-grained action understanding [15, 22, 88] is drawing more research attention. However, these 2D datasets lack 3D annotations. As for RGB-D datasets, earlier works are small in scale [30, 53, 101]. As a remedy, the latest NTU RGB-D series [58, 87, 97] features 60-120 actions. However, the majority of the actions are focused on the upper body. We develop a larger and more complete action set in HuMMan.

**2D and 3D Keypoint Detection.** Estimation of a human pose is a vital task in computer vision, and a popular pose representation is human skeletal keypoints. The field is categorized by output format: 2D [12, 50, 74, 92] and 3D [66, 80, 110–112, 118] keypoint detection, or by the number of views: single-view [12, 66, 74, 80, 92, 111, 118] and multi-view pose estimation [31, 33, 83]. For 2D keypoint detection, single-frame datasets such as MPII [3] and COCO [56] provide diverse images with 2D keypoints annotations, whereas video datasets such as J-HMDB [35], Penn Action [114] and PoseTrack [2] provide sequences of 2D keypoints. However, they lack 3D ground truths. In contrast, 3D keypoint datasets are typically built indoor data to accommodate sophisticated equipment, such as Human3.6M [32], CMU Panoptic [38], MPI-INF-3DHP [68], TotalCapture [98], and AIST++ [52]. Compared to these datasets, HuMMan not only supports 2D and 3D keypoint detection but also textured mesh reconstruction assist in more holistic modeling of humans.

**3D Parametric Human Recovery.** Also known as human pose and shape estimation, 3D parametric human recovery leverages human parametric model representation (such as SMPL [61], SMPL-X [78], STAR [76] and GHUM [104]) that achieves sophisticated mesh reconstruction with a small amount of parameters. Existing methods take keypoints [5, 78, 115], images [21, 23, 47, 48, 51, 75, 79], videos [13, 40, 62, 69, 72, 93] as the input to obtain the parameters. Recently, point clouds have become more popular [4, 28, 36, 57, 102] for both parametric human and clothing recovery. Apart from those that provide keypoints, various datasets also provide ground-truth SMPL parameters. MoSh [60] is applied on Human3.6M [32] to generate SMPL annotations. CMU Panoptic [38] and HUMBI [109]

leverages keypoints from multiple camera views. 3DPW [99] combines a mobile phone and inertial measurement units (IMUs). Synthetic dataset such as AGORA [77] renders high-quality human scans in virtual environments and fits SMPL to the original mesh. Video games have also become an alternative source of data [9, 10]. In addition to SMPL parameters that do not model clothes or texture, HuMMan also provides textured meshes of clothed subjects.

**Textured Mesh Reconstruction.** To reconstruct the 3D surface, common methods include multi-view stereo [19], volumetric fusion [34, 73, 108], Poisson surface reconstruction [43, 45], and neural surface reconstruction [81, 86]. To reconstruct texture for the human body, popular approaches include texture mapping or montage [20], deep neural rendering [59], deferred neural rendering [94], and NeRF-like methods [70]. Unfortunately, existing datasets for textured human mesh reconstruction typically provide no sequential data [107, 119], which is valuable to the reconstruction of animatable avatars [84, 103]. Moreover, many have only a limited number of subjects [1, 6, 24–26, 63, 82, 113]. In contrast, HuMMan includes diverse subjects with high-resolution body scans and a large amount of dynamic 3D sequences.

## 3. Hardware Setup

We customize an octagonal prism-shaped multi-layer framework to accommodate calibrated and synchronized sensors. The system is 1.7 m in height and 3.4 m in side length of its octagonal cross-section as illustrated in Fig. 2.

### 3.1. Sensors

**RGB-D Sensors.** Azure Kinect is popular with both academia and the industry with a color resolution of 1920×1080, and a depth resolution of 640×576. We deploy ten Kinects to capture multi-view RGB-D sequences. The Kinects are strategically placed to ensure a uniform spacing, and a wide coverage such that any body part of the subject, even in most expressive poses, is visible to at least two sensors. We develop a program that interfaces with Kinect's SDK to obtain a data throughput of 74.4 MB per frame and 2.2 GB per second at 30 FPS before data compression.

**Mobile Device.** An iPhone 12 Pro Max is included in the sensor suite to allow for the study on a mobile device. Besides the regular color images of resolution 1920×1440, the built-in LiDAR produces depth maps of resolution 256×192. We develop an iOS app upon ARKit to retrieve the data.

**High-Resolution Scanner.** To supplement our sequential data with high-quality body shape information, a professional handheld 3D scanner, Artec Eva, is used to produce

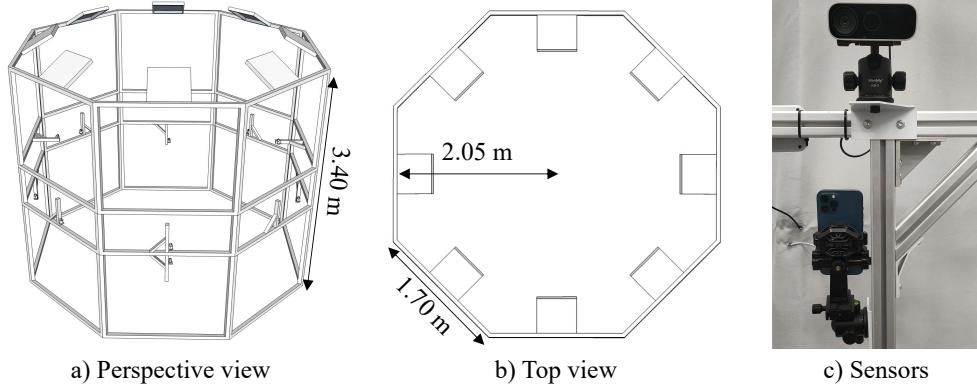a) Perspective view     b) Top view     c) Sensors

Figure 2. Hardware setup. a) and b) we build a octagonal prism-shaped framework to accommodate the data collection system. c) sensors used to collect sequential data include ten Azure Kinects and an iPhone 12 Pro Max. Besides, an Artec Eva is used to produce high-resolution static scans of the subjects.

a body scan of resolution up to 0.2 mm and accuracy up to 0.1 mm. A typical scan consists of $300k$ to $500k$ faces and $100k$ to $300k$ vertices, with a 4K (4096×4096) resolution texture map.

### 3.2. Two-Stage Calibration

**Image-based Calibration.** To obtain a coarse calibration, we first perform image-based calibration following the general steps in Zhang's method [116]. However, we highlight that Kinect's active IR depth cameras encounter over-exposure with regular chessboards. Hence, we customize a light absorbent material to cover the black squares of the chessboard pattern. In this way, we acquire reasonably accurate extrinsic calibration for Kinects and iPhones.

**Geometry-based Calibration.** Image-based calibration is unfortunately not accurate enough to reconstruct good-quality mesh. Hence, we propose to take advantage of the depth information in a geometry-based calibration stage. We empirically verify that image-based calibration serves as a good initialization for geometry-based calibration. Hence, we randomly place stacked cubes inside the framework. After that, we convert captured depth maps to point clouds and apply multi-way ICP registration [14] to refine the calibration.

### 3.3. Synchronization

**Kinects.** As the Azure Kinect implements the Time-of-Flight principle, it actively illuminates the scene multiple times (nine exposures in our system) for depth computation. To avoid interference between individual sensors, we use the synchronization cables to propagate a unified clock in a daisy chain fashion, and reject any image that is 33 ms or above out of synchronization. We highlight that there is only a 1450-us interval between exposures of 160 us;

our system of ten Kinects reaches the theoretical maximum number.

**Kinect-iPhone.** Due to hardware limitations, we cannot apply the synchronization cable to the iPhone. We circumvent this challenge by implementing a TCP-based communication protocol that computes an offset between the Kinect clock and the iPhone ARKit clock. As iPhone is recording at 60 FPS, we then use the offset to map the closest iPhone frames to Kinect frames. Our test shows the synchronization error is constrained below 33 ms.

## 4. Toolchain

To handle the large volume of data, we develop an automatic toolchain to provide annotations such as keypoints and SMPL parameters. Moreover, dynamic sequences of textured mesh are also reconstructed. The pipeline is illustrated in Fig. 3. Note that there is a human inspection stage to reject low-quality data with erroneous annotations.

### 4.1. Keypoint Annotation

There are two stages of keypoint annotation (I and II) in the toolchain. For stage I, virtual cameras are placed around the minimally clothed body scan to render multi-view images. For stage II, the color images from multi-view RGB-D are used. The core ideas of the keypoint annotation are demonstrated below, with the detailed algorithm in the Appendix.

**Multi-view 2D Keypoint Detection.** We employ the whole-body pose model that includes body, hand and face 2D keypoints $\hat{\mathcal{P}}_{2D} \in \mathbb{R}^{P \times 2}$, where $P = 133$. A large deep learning model HRNet-w48 [92] is used which achieves AP 66.1 and AR 74.3 on COCO whole-body benchmark [37].

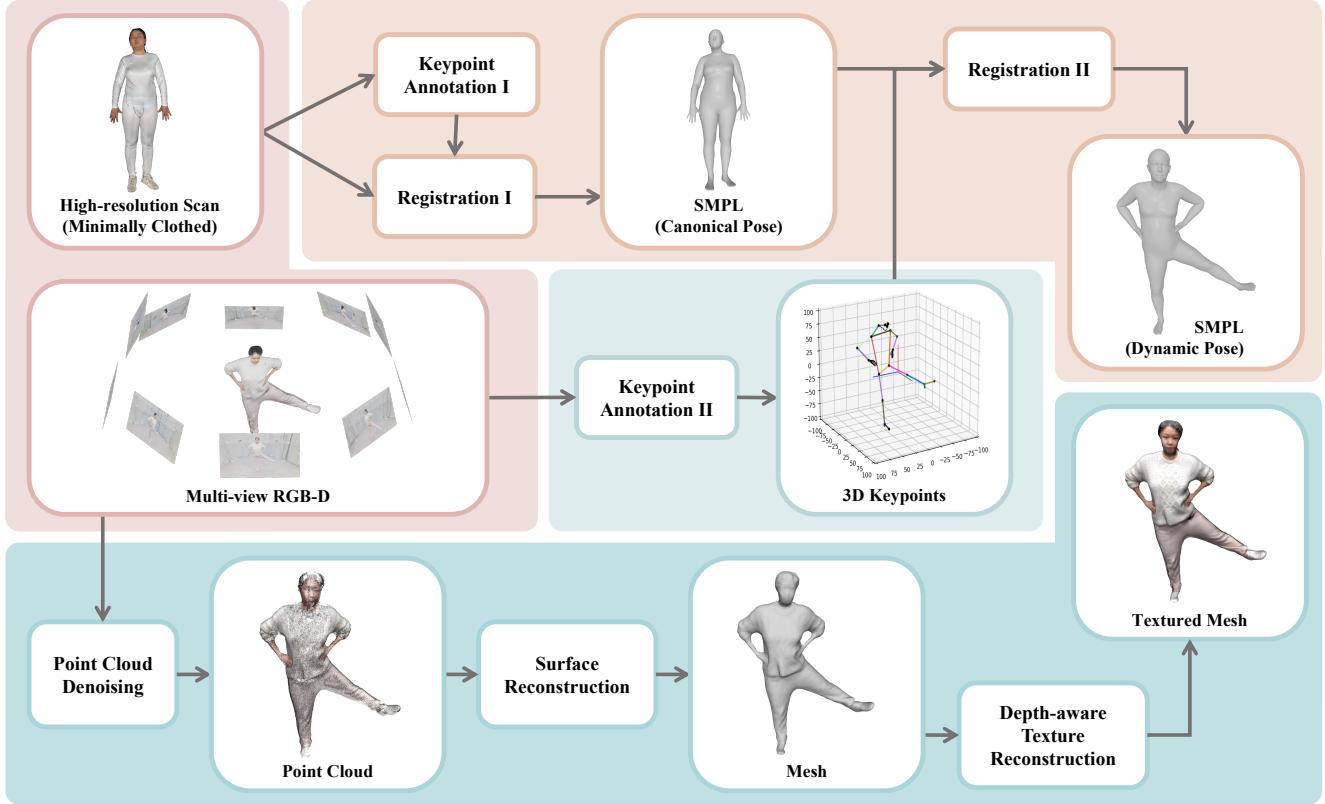**3D Keypoint Triangulation.** As the camera intrinsic and

