
FreeMan: Towards Benchmarking 3D Human Pose Estimation in the Wild

Jiong Wang^{1,2*}, Fengyu Yang^{1*}, Wenbo Gou¹, Bingliang Li¹, Danqi Yan¹,
Ailing Zeng³, Yijun Gao², Junle Wang², Ruimao Zhang^{1†}

¹The Chinese University of Hong Kong, Shenzhen ²Tencent ³IDEA

Abstract

1 Estimating the 3D structure of the human body from natural scenes is a fundamental
2 aspect of visual perception. This task carries great importance for fields like
3 AIGC and human-robot interaction. In practice, 3D human pose estimation in
4 real-world settings is a critical initial step in solving this problem. However,
5 the current datasets, often collected under controlled laboratory conditions using
6 complex motion capture equipment and unvarying backgrounds, are insufficient.
7 The absence of real-world datasets is stalling the progress of this crucial task. To
8 address facilitate the development of 3D pose estimation, we present FreeMan,
9 the first large-scale, real-world multi-view dataset. FreeMan was captured by
10 synchronizing 8 smartphones across diverse scenarios. It comprises 11M frames
11 from 8000 sequences, viewed from different perspectives. These sequences cover
12 40 subjects across 10 different scenarios, each with varying lighting conditions. We
13 have also established an automated, precise labeling pipeline that allows for large-
14 scale processing efficiently. We provide comprehensive evaluation baselines for a
15 range of tasks, underlining the significant challenges posed by FreeMan. Further
16 evaluations of standard indoor/outdoor human sensing datasets reveal that FreeMan
17 offers robust representation transferability in real and complex scenes. FreeMan is
18 now publicly available at <https://github.com/wangjiongw/FreeMan>.

19

1 Introduction

20 Estimating 3D human poses from real scene input is a longstanding yet active research topic since
21 its huge potential in real applications, such as animation creation [1, 2], virtual reality [3, 4], the
22 metaverse [5–7] and human-robot interaction [8]. Specifically, it aims to identify and determine
23 the spatial positions and orientations of the human body’s parts in 3D space from input data such
24 as the image or the video. Despite numerous models proposed in recent years [9–11], practical
25 implementation in real scenes remains challenging due to the viewpoint variability, occasions,
26 human scale variation and complex background. Some challenges may stem from the disparity
27 between the recent benchmarks and real-world scenarios. As shown in Fig. 1, the widely recognized
28 Human3.6M [12], along with the currently largest dataset HuMMan [13], are usually collected in
29 laboratory settings utilizing intricate equipment, which maintains constant camera parameters and
30 offers minimal variation in background conditions. The effectiveness of the trained models when
31 trained using these datasets often experiences a significant decline in real-world environments.

32 From a data-oriented perspective, we have identified several constraints that hinder the performance
33 of the existing models. **(1) Insufficient Scene Diversity.** Existing datasets, as shown in Tab. 1, are

*The first two authors contributed equally to this work. Email: {jiongwang, fengyuyang}@link.cuhk.edu.cn

†The corresponding author is Ruimao Zhang. Email: zhangruimao@cuhk.edu.cn



Figure 1: The left displays sample frames from Human3.6M [12] and HuMMAN [13], which were collected under laboratory conditions, and contrasted with our FreeMan dataset that was collected in real-world scenarios. Frames from FreeMan have been cropped into a square format for visualization purposes, with the original resolution being 1920×1080 pixels. The right-hand side demonstrates the test results on 3DPW of the HMR model [16] when trained on these three datasets. Notably, the model trained using FreeMan is able to adapt flawlessly to real-world conditions, demonstrating its superior generalization ability.

Dataset	Environment	#Subj	#Action	#Scene	#Seq	#Frame	#Camera	FPS
HumanEva[17]	Laboratory	4	6	1	168	80K	7	30
CMU Panoptic[18]	Laboratory	8	5	1	65	154M	31	30
MPI-INF-3DHP[19]	Real Scene	8	8	1	16	1.3M	14	30
MuCo-3DHP[14]	Real Scene	7	-	4	60	51K	1	30
3DPW[15]	Real Scene	7	47	4	60	51K	1	30
Mirrored Human[20]	Laboratory	-	-	-	-	1.5M	1	30
Human3.6M[12]	Laboratory	9	15	1	840	3.6M	4 (Fixed)	30
AIST+ [21]	Laboratory	30	10	1	1408	10.1M	9 (Fixed)	30
HuMMAN[13]	Laboratory	1000	500	1	400K	60M [†]	11 (Fixed)	30
HuMMAN-released[13]	Laboratory	132	20	1	4466	278K [†]	11 (Fixed)	30
FreeMan	Real Scene	40	123	10 [‡]	8000	11.3M	8 (Movable)	30 / 60

Table 1: Comparison of our proposed FreeMan dataset with existing 3D Human Pose datasets. Only HD Cameras counted for FreeMan. [†] Only 1% of the HuMMAN dataset (600K frames) is made publicly available. [‡] FreeMan includes 10 type of scenes that corresponds to 29 locations. *Fixed* means cameras are fixed within the whole dataset, while our cameras are *movable* and camera poses vary among video sequences.



Figure 2: Equipment setting of data collection using 8 cameras.

mainly collected in controlled laboratory conditions, which may not be optimal for robust model training due to static lighting conditions and uniform backgrounds. This limitation becomes especially crucial when the objective is to estimate 3D pose in real-world scenarios, where scene contexts exhibit substantial variability. In certain datasets, even though the data is collected from outdoor scenes, *e.g.*, MuCo [14] and 3DPW [15] in Tab. 1, the variety of scenarios remains remarkably limited. This constraint significantly hampers the applicability of trained models across a broader range of situations. **(2) Limited Actions and Body Scales.** In existing datasets, the range of human actions tends to be rather limited. Even in the currently largest dataset, HuMMAN [13], the variety of actions in the publicly available data is quite restricted. Additionally, these large datasets typically employ fixed cameras to capture data from various perspectives. The distance from the camera to the actor is relatively constant, which results in a relatively fixed human body scale across different videos. **(3) Restricted Scalability.** The annotation of current datasets primarily relies on expensive manual processing, which greatly restricts the scalability of the datasets. Especially when the camera used for collection is movable, how to effectively align data from different cameras and perform efficient annotation remains an open issue.

To address these above issues, this work presents FreeMan, a novel large-scale benchmark for 3D human pose estimation in the wild. FreeMan contains 11M frames in 8000 sequences captured by 