4

Figure 3. Our toolchain produces multiple annotation formats such as 3D keypoint sequences, SMPL parameter sequences, and textured mesh sequences

extrinsic parameters are available, we triangulate 3D keypoints $\mathcal{P}_{3D} \in \mathbb{R}^{P \times 3}$ with the multi-view 2D estimated keypoints $\hat{\mathcal{P}}_{2D}$. However, 2D keypoints from any single view may not be always reliable. Hence, we use the following strategies to improve the quality of 3D keypoints. 1) *Keypoint selection*. To avoid the influence of poor-quality estimated 2D keypoints, we use a threshold $\tau_k$ to remove keypoints with a low confidence score. 2) *Camera selection*. As our system consists of ten Kinects, we exploit the redundancy to remove low-quality views. We only keep camera views with reprojection errors that are top-$k$ smallest [41] and no larger than a threshold $\tau_c$. 3) *Smoothness constraint*. Due to inevitable occlusion in the single view, the estimated 2D keypoints often have jitters. To alleviate the issue, we develop a smoothness loss to minimize the difference between consecutive triangulated 3D keypoints. Note that we design the loss weight to be inversely proportional to average speed, in order to remove jitters without compromising the ability to capture fast body motions. 4) *Bone length constraint*. As human bone length is constant, the per-frame bone length is constrained towards the median bone length $\mathcal{B}$ pre-computed from the initial triangulated 3D keypoints. The constraints are formulated as Eq. 1:

$$
\begin{aligned}
E_{tri} = & \lambda_1 \sum_{t=0}^{T-1} \|\mathcal{P}_{3D}(t+1) - \mathcal{P}_{3D}(t)\| + \\
& \lambda_2 \sum_{(i,j) \in \mathcal{I}_{\mathcal{B}}} \|\mathcal{B}_{i,j} - f_{\mathcal{B}}(\mathcal{P}_{3D}(i,j))\|
\end{aligned}
\tag{1}
$$

where $\mathcal{I}_{\mathcal{B}}$ contains the indices of connected keypoints and $f_{\mathcal{B}}(\cdot)$ calculates the average bone length of a given 3D keypoint sequence. Note that 3) and 4) are jointly optimized.

**2D Keypoint Projection.** To obtain high-quality 2D keypoints $\mathcal{P}_{2D} \in \mathbb{R}^{P \times 2}$, we project the triangulated 3D keypoints to image space via calibrated camera parameters. Note that this step is only needed for stage II keypoint annotation.

**Keypoint Quality.** We use $\mathcal{P}_{2D}$ and $\mathcal{P}_{3D}$ as keypoint annotations for 2D Pose Estimation and 3D Pose Estimation, respectively. To gauge the accuracy of the automatic keypoint annotation pipeline, we manually annotate a subset of data. The average Euclidean distance between annotated 2D keypoints and reprojected 2D keypoints $\mathcal{P}_{2D}$ is 15.13 pixels.
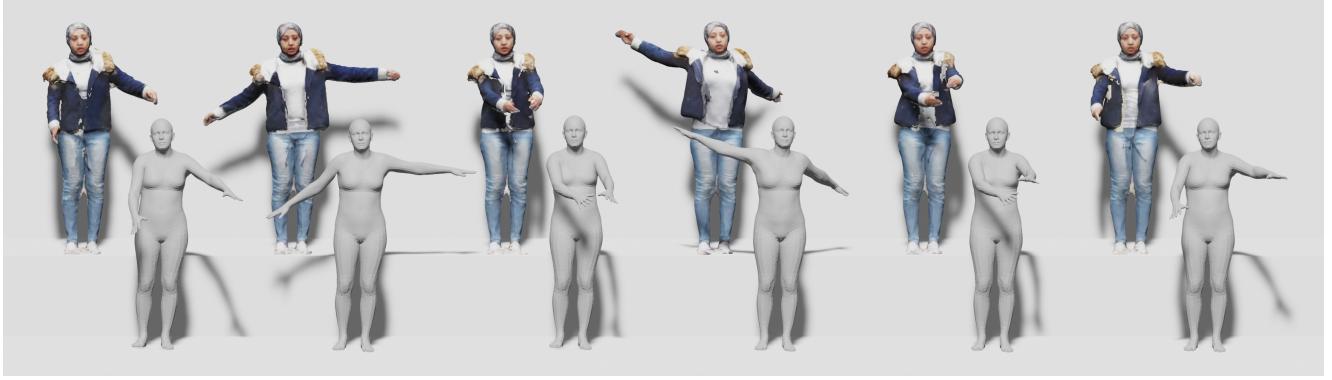
5

Figure 4. HuMMan provides synchronized sequences of multiple data formats and annotations. Here we demonstrate textured mesh sequences and SMPL parameter sequences

## 4.2. Human Parametric Model Registration

We select SMPL [61] as the human parametric model for its popularity. There are two stages of registration (I and II). Stage I is used to obtain accurate shape parameters from the static high-resolution scan, whereas stage II is used to obtain pose parameters from the dynamic sequence, with shape parameters from stage I. The registration is formulated as an optimization task to obtain SMPL pose parameters $\theta \in \mathbb{R}^{n \times 72}$, shape parameters $\beta \in \mathbb{R}^{n \times 10}$ (stage I only) and translation parameters $t \in \mathbb{R}^{n \times 3}$ where $n$ is the number of frames ($n = 1$ for stage I), with the following energy terms and constraints. We show a sample sequence of SMPL models with reconstructed textured mesh in Fig. 4.

**Keypoint Energy.** SMPLify [5] estimates camera parameters to leverage 2D keypoint supervision, which may be prone to depth and scale ambiguity. Hence, we develop the keypoint energy on 3D keypoints. For simplicity, we denote $P_{3D}$ as $P$. $\mathcal{J}$ is the joint regressor and $\mathcal{M}$ is the parametric model. We formulate the energy term:

$$E_{\mathcal{P}}(\theta, \beta, t) = \frac{1}{|\mathcal{P}|} \sum_{p}^{\mathcal{P}} \|\mathcal{J}(\mathcal{M}(\theta, \beta, t)) - p\| \quad (2)$$

**Surface Energy.** To supplement 3D keypoints that do not provide sufficient constraint for shape parameters, we add an additional surface energy term for registration on the high-resolution minimally clothed scans in stage I only. We use bi-directional Chamfer distance to gauge the difference between two mesh surfaces:

$$E_S = \frac{1}{|\mathcal{V}_H|} \sum_{v_H \in \mathcal{V}_H} \min_{v_S \in \mathcal{V}_S} \|v_H - v_S\| +$$
$$\frac{1}{|\mathcal{V}_S|} \sum_{v_S \in \mathcal{V}_S} \min_{v_H \in \mathcal{V}_H} \|v_H - v_S\| \quad (3)$$

where $\mathcal{V}_H$ and $\mathcal{V}_S$ are the mesh vertices of the high-resolution scan and SMPL.

**Shape Consistency.** Unlike existing work [77] that enforces an inter-beta energy term due to the lack of minimally clothed scan of each subject, we obtain accurate shape parameters from the high-resolution scan that allow us to apply constant beta parameters in the registration in stage II.

**Full-body Joint Angle Prior.** Joint rotation limitations serve as an important constraint to prevent unnaturally twisted poses. We extend existing work [5, 78] that only applies constraints on elbows and knees to all $J = 23$ joints in SMPL. The constraint is formulated as a strong penalty outside the plausible rotation range (with more details included in the Appendix):

$$E_a = \frac{1}{J \times 3} \sum_{j}^{J \times 3} exp(\max(\theta_i - \theta_i^u, 0) + \max(\theta_i^l - \theta_i, 0)) - 1$$
$$(4)$$

where $\theta_i^u$ and $\theta_i^l$ are the upper and lower limit of a rotation angle. Note that each joint rotation is converted to three Euler angles which can be interpreted as a series of individual rotations to decouple the original axis-angle representation.

## 4.3. Textured Mesh Reconstruction

**Point Cloud Reconstruction and Denoising.** We convert depth maps to point clouds and transform them into a world coordinate system with camera extrinsic parameters. However, the depth images captured by Kinect contain noisy pixels, which are prominent at subject boundaries where the depth gradient is large. To solve this issue, we first generate a binary boundary mask through edge finding with Laplacian of Gaussian Filters. Since our cameras have highly overlapped views to supplement points for one another, we apply a more aggressive threshold to remove boundary pixels. After the point cloud is reconstructed from the denoised depth images, we apply Statistical Outlier Removal [27] to further remove sprinkle noises.

**Geometry and Depth-aware Texture Reconstruction.**

6

Figure 5. Examples of SMPL registered on high-resolution static body scans for accurate shape parameters. The subjects are instructed to wear tight clothes for this scan. Note that each subject has another naturally clothed scan



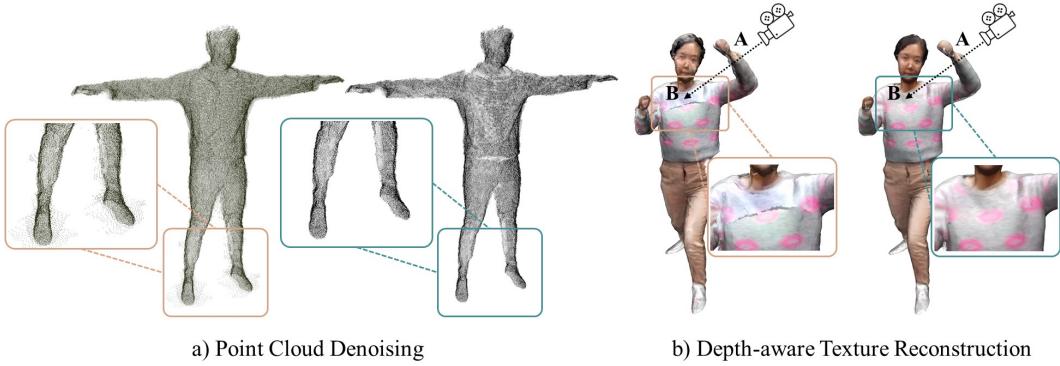a) Point Cloud Denoising                    b) Depth-aware Texture Reconstruction

Figure 6. Key steps to textured mesh reconstruction. a) Point cloud denoising removes noisy points. b) Depth-aware texture reconstruction prevents texture miss projection artifacts (such as projecting texture at point A to point B) due to misalignment between the actual subject and the reconstructed geometry

With complete and dense point cloud reconstructed, we apply Poisson Surface Reconstruction with envelope constraints [44] to reconstruct the watertight mesh. However, due to inevitable self-occlusion in complicated poses, interpolation artifacts arise from missing depth information, which leads to a shrunk or a dilated geometry. These artifacts are negligible for geometry reconstruction. However, a prominent artifact appears when projecting a texture onto the mesh even if the inconsistency between the true surface and the reconstructed surface is small. Hence, we extend MVS-texturing [100] to be depth-aware in texture reconstruction. We render the reconstructed mesh back into the camera view and compare the rendered depth map with the original depth map to generate the difference mask. We then mask out all the misalignment regions where the depth difference exceeds a threshold $\tau_d$. The masked regions do not contribute to texture projection. As shown in Fig. 6(b), the depth-aware texture reconstruction is more accurate and visually pleasing.

## 5. Action Set

Understanding human actions is a long-standing computer vision task. In this section, we elaborate on the two principles, following which we design the action set of 500 actions: *completeness* and *unambiguity*. More details are included in the Appendix.

**Completeness.** We build the action set to cover plausible human movements as much as possible. Compared to the popular 3D action recognition dataset NTU-RGBD-120 [58] whose actions are focused on upper body movements, we employ a hierarchical design to first divide possible actions into upper extremity, lower limbs, and whole-body movements. Such design allows us to achieve a balance between various body parts instead of over-emphasizing a specific group of movements. Note that we define whole body movements to be actions that require multiple body parts to collaborate, including different poses of the body trunk (*e.g.* lying down and sprawling). Fig. 7(c) demonstrates the action hierarchy and examples of interesting actions that are vastly diverse.

(a) Schematic Diagram (Front View)

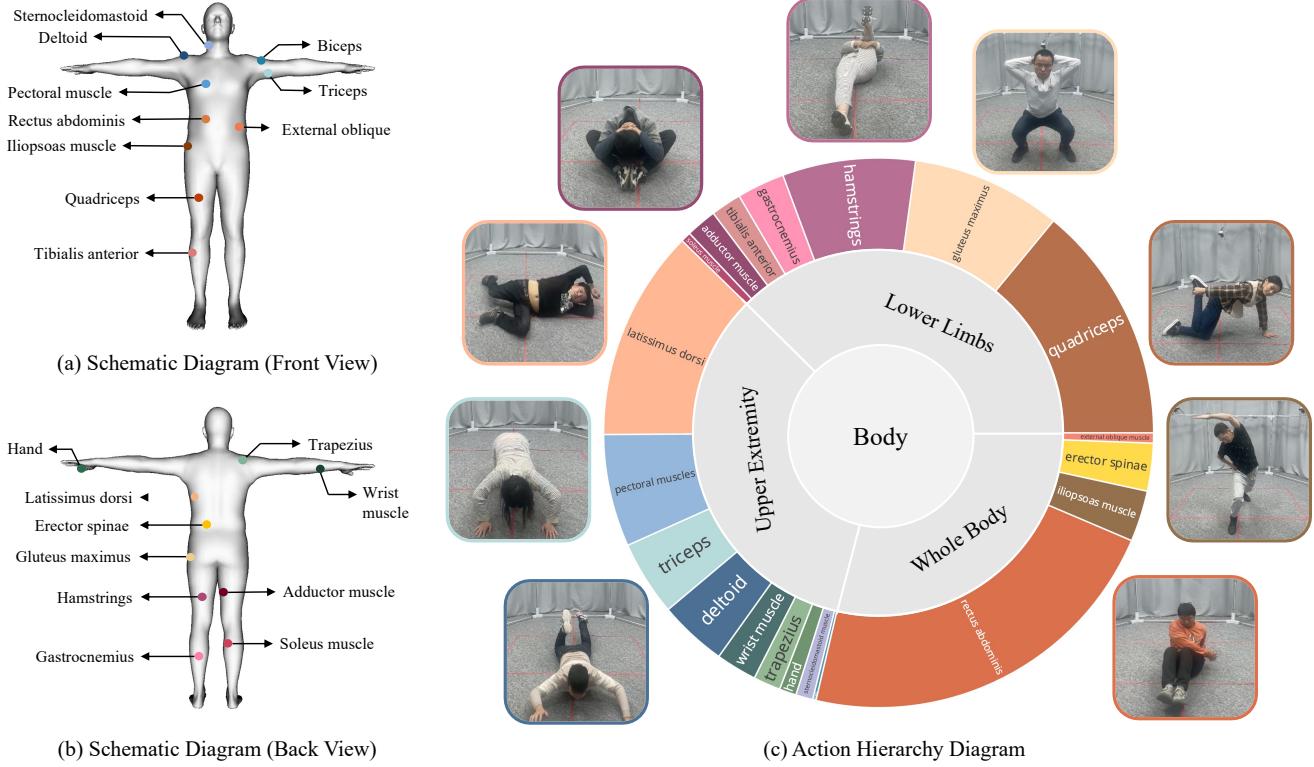(b) Schematic Diagram (Back View)

(c) Action Hierarchy Diagram

Figure 7. Schematic diagram of muscles from a) front and b) back views. c) HuMMan categorizes 500 actions hierarchically, first by body parts to achieve *complete* body coverage, then by driving muscles for *unambiguous* action definition

**Unambiguity.** Instead of providing a general description of the motions [11, 32, 42, 68, 71, 91, 99], we argue that the action classes should be clearly defined and are easy to identify and reproduce. Inspired by the fact that all human actions are the result of muscular contractions, we propose a *muscle-driven* strategy to systematically design the action set from the perspective of human anatomy. As illustrated in Fig. 7(a)(b), 20 major muscles are identified by professionals in fitness and yoga training, who then put together a list of standard movements associated with these muscles. Moreover, we cross-check with the action definitions from existing datasets [7, 11, 15, 22, 38, 42, 54, 58] to ensure a wide coverage is achieved.

## 6. Subjects

**Diversity.** HuMMan consists of 1000 subjects with a wide coverage of genders, ages, body shapes (heights, weights), and ethnicity. The subjects are instructed to wear their personal daily clothes to achieve a large collection of natural appearances. We demonstrate examples of high-resolution scans of the subjects in Fig. 8. We include statistics in the Appendix.

**Ethics.** HuMMan involves a large number of human sub-

jects so that we pay special attention to address ethic concerns. The recruitment process is conducted on an entirely voluntary basis. Actors and actresses who participate in HuMMan are well-informed, with legal agreements signed to acknowledge that the data will be made public for research purposes.

## 7. Experiments

In this section, we evaluate popular methods from various research fields on HuMMan. To constrain the training within a reasonable computation budget, we sample 10% of data and split them into training and testing sets for both Kinects and iPhone. The details are included in the Appendix.

Table 2. **Action Recognition**

| Method | Top-1 (%)↑ | Top-5 (%)↑ |
|---|---|---|
| ST-GCN | 72.5 | 94.3 |
| 2s-AGCN | 74.1 | 95.4 |

**Action Recognition.** HuMMan provides action labels and 3D skeletal positions, which can verify its usefulness

Figure 8. HuMMan contains 1000 subjects with diverse appearances. For each subject, a naturally clothed high-resolution scan is provided

on 3D action recognition. Specifically, we train popular graph-based methods (STGCN [105] and 2s-AGCN [89]) on HuMMan. Results are shown in Table 2. Compared to NTU RGB+D, a large-scale 3D action recognition dataset and a standard benchmark that contains 120 actions [58], HuMMan may be more challenging since 2s-AGCN [89] achieves Top-1 accuracy of 88.9% and 82.9% on NTU RGB+D 60 and 120 respectively, but 74.1% only on HuM-Man. The difficulties come from the whole-body coverage design in our action set, instead of over-emphasis on certain body parts (*e.g.* NTU RGB+D has a large proportion of upper body movements). Moreover, we observe a significant gap between Top-1 and Top-5 accuracy (~30%). We attribute this phenomenon to the fact that there are plenty of *intra-actions* in HuMMan. For example, there are similar variants of push-ups such as quadruped push-ups, kneeling push-ups, and leg push-ups. This challenges the model to pay more attention to the fine-grained differences in these actions. Hence, we find HuMMan would serve as an indicative benchmark for fine-grained action understanding.

Table 3. **3D Keypoint Detection**. PA: PA-MPJPE

| Train | Test | MPJPE ↓ | PA ↓ |
|---|---|---|---|
| FCN [66] | | | |
| HuMMan | HuMMan | 78.5 | 46.3 |
| H36M | AIST++ | 133.9 | 73.1 |
| HuMMan | AIST++ | 116.4 | 67.2 |
| Video3D [80] | | | |
| HuMMan | HuMMan | 73.1 | 43.5 |
| H36M | AIST++ | 128.5 | 72.0 |
| HuMMan | AIST++ | 109.2 | 63.5 |

**3D Keypoint Detection.** With the well-annotated 3D keypoints, HuMMan supports 3D keypoint detection. We employ popular 2D-to-3D lifting backbones [66, 80] as single-frame and multi-frame baselines on HuMMan. We experiment with different training and test settings to obtain the baseline results in Table 3. First, in-domain training and testing on HuMMan are provided. The values are slightly higher than the same baselines on Human3.6M [32] (on which FCN obtains MPJPE of 53.4 mm). Second, methods trained on HuMMan tend to generalize better than on Human3.6M. This may be attributed to HuMMan's diverse collection of subjects and actions.

Table 4. **3D Parametric Human Recovery**. Image- and point cloud-based methods are evaluated

| Method | MPJPE ↓ | PA-MPJPE ↓ |
|---|---|---|
| HMR | 54.78 | 36.14 |
| VoteHMR | 144.99 | 106.32 |

**3D Parametric Human Recovery.** HuMMan provides SMPL annotations, RGB and RGB-D sequences. Hence, we evaluate HMR [39], not only one of the first deep learning approaches towards 3D parametric human recovery but a fundamental component for follow-up works [46, 48], to represent image-based methods. In addition, we employ VoteHMR [57], a recent work that takes point clouds as the input. In Table 4, we find that HMR has achieved low MPJPE and PA-MPJPE, which may be attributed to the clearly defined action set and the training set already includes all action classes. However, VoteHMR is not performing well. We argue that existing point cloud-based methods [36, 57, 102] rely heavily on synthetic data for training and evaluation, whereas HuMMan provides gen-

Figure 9. We compare Function4D with HuMMan in textured mesh reconstruction

uine point clouds from commercial RGB-D sensors that remain challenging.

**Textured Mesh Reconstruction.** To fully demonstrate the capacity of HuMMan, we also provide the results of Function4D [107] as a baseline for textured mesh reconstruction since it combines both volumetric fusion and implicit surface reconstruction for volumetric capture in real-time. The results of Function4D, using 4 (ID: 0,3,6,9) views, are shown in Fig. 9. Note that benefiting from the multimodality signals in HuMMan, various surface reconstruction methods like PIFu [86](using only RGB as input for textured human mesh reconstruction), 3D Self-Portrait [55] (using single-view RGBD video for 3D portrait reconstruction), and CON [81] (using multi-view depth point cloud as input for complete mesh reconstruction) are also supported.

Table 5. **Mobile Device**. The models are trained with different training sets, and evaluated on HuMMan iPhone test set. Kin.: Kinect training set. iPh.: iPhone training set. PA: PA-MPJPE

| Method | Kin. | iPh. | MPJPE ↓ | PA ↓ |
|--------|------|------|---------|------|
| HMR | ✓ | - | 97.81 | 52.74 |
| HMR | - | ✓ | 72.62 | 41.86 |
| VoteHMR | ✓ | - | 255.71 | 162.00 |
| VoteHMR | - | ✓ | 83.18 | 61.69 |

**Mobile Device.** It is under-explored that if model trained with the regular device is readily transferable to the mobile device. In Table 5, we study the performance gaps across devices. For the image-based method, we find that there exists a considerable domain gap across devices, despite that they have similar resolutions. Moreover, for the point cloud-based method, the domain gap is much more significant as the mobile device tends to have much sparser point clouds as a result of lower depth map resolution. Hence, it remains a challenging problem to transfer knowledge across devices, especially for point cloud-based methods.

## 8. Discussion

We present HuMMan, a large-scale 4D human dataset that features multi-modal data and annotations, inclusion of mobile device, a comprehensive action set, and support for multiple tasks. Our experiments point out interesting directions that await future research, such as fine-grained action recognition, point cloud-based parametric human estimation, dynamic mesh sequence reconstruction, transferring knowledge across devices, and potentially, multi-task joint training. We hope HuMMan would facilitate the development of better algorithms for sensing and modeling humans.

## References

[1] Thiemo Alldieck, Marcus Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. Video based reconstruction of 3d people models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8387–8397, 2018. 3, 17

[2] Mykhaylo Andriluka, Umar Iqbal, Eldar Insafutdinov, Leonid Pishchulin, Anton Milan, Juergen Gall, and Bernt Schiele. Posetrack: A benchmark for human pose estimation and tracking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5167–5176, 2018. 2, 3, 17

[3] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *Proceedings of the IEEE Conference on computer Vision and Pattern Recognition*, pages 3686–3693, 2014. 2, 3, 17

[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*, pages 311–329. Springer, 2020. 3

[5] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In *European conference on computer vision*, pages 561–578. Springer, 2016. 3, 6

[6] Federica Bogo, Javier Romero, Gerard Pons-Moll, and Michael J Black. Dynamic faust: Registering human bodies in motion. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6233–6242, 2017. 2, 3, 17

[7] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the ieee conference on computer vision and pattern recognition*, pages 961–970, 2015. 8

[8] Zhongang Cai, Junzhe Zhang, Daxuan Ren, Cunjun Yu, Haiyu Zhao, Shuai Yi, Chai Kiat Yeo, and Chen Change Loy. Messytable: Instance association in multiple camera views. In *European Conference on Computer Vision*, pages 1–16. Springer, 2020. 21

[9] Zhongang Cai, Mingyuan Zhang, Jiawei Ren, Chen Wei, Daxuan Ren, Jiatong Li, Zhengyu Lin, Haiyu Zhao, Shuai Yi, Lei Yang, et al. Playing for 3d human recovery. *arXiv preprint arXiv:2110.07588*, 2021. 3

[10] Zhe Cao, Hang Gao, Karttikeya Mangalam, Qi-Zhi Cai, Minh Vo, and Jitendra Malik. Long-term human motion prediction with scene context. In *European Conference on Computer Vision*, pages 387–404. Springer, 2020. 3

[11] Joao Carreira, Eric Noland, Chloe Hillier, and Andrew Zisserman. A short note on the kinetics-700 human action dataset. *arXiv preprint arXiv:1907.06987*, 2019. 3, 8, 17

[12] Yilun Chen, Zhicheng Wang, Yuxiang Peng, Zhiqiang Zhang, Gan g Yu, and Jian Sun. Cascaded pyramid network for multi-person pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7103–7112, 2018. 3, 20, 21

[13] Hongsuk Choi, Gyeongsik Moon, and Kyoung Mu Lee. Beyond static features for temporally consistent 3d human pose and shape from a video. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3

[14] Sungjoon Choi, Qian-Yi Zhou, and Vladlen Koltun. Robust reconstruction of indoor scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5556–5565, 2015. 4

[15] Jihoon Chung, Cheng-hsin Wuu, Hsuan-ru Yang, Yu-Wing Tai, and Chi-Keung Tang. Haa500: Human-centric atomic action dataset with curated videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13465–13474, 2021. 2, 3, 8, 17

[16] Qi Fang, Qing Shuai, Junting Dong, Hujun Bao, and Xiaowei Zhou. Reconstructing 3d human pose by watching humans in the mirror. In *CVPR*, 2021. 17

[17] Christoph Feichtenhofer. X3d: Expanding architectures for efficient video recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 203–213, 2020. 3

[18] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6202–6211, 2019. 3

[19] Yasutaka Furukawa and Jean Ponce. Accurate, dense, and robust multiview stereopsis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(8):1362–1376, 2010. 3

[20] Ran Gal, Yonatan Wexler, Eyal Ofek, Hugues Hoppe, and Daniel Cohen-Or. Seamless montage for texturing models. *Computer Graphics Forum*, 29(2):479–486, 2010. 3

[21] Georgios Georgakis, Ren Li, Srikrishna Karanam, Terrence Chen, Jana Košecká, and Ziyan Wu. Hierarchical kinematic human mesh recovery. In *European Conference on Computer Vision*, pages 768–784. Springer, 2020. 3

[22] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, et al. Ava: A video dataset of spatio-temporally localized atomic visual actions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6047–6056, 2018. 2, 3, 8, 17

[23] Riza Alp Guler and Iasonas Kokkinos. Holopose: Holistic 3d human reconstruction in-the-wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10884–10894, 2019. 3

[24] Marc Habermann, Lingjie Liu, Weipeng Xu, Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. Real-time deep dynamic characters. *ACM Transactions on Graphics (TOG)*, 40(4):1–16, 2021. 3, 17

[25] Marc Habermann, Weipeng Xu, Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. Livecap: Real-time human performance capture from monocular video. *ACM Transactions On Graphics (TOG)*, 38(2):1–17, 2019. 3, 17

[26] Marc Habermann, Weipeng Xu, Michael Zollhofer, Gerard Pons-Moll, and Christian Theobalt. Deepcap: Monocular human performance capture using weak supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5052–5063, 2020. 3, 17

[27] Victoria Hodge and Jim Austin. A survey of outlier detection methodologies. *Artificial intelligence review*, 22(2):85–126, 2004. 6

[28] Fangzhou Hong, Liang Pan, Zhongang Cai, and Ziwei Liu. Garment4d: Garment reconstruction from point cloud sequences. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021. 3

[29] Fangzhou Hong, Liang Pan, Zhongang Cai, and Ziwei Liu. Versatile multi-modal pre-training for human-centric perception. *arXiv preprint arXiv:2203.13815*, 2022. 2

[30] Jian-Fang Hu, Wei-Shi Zheng, Jianhuang Lai, and Jianguo Zhang. Jointly learning heterogeneous features for rgb-d activity recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5344–5352, 2015. 2, 3, 17

[31] Fuyang Huang, Ailing Zeng, Minhao Liu, Qiuxia Lai, and Qiang Xu. Deepfuse: An imu-aware network for real-time 3d human pose estimation from multi-view image. *arXiv preprint arXiv:1912.04071*, 2019. 3

[32] 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 transactions on pattern analysis and machine intelligence*, 36(7):1325–1339, 2013. 1, 2, 3, 8, 9, 17, 18

[33] Karim Iskakov, Egor Burkov, Victor Lempitsky, and Yury Malkov. Learnable triangulation of human pose. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7718–7727, 2019. 3

[34] Shahram Izadi, David Kim, Otmar Hilliges, David Molyneaux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew Davison, et al. Kinectfusion: real-time 3d reconstruction and interaction using a moving depth camera. In *Proceedings of the 24th annual ACM symposium on User interface software and technology*, pages 559–568, 2011. 3

[35] Hueihan Jhuang, Juergen Gall, Silvia Zuffi, Cordelia Schmid, and Michael J Black. Towards understanding action recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 3192–3199, 2013. 3, 17

[36] Haiyong Jiang, Jianfei Cai, and Jianmin Zheng. Skeleton-aware 3d human shape reconstruction from point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5431–5441, 2019. 3, 9

[37] Sheng Jin, Lumin Xu, Jin Xu, Can Wang, Wentao Liu, Chen Qian, Wanli Ouyang, and Ping Luo. Whole-body human pose estimation in the wild. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 4

[38] 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 *Proceedings of the IEEE International Conference on Computer Vision*, pages 3334–3342, 2015. 2, 3, 8, 17

[39] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7122–7131, 2018. 9

[40] Angjoo Kanazawa, Jason Y Zhang, Panna Felsen, and Jitendra Malik. Learning 3d human dynamics from video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5614–5623, 2019. 3

[41] Pierre Karashchuk, Katie L Rupp, Evyn S Dickinson, Sarah Walling-Bell, Elischa Sanders, Eiman Azim, Bingni W Brunton, and John C Tuthill. Anipose: a toolkit for robust markerless 3d pose estimation. *Cell reports*, 36(13):109730, 2021. 5

[42] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014. 3, 8, 17

[43] Misha Kazhdan, Ming Chuang, Szymon Rusinkiewicz, and Hugues Hoppe. Poisson surface reconstruction with envelope constraints. *Computer Graphics Forum (Proc. Symposium on Geometry Processing)*, 39(5), July 2020. 3

[44] Misha Kazhdan, Ming Chuang, Szymon Rusinkiewicz, and Hugues Hoppe. Poisson surface reconstruction with envelope constraints. In *Computer graphics forum*, volume 39, pages 173–182. Wiley Online Library, 2020. 7

[45] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Trans. Graph.*, 32(3), jul 2013. 3

[46] Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. Vibe: Video inference for human body pose and shape estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5253–5263, 2020. 9

[47] Muhammed Kocabas, Chun-Hao P Huang, Otmar Hilliges, and Michael J Black. Pare: Part attention regressor for 3d human body estimation. *arXiv preprint arXiv:2104.08527*, 2021. 3

[48] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2252–2261, 2019. 3, 9

[49] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large video database for human motion recognition. In *2011 International conference on computer vision*, pages 2556–2563. IEEE, 2011. 3, 17

[50] Jiefeng Li, Siyuan Bian, Ailing Zeng, Can Wang, Bo Pang, Wentao Liu, and Cewu Lu. Human pose regression with residual log-likelihood estimation. In *ICCV*, 2021. 3

[51] Jiefeng Li, Chao Xu, Zhicun Chen, Siyuan Bian, Lixin Yang, and Cewu Lu. Hybrik: A hybrid analytical-neural inverse kinematics solution for 3d human pose and shape estimation. In *CVPR*, pages 3383–3393. Computer Vision Foundation / IEEE, 2021. 3

[52] Ruilong Li, Shan Yang, David A Ross, and Angjoo Kanazawa. Ai choreographer: Music conditioned 3d dance generation with aist++. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13401–13412, 2021. 2, 3, 17

[53] Wanqing Li, Zhengyou Zhang, and Zicheng Liu. Action recognition based on a bag of 3d points. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition-Workshops*, pages 9–14. IEEE, 2010. 3, 17

[54] Yong-Lu Li, Xinpeng Liu, Xiaoqian Wu, Yizhuo Li, Zuoyu Qiu, Liang Xu, Yue Xu, Hao-Shu Fang, and Cewu Lu. Hake: A knowledge engine foundation for human activity understanding, 2022. 8

[55] Zhe Li, Tao Yu, Zerong Zheng, and Yebin Liu. Robust and accurate 3d self-portraits in seconds. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021. 10

[56] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 1, 2, 3, 17, 21

[57] Guanze Liu, Yu Rong, and Lu Sheng. Votehmr: Occlusion-aware voting network for robust 3d human mesh recovery from partial point clouds. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 955–964, 2021. 3, 9

[58] Jun Liu, Amir Shahroudy, Mauricio Perez, Gang Wang, Ling-Yu Duan, and Alex C Kot. Ntu rgb+ d 120: A large-scale benchmark for 3d human activity understanding. *IEEE transactions on pattern analysis and machine intelligence*, 42(10):2684–2701, 2019. 2, 3, 7, 8, 9, 17

[59] Stephen Lombardi, Jason Saragih, Tomas Simon, and Yaser Sheikh. Deep appearance models for face rendering. 37(4), jul 2018. 3

[60] Matthew Loper, Naureen Mahmood, and Michael J Black. Mosh: Motion and shape capture from sparse markers. *ACM Transactions on Graphics (TOG)*, 33(6):1–13, 2014. 3, 18

[61] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. SMPL: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):1–16, 2015. 3, 6

[62] Zhengyi Luo, S Alireza Golestaneh, and Kris M Kitani. 3d human motion estimation via motion compression and refinement. In *Proceedings of the Asian Conference on Computer Vision*, 2020. 3

[63] Qianli Ma, Jinlong Yang, Anurag Ranjan, Sergi Pujades, Gerard Pons-Moll, Siyu Tang, and Michael J Black. Learning to dress 3d people in generative clothing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6469–6478, 2020. 1, 2, 3, 17

[64] Spandan Madan, Timothy Henry, Jamell Dozier, Helen Ho, Nishchal Bhandari, Tomotake Sasaki, Frédo Durand, Hanspeter Pfister, and Xavier Boix. When and how do cnns generalize to out-of-distribution category-viewpoint combinations? *arXiv preprint arXiv:2007.08032*, 2020. 21

[65] Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. Amass: Archive of motion capture as surface shapes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5442–5451, 2019. 2, 17

[66] Julieta Martinez, Rayat Hossain, Javier Romero, and James J Little. A simple yet effective baseline for 3d human pose estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2640–2649, 2017. 3, 9, 21

[67] Joumana Medlej. Human anatomy fundamentals: Flexibility and joint limitations. 17

[68] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *2017 international conference on 3D vision (3DV)*, pages 506–516. IEEE, 2017. 2, 3, 8, 17

[69] Dushyant Mehta, Oleksandr Sotnychenko, Franziska Mueller, Weipeng Xu, Mohamed Elgharib, Pascal Fua, Hans-Peter Seidel, Helge Rhodin, Gerard Pons-Moll, and Christian Theobalt. Xnect: Real-time multi-person 3d motion capture with a single rgb camera. *ACM Transactions on Graphics (TOG)*, 39(4):82–1, 2020. 3

[70] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 3

[71] Mathew Monfort, Alex Andonian, Bolei Zhou, Kandan Ramakrishnan, Sarah Adel Bargal, Tom Yan, Lisa Brown, Quanfu Fan, Dan Gutfreund, Carl Vondrick, et al. Moments in time dataset: one million videos for event understanding. *IEEE transactions on pattern analysis and machine intelligence*, 42(2):502–508, 2019. 3, 8, 17

[72] Gyeongsik Moon and Kyoung Mu Lee. I2l-meshnet: Image-to-lixel prediction network for accurate 3d human pose and mesh estimation from a single RGB image. In *ECCV (7)*, volume 12352 of *Lecture Notes in Computer Science*, pages 752–768. Springer, 2020. 3

[73] Richard A Newcombe, Dieter Fox, and Steven M Seitz. Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 343–352, 2015. 3

[74] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *European conference on computer vision*, pages 483–499. Springer, 2016. 3

[75] Mohamed Omran, Christoph Lassner, Gerard Pons-Moll, Peter Gehler, and Bernt Schiele. Neural body fitting: Unifying deep learning and model based human pose and shape estimation. In *2018 international conference on 3D vision (3DV)*, pages 484–494. IEEE, 2018. 3

[76] Ahmed A. A. Osman, Timo Bolkart, and Michael J. Black. STAR: sparse trained articulated human body regressor. In *ECCV (6)*, volume 12351 of *Lecture Notes in Computer Science*, pages 598–613. Springer, 2020. 3

[77] Priyanka Patel, Chun-Hao P Huang, Joachim Tesch, David T Hoffmann, Shashank Tripathi, and Michael J Black. AGORA: Avatars in geography optimized for regression analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13468–13478, 2021. 3, 6

[78] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10975–10985, 2019. 3, 6

[79] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3d human pose and shape from a single color image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 459–468, 2018. 3

[80] Dario Pavllo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3d human pose estimation in video with temporal convolutions and semi-supervised training. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7753–7762, 2019. 3, 9, 21

[81] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *ECCV*, 2020. 3, 10

[82] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes

for novel view synthesis of dynamic humans. In *CVPR*, 2021. 2, 3, 17

[83] Haibo Qiu, Chunyu Wang, Jingdong Wang, Naiyan Wang, and Wenjun Zeng. Cross view fusion for 3d human pose estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4342–4351, 2019. 3

[84] Amit Raj, Julian Tanke, James Hays, Minh Vo, Carsten Stoll, and Christoph Lassner. Anr-articulated neural rendering for virtual avatars. In *arXiv:2012.12890*, 2020. 3

[85] Linden Research. Suggested bvh joint rotation limits. 17

[86] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2304–2314, 2019. 3, 10

[87] Amir Shahroudy, Jun Liu, Tian-Tsong Ng, and Gang Wang. Ntu rgb+ d: A large scale dataset for 3d human activity analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1010–1019, 2016. 2, 3, 17

[88] Dian Shao, Yue Zhao, Bo Dai, and Dahua Lin. Finegym: A hierarchical video dataset for fine-grained action understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2616–2625, 2020. 2, 3, 17

[89] Lei Shi, Yifan Zhang, Jian Cheng, and Hanqing Lu. Skeleton-based action recognition with multi-stream adaptive graph convolutional networks. *arXiv preprint arXiv:1912.06971*, 2019. 3, 9

[90] Lei Shi, Yifan Zhang, Jian Cheng, and Hanqing Lu. Skeleton-based action recognition with multi-stream adaptive graph convolutional networks. *IEEE Transactions on Image Processing*, 29:9532–9545, 2020. 3

[91] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 2, 3, 8, 17

[92] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5693–5703, 2019. 3, 4, 20, 21

[93] Yu Sun, Yun Ye, Wu Liu, Wenpeng Gao, Yili Fu, and Tao Mei. Human mesh recovery from monocular images via a skeleton-disentangled representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5349–5358, 2019. 3

[94] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Trans. Graph.*, 38(4), jul 2019. 3

[95] Du Tran, Heng Wang, Lorenzo Torresani, and Matt Feiszli. Video classification with channel-separated convolutional networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5552–5561, 2019. 3

[96] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proceedings*

*of the IEEE conference on Computer Vision and Pattern Recognition*, pages 6450–6459, 2018. 3

[97] Neel Trivedi, Anirudh Thatipelli, and Ravi Kiran Sarvadevabhatla. Ntu-x: An enhanced large-scale dataset for improving pose-based recognition of subtle human actions. *arXiv preprint arXiv:2101.11529*, 2021. 2, 3, 17

[98] Matthew Trumble, Andrew Gilbert, Charles Malleson, Adrian Hilton, and John P Collomosse. Total capture: 3d human pose estimation fusing video and inertial sensors. In *BMVC*, volume 2, pages 1–13, 2017. 3, 17

[99] Timo von Marcard, Roberto Henschel, Michael J Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 601–617, 2018. 1, 2, 3, 8, 17

[100] Michael Waechter, Nils Moehrle, and Michael Goesele. Let there be color! — Large-scale texturing of 3D reconstructions. In *Proceedings of the European Conference on Computer Vision*. Springer, 2014. 7

[101] Jiang Wang, Xiaohan Nie, Yin Xia, Ying Wu, and Song-Chun Zhu. Cross-view action modeling, learning and recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2649–2656, 2014. 3, 17

[102] Shaofei Wang, Andreas Geiger, and Siyu Tang. Locally aware piecewise transformation fields for 3d human mesh registration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7639–7648, 2021. 3, 9

[103] Donglai Xiang, Fabian Prada, Timur Bagautdinov, Weipeng Xu, Yuan Dong, He Wen, Jessica Hodgins, and Chenglei Wu. Modeling clothing as a separate layer for an animatable human avatar. *ACM Trans. Graph.*, 40(6), dec 2021. 3

[104] 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. 3

[105] Sijie Yan, Yuanjun Xiong, and Dahua Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. *arXiv preprint arXiv:1801.07455*, 2018. 3, 9

[106] Changqian Yu, Bin Xiao, Changxin Gao, Lu Yuan, Lei Zhang, Nong Sang, and Jingdong Wang. Lite-hrnet: A lightweight high-resolution network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10440–10450, 2021. 21

[107] Tao Yu, Zerong Zheng, Kaiwen Guo, Pengpeng Liu, Qionghai Dai, and Yebin Liu. Function4d: Real-time human volumetric capture from very sparse consumer rgbd sensors. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR2021)*, June 2021. 2, 3, 10, 17

[108] 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 *IEEE*

*Conference on Computer Vision and Pattern Recognition(CVPR)*, pages 7287–7296, Salt Lake City, June 2018. IEEE. 3

[109] Zhixuan Yu, J. S. Yoon, I. Lee, Prashanth Venkatesh, Jaesik Park, J. Yu, and H. Park. Humbi: A large multiview dataset of human body expressions. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2987–2997, 2020. 1, 2, 3, 17

[110] Ailing Zeng, Xiao Sun, Fuyang Huang, Minhao Liu, Qiang Xu, and Stephen Ching-Feng Lin. Srnet: Improving generalization in 3d human pose estimation with a split-and-recombine approach. In *ECCV*, 2020. 3

[111] Ailing Zeng, Xiao Sun, Lei Yang, Nanxuan Zhao, Minhao Liu, and Qiang Xu. Learning skeletal graph neural networks for hard 3d pose estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, 2021. 3

[112] Ailing Zeng, Lei Yang, Xuan Ju, Jiefeng Li, Jianyi Wang, and Qiang Xu. Smoothnet: A plug-and-play network for refining human poses in videos. *arXiv preprint arXiv:2112.13715*, 2021. 3

[113] Chao Zhang, Sergi Pujades, Michael J Black, and Gerard Pons-Moll. Detailed, accurate, human shape estimation from clothed 3d scan sequences. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4191–4200, 2017. 2, 3, 17

[114] Weiyu Zhang, Menglong Zhu, and Konstantinos G Derpanis. From actemes to action: A strongly-supervised representation for detailed action understanding. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2248–2255, 2013. 3, 17

[115] Yuxiang Zhang, Zhe Li, Liang An, Mengcheng Li, Tao Yu, and Yebin Liu. Lightweight multi-person total motion capture using sparse multi-view cameras. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5560–5569, October 2021. 3

[116] Zhengyou Zhang. A flexible new technique for camera calibration. *IEEE Transactions on pattern analysis and machine intelligence*, 22(11):1330–1334, 2000. 4

[117] Hang Zhao, Antonio Torralba, Lorenzo Torresani, and Zhicheng Yan. Hacs: Human action clips and segments dataset for recognition and temporal localization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8668–8678, 2019. 3, 17

[118] Long Zhao, Xi Peng, Yu Tian, Mubbasir Kapadia, and Dimitris N Metaxas. Semantic graph convolutional networks for 3d human pose regression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3425–3435, 2019. 3

[119] Zerong Zheng, Tao Yu, Yixuan Wei, Qionghai Dai, and Yebin Liu. Deephuman: 3d human reconstruction from a single image. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 1, 3, 17

## A. Appendix

We provide a more complete dataset comparison (Section B), and additional details of data collection (Section C), hardware (Section D), toolchain (Section E), action set (Section F), subjects (Section G) and experiments (Section H).

## B. A More Complete Dataset Comparison

In Table 6, we provide a more thorough comparison of HuMMan with similar datasets for 1) action recognition, 2) 2D and 3D keypoint detection, 3) 3D parametric human recovery, and 4) mesh reconstruction. We only include real datasets in the Table but there are also popular synthetic datasets [].

## C. Additional Details of Data Collection

The data collection has two stages for each subject. **1)** each subject receives two high-resolution scans, one with natural clothes on and the other with a tight-fitting suit on, both captured by the Artex Eva 3D Scanner. To ensure the high quality of the scans, the subjects are instructed to stand in a special pose (the *canonical pose*) on a turntable, that allows for a 360-degree full-body scanning with minimal self-occlusion. Each high-resolution scan includes an MTL information file, an OBJ mesh file, and a BMP texture file. **2)** After that static body scanning, the subject enters the framework and follows instructions to perform 40-60 actions, randomly sampled from the action set that contains 500 actions. Each action that a subject performs is a *sequence*, that consists of ten Kinect RGB-D sequences and an iPhone RGB-D sequence. We show sample frames collected with our hardware setup in Fig. 10. Each sequence takes 5-15 seconds and 150-450 frames at 30 FPS per view. We compress all sequential data in a custom data format *SMC* that is developed based on HDF5 format. The SMC file also contains additional information such as camera parameters, subject ID, and action ID.

## D. Additional Details of Hardware

### D.1. Sensors

We provide more details on the RGB-D sensor (Azure Kinect). We set operating mode to *NFOV unbinned* for the depth cameras, which results in the largest view overlap with the color camera and the densest point clouds. The depth camera in this mode has an FOV of $90° \times 59°$. The operating range of the depth sensor in this mode is between 0.5 m to 3.86 m. The typical systematic error of the depth sensor is less than 11 mm + 0.1% of distance with a standard deviation of less than 17 mm. In view of the limited FOV and depth error-distance relation, we design our aluminum framework such that the subject is around 2 m away

from the Kinects: at that distance, the FOV can accommodate the subject's whole body, without incurring any extra depth error.

### D.2. Synchronization

Our data sampling program runs on a workstation, and it 1) integrates the Kinect SDK, and 2) communicates with the iPhone app developed based on ARKit through TCP. Since there is no existing hardware approach to Kinect-iPhone synchronization, we develop a method to compute the difference between Kinect clock and iPhone ARKit clock $t_{K \to A}$. Hence, we first obtain the offset from the workstation to the Kinects $t_{K \to W}$ as

$$t_{K \to W} = t_W - t_K$$

where $t_K$ is the Kinect clock time and $t_W$ is the workstation's system time, obtained at the same moment. We also send a message to the iPhone app, which records down the iPhone system clock $t_I$ upon receiving the message and sends back a message to the workstation to complete a round trip. We compute the offset from the iPhone system clock to the workstation system clock $t_{W \to I}$ as

$$t_{W \to I} = t_I - t_W - \frac{t_{round}}{2}$$

where $t_{round}$ is the round trip time taken. Note that there is an additional offset between the ARKit clock and the iPhone system clock $t_{I \to A}$, computed as

$$t_{I \to A} = t_A - t_I$$

where $t_A$ is the ARKit clock. Finally, the required clock difference $t_{K \to A}$ is

$$t_{K \to A} = t_{K \to W} + t_{W \to I} + t_{I \to A}$$

### D.3. Point Clouds

Both Kinect and iPhone produce depth maps that can be converted to point clouds. However, iPhone's point cloud is much sparser than Kinect's. We show unprocessed raw point clouds produced by the two types of sensors in Fig. 11. In addition, iPhone does not report the LiDAR accuracy; we empirically find that iPhone point clouds are noisier, especially at the object boundaries, than Kinect point clouds.

## E. Additional Details of Toolchain

### E.1. Keypoint Annotation

The overall pipeline for keypoint annotation is summarized in Algorithm 1.

Table 6. A more complete comparison of HuMMan with published datasets. Subj: subjects; Act: actions; Seq: sequences; Video: sequential data, not limited to RGB sequences; Mobile: mobile device in the sensor suite; D/PC: depth image or point cloud, only genuine point cloud collected from depth sensors are considered; Act: action label; K2D: 2D keypoints; K3D: 3D keypoints; Param: statistical model (*e.g.* SMPL) parameters; Txtr: texture. -: not applicable or not reported

| Dataset | #Subj | #Act | #Seq | #Frame | Video | Mobile | Modalities | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | RGB | D/PC | Act | K2D | K3D | Param | Mesh | Txtr |
| Action Recognition | | | | | | | | | | | | | | |
| HMDB51 [49] | - | 51 | 7k | - | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| UCF101 [91] | - | 101 | 13k | - | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| Sports1M [42] | - | 487 | 1M | - | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| AVA [22] | - | 80 | 437 | - | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| Kinectics 700 [11] | - | 700 | 650k | - | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| HACS [117] | - | 200 | 1.55M | - | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| Moments-In-Time [71] | - | 339 | 1M | - | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| FineGym [88] | - | 530 | 32k | - | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| HAA500 [15] | - | 500 | 10k | 591k | ✓ | - | ✓ | - | ✓ | - | - | - | - | - |
| MSR-Action3D [53] | 10 | 20 | 567 | - | ✓ | - | - | ✓ | ✓ | - | ✓ | - | - | - |
| Northwestern-UCLA [101] | 10 | 10 | 1.47k | >23k | ✓ | - | ✓ | ✓ | ✓ | - | ✓ | - | - | - |
| SYSU 3DHOI [30] | 40 | 12 | 65 | - | ✓ | - | ✓ | ✓ | ✓ | - | ✓ | - | - | - |
| NTU RGB+D [87] | 40 | 60 | 56k | - | ✓ | - | ✓ | ✓ | ✓ | - | ✓ | - | - | - |
| NTU RGB+D 120 [58] | 106 | 120 | 114k | - | ✓ | - | ✓ | ✓ | ✓ | - | ✓ | - | - | - |
| NTU RGB+D X [97] | 106 | 120 | 113k | - | ✓ | - | ✓ | ✓ | ✓ | - | ✓ | ✓ | - | - |
| 2D/3D Keypoint Detection and 3D Parametric Human Recovery | | | | | | | | | | | | | | |
| J-HMDB [35] | - | 21 | 928 | 33.18k | ✓ | - | ✓ | - | ✓ | ✓ | - | - | - | - |
| Penn Action [114] | - | 15 | 2.32k | - | ✓ | - | ✓ | - | ✓ | ✓ | - | - | - | - |
| MPII [3] | - | 410 | - | 24k | - | - | ✓ | - | ✓ | ✓ | - | - | - | - |
| COCO [56] | - | - | - | 104k | - | - | ✓ | - | - | ✓ | - | - | - | - |
| PoseTrack [2] | - | - | >1.35k | >46k | ✓ | - | ✓ | - | - | ✓ | - | - | - | - |
| Human3.6M [32] | 11 | 17 | 839 | 3.6M | ✓ | - | ✓ | ✓ | ✓ | ✓ | ✓ | - | - | - |
| CMU Panoptic [38] | 8 | 5 | 65 | 154M | ✓ | - | ✓ | ✓ | - | ✓ | ✓ | - | - | - |
| MPI-INF-3DHP [68] | 8 | 8 | 16 | 1.3M | ✓ | - | ✓ | - | - | ✓ | ✓ | - | - | - |
| TotalCapture [98] | 5 | 5 | 60 | 1.89M | ✓ | - | ✓ | - | - | ✓ | ✓ | - | - | - |
| 3DPW [99] | 7 | - | 60 | 51k | ✓ | ✓ | ✓ | - | - | - | - | ✓ | - | - |
| AMASS [65] | 344 | - | >11k | >16.88M | ✓ | - | - | - | - | - | - | ✓ | ✓ | - |
| Mirrored-Human [16] | - | 56 | 56 | >1.5M | ✓ | - | - | - | ✓ | ✓ | ✓ | ✓ | - | - |
| AIST++ [52] | 30 | - | 1.40k | 10.1M | ✓ | - | ✓ | - | - | ✓ | ✓ | ✓ | - | - |
| Mesh Reconstruction | | | | | | | | | | | | | | |
| ZJU LightStage [82] | 6 | 6 | 9 | >1k | ✓ | - | ✓ | - | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| CAPE [63] | 15 | - | >600 | >140k | ✓ | - | - | - | ✓ | - | ✓ | ✓ | ✓ | - |
| BUFF [113] | 6 | 3 | >30 | >13.6k | ✓ | - | ✓ | ✓ | ✓ | - | ✓ | ✓ | ✓ | ✓ |
| DFAUST [6] | 10 | >10 | >100 | >40k | ✓ | - | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| People Snapshot [1] | 9 | - | 24 | 15k | ✓ | - | ✓ | - | - | - | ✓ | ✓ | ✓ | ✓ |
| LiveCap [25] | 7 | 11 | 11 | 36k | ✓ | - | ✓ | - | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| DynaCap [24] | 4 | 5 | 5 | 35k | ✓ | - | ✓ | - | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| DeepCap [26] | 4 | 17 | 17 | 26k | ✓ | ✓ | ✓ | - | ✓ | ✓ | ✓ | - | ✓ | ✓ |
| HUMBI [109] | 772 | - | - | ~26M | ✓ | - | ✓ | - | - | ✓ | ✓ | ✓ | ✓ | ✓ |
| THuman [119] | 200 | - | - | >6k | - | - | ✓ | ✓ | - | - | - | - | ✓ | ✓ |
| THuman2.0 [107] | 200 | - | - | >500 | - | - | - | - | - | - | - | ✓ | ✓ | ✓ |
| Multi-task | | | | | | | | | | | | | | |
| **HuMMan (ours)** | 1000 | 500 | 400k | 60M | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

## E.2. Full-body Angle Prior

It is surprisingly difficult to find literature that provides a complete analysis of joint movement ranges, especially rotations in three degrees of freedom (DOF). Hence, we take references from artists' guidelines on human anatomy [67] and 3D modelers' suggested practices [85], to simplify the constraint such that the three DOF movement range is bounded by the maximum ranges in each of the DOF. Despite that this formulation is not perfect, it provides con-

Figure 10. HuMMan deploys ten Azure Kinects and an iPhone 12 Pro Max for multi-view sequential data collection. We show several synchronized RGB frames captured with our hardware setup. The numbers are device IDs

straints that are otherwise completely absent. To easily apply the per-axis ranges, we convert the axis-angle representation into Euler angles and define the Z-axis to be aligned with the child bone of the joint in the kinematic tree (for example, *forearm* is the child bone of the joint *elbow*). We then define the X-axis as the axis around which the largest rotation is achieved. Y-axis is finally defined with X- and Z-axis fixed. All values undergo manual inspection and are adjusted empirically.

### E.3. Annotation Quality of SMPL Parameters.

To evaluate the body shape, we compute the per-vertex error on the high-resolution scan that is the uni-directional Chamfer distance from registered SMPL mesh vertices to the high-resolution scan vertices. Note that high-resolution scans have been scaled to the real height of scanned persons.

The mean per-vertex error is 0.16 mm. We also visualize the registration quality in Fig. 12. To evaluate the body pose, we compute the per-joint error as the L2 Euclidean distance between 3D keypoints and 3D joints of registered SMPL on the dynamic sequences. The mean per-joint error is 38.18 mm. Note that the error is largely attributed to the difference in the joint definition of the keypoint detector and the parametric model. As a reference, registration with an accurate optical marker system [32, 60] yields a per-joint error of 29.34 mm.

## F. Additional Details of Action Set

In HuMMan, we design a hierarchical structure for a systematic coverage of different body parts to collate a *complete* and *unambiguous* action set. Specifically, we have *body* at the center as the first order. The second order

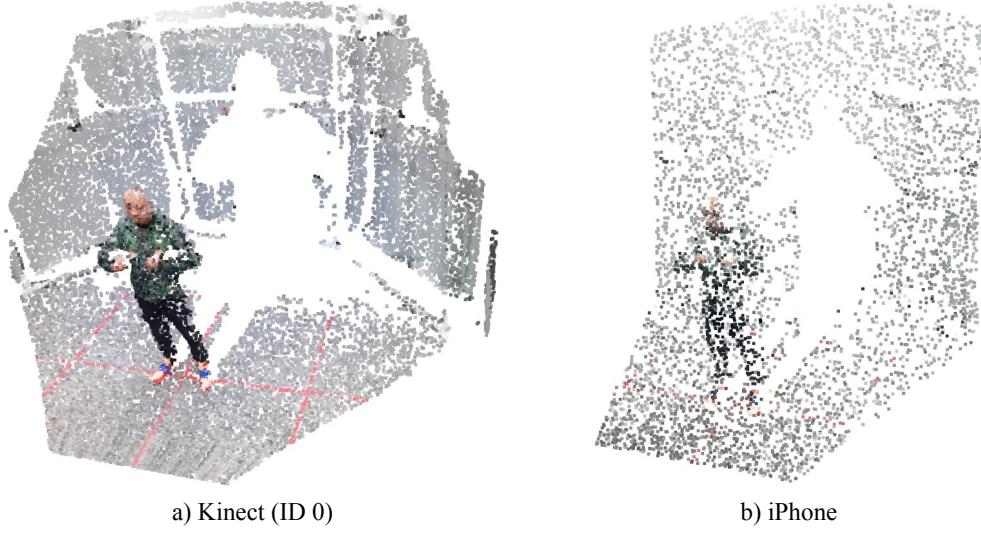a) Kinect (ID 0)                                     b) iPhone

Figure 11. The point clouds produced by the Kinect and the iPhone are different: the latter is significantly sparser. Note that the point clouds shown here are raw (not filtered or denoised). For visual comparison purpose, both point clouds are downsampled by the same factor of 10

---

**Algorithm 1** Keypoint Annotation

**Input:** Detected 2D Keypoints $\hat{\mathcal{P}}_{2D}$, camera parameters set $\mathcal{C}$, keypoint threshold $\tau_k$, reprojection minimal threshold $\tau_{min}$, reprojection maximum threshold $\tau_{max}$, camera threshold step $\Delta_c$, best camera number $N_c$.

**Output:** 3D Keypoints $\mathcal{P}_{3D}$, 2D Keypoints $\mathcal{P}_{2D}$

1: $\tau_c = \tau_{min}$, $\hat{\mathcal{C}} = \emptyset$
2: $\bar{\mathcal{P}}_{2D} = \text{FILTERKEYPOINTS}(\hat{\mathcal{P}}_{2D}, \tau_k)$
3: **while** $\tau_c \leq \tau_{max}$ **do**
4:     $\mathcal{P}_{3D} = \text{TRIANGULATE}(\bar{\mathcal{P}}_{2D}, \mathcal{C})$
5:     $\mathcal{P}_{2D} = \text{REPROJECTION}(\mathcal{P}_{3D})$
6:     **while** $\tau_c \leq \tau_{max}$ and $|\hat{\mathcal{C}}| < 3$ **do**
7:         $\hat{\mathcal{C}} = \text{SELECTCAM}(\mathcal{P}_{2D}, \bar{\mathcal{P}}_{2D}, \tau_c, N_c)$
8:         $\tau_c = \tau_c + \Delta_c$
9:     **end while**
10:    **if** $\mathcal{C} == \hat{\mathcal{C}}$ **then**
11:        **return** $\mathcal{P}_{3D}, \mathcal{P}_{2D}$
12:    **else**
13:        $\mathcal{C} = \hat{\mathcal{C}}$
14:    **end if**
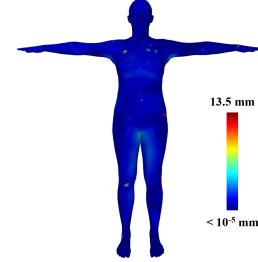15: **end while**
16: **return** Fail



Figure 12. The registration accuracy on high-resolution mesh (minimally clothed). The metric is mean uni-directional Chamfer distance (from SMPL vertices to high-resolution mesh vertices). Our registration (and subsequently the body shape obtained) is mostly accurate

consists of *whole body*, *upper extremity* and *lower limbs* that categorize actions by major body parts. After that, we propose a *muscle-driven* strategy to further split each major body part into main muscle groups according to human anatomy as the third order. Finally, we involve domain experts to design a serie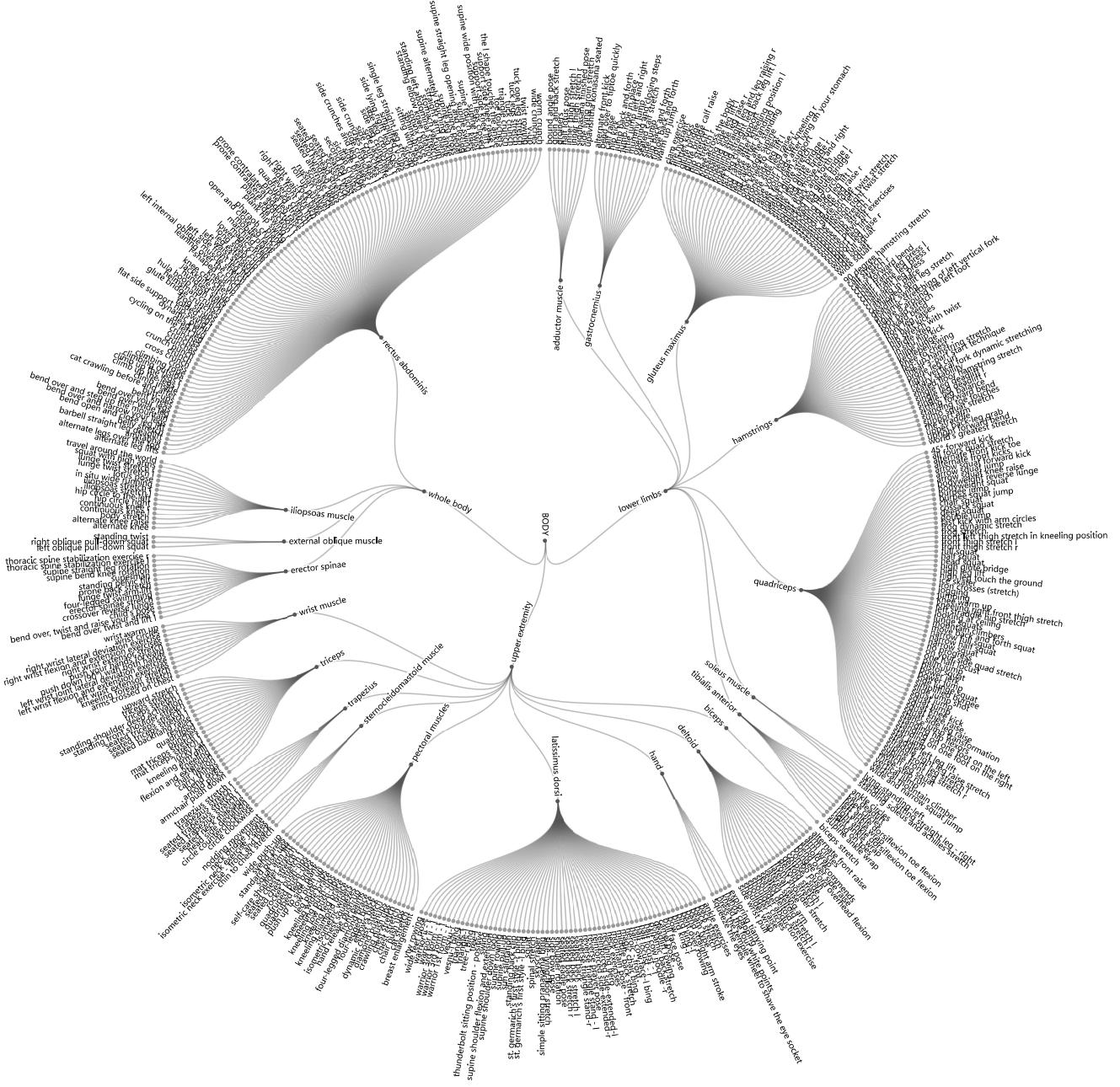s of action variants associated with each muscle in the fourth order. The full action hierarchy is demonstrated in Fig. 13.

## G. Additional Details of Subjects

HuMMan consists of 1000 subjects. To evaluate the diversity, we include key statistics (gender, age, height and weight) of the subjects in Fig. 14.

## H. Additional Details of Experiments

### H.1. Splits and Protocols

HuMMan contains a massive scale of subjects (1000), actions (500), sequences (400k) and frames (60M). To constrain training and testing within a reasonable computation budget, we sample only 10% of the data. We then develop

Figure 13. The complete set of 500 actions

three protocols to split iPhone and Kinect data into training and test sets. **Protocol 1 (P1)**: split by subjects, the training and test set are mutually exclusive and contain 70% and 30% of the subjects respectively. P1 is used for all experiments in the main paper. **Protocol 2 (P2)**: split by actions. We split actions into three categories according to major body parts involved: *upper extremity*, *lower limbs*, and *whole body*. Training is conducted on one category whereas the test is conducted on the other two. **Protocol**

**3 (P3)**: split by views. Model is trained on only one view (the *front* view, or the view of the iPhone and the Kinect with ID 0) and tested on all views.

## H.2. 2D Keypoint Detection

We study 2D keypoint detection baselines on HuMMan primarily for 2D-to-3D keypoint lifting. CPN [12] is a cascaded pyramid network to improve hard keypoints detection. HRNet [92] is a novel high-resolution network that ob-
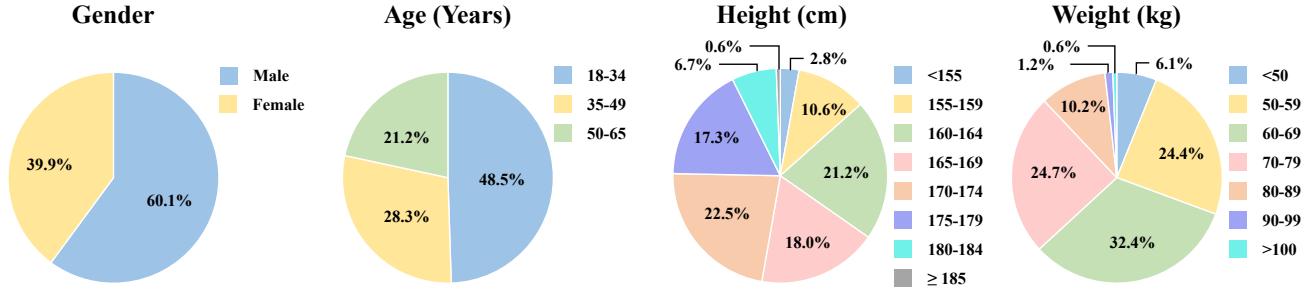
Figure 14. Statistics of HuMMan subjects

Table 7. 2D Keypoint Detection under Protocol 1. Input image is resized to 384×288

| Method | $AP^{50}$ ↑ | $AP^{75}$ ↑ |
|---|---|---|
| CPN [12] | 0.86 | 0.93 |
| HRNet [92] | 0.91 | 0.97 |
| Lite-HRNet [106] | 0.87 | 0.93 |

Table 8. 3D keypoint detection under Protocol 2 on Kinect splits. FCN is used as the base model.

| Training | Testing | MPJPE ↓ | PA-MPJPE ↓ |
|---|---|---|---|
| Lower Limbs | Upper Extremity | 70.3 | 55.7 |
| Lower Limbs | Whole Body | 97.5 | 72.3 |
| Upper Extremity | Lower Limbs | 75.8 | 55.1 |
| Upper Extremity | Whole Body | 99.6 | 72.5 |
| Whole Body | Lower Limbs | 77.4 | 56.2 |
| Whole Body | Upper Extremity | 86.2 | 66.4 |
| Mean Error | | 84.4 | 63.0 |

Table 9. 3D keypoint detection under Protocol 3 on Kinect splits. FCN is used as the base model. The model is trained on View 0 and tested on all views.

| View | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | Mean |
|---|---|---|---|---|---|---|---|---|---|---|---|
| MPJPE ↓ | 66.4 | 97.2 | 167.1 | 172.0 | 247.2 | 268.4 | 245.1 | 175.3 | 165.4 | 95.9 | 170.0 |
| PA-MPJPE ↓ | 41.2 | 67.5 | 100.9 | 103.5 | 112.3 | 118.7 | 111.8 | 103.9 | 100.2 | 67.1 | 92.7 |

tains high performance on COCO dataset [56], and LiteHR-Net is an efficient version of HRNet. The comparison results are listed in Table 7. Because 2D keypoints are often used as an intermediate representation of 3D keypoints in a two-stage manner [66, 80], the good performance in this task can be helpful to the estimation of subsequent 3D.

### H.3. 3D Keypoint Detection

3D keypoint detection benchmarks under P1 setting are presented in the main paper and additional benchmarks under P2 and P3 are provided here. In Table 8, we show results on the cross-action (P2) performance of the FCN

Table 10. 3D parametric human recovery under Protocol 2 on Kinect splits. HMR is used as the base model.

| Training | Testing | MPJPE ↓ | PA-MPJPE ↓ |
|---|---|---|---|
| Lower Limbs | Upper Extremity | 77.2 | 57.0 |
| Lower Limbs | Whole Body | 109.8 | 77.9 |
| Upper Extremity | Lower Limbs | 80.6 | 56.5 |
| Upper Extremity | Whole Body | 114.2 | 73.3 |
| Whole Body | Lower Limbs | 85.4 | 61.9 |
| Whole Body | Upper Extremity | 98.3 | 72.6 |
| Mean Error | | 94.2 | 66.5 |

method [66]. Compared with Protocol 1, we observe that training with fewer actions and testing on unseen actions degrade the precision significantly, especially for cross-evaluation on the *whole body* category which seems to have a large action distribution misalignment with the other two categories. Furthermore, deep learning models are sensitive to viewing angles [8, 64], we thus report results of cross-view (P3) in Table 9. When the model is only trained on one view (*i.e.*, View 0), we observe a considerable domain gap across different views as the errors increase as the deviation from the test view from the training view increases. The experiment results indicate that cross-view 3D keypoint detection is challenging.

### H.4. 3D Parametric Human Recovery

Table 11. 3D parametric human recovery under Protocol 3 on Kinect splits. HMR is used as the base model. The model is trained on View 0 and tested on all views.

| View | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | Mean |
|---|---|---|---|---|---|---|---|---|---|---|---|
| MPJPE ↓ | 61.9 | 122.9 | 223.9 | 206.2 | 343.9 | 421.0 | 334.0 | 208.0 | 199.0 | 123.5 | 224.4 |
| PA-MPJPE ↓ | 40.2 | 71.9 | 123.7 | 115.0 | 124.4 | 133.1 | 127.2 | 123.1 | 118.0 | 73.3 | 105.0 |

In addition to P1 benchmarks for 3D parametric human recovery presented in the main paper, we also provide more benchmarks under P2 and P3. In Table 10, we evaluate the cross-action (P2) performance of the HMR baseline. We

find that testing on unseen poses is challenging (compared to P1 benchmark results). Moreover, *whole body* actions seem to have a distribution that is further away from *lower limbs* and *upper extremity* actions. In Table 11, we study the cross-view setting (P3), which is even worse than the cross-action setting. The HMR baseline is trained on View 0, and gives a clear trend that the greater the viewing angle difference, the larger the errors. View 5 is directly opposite View 0 and yields the largest error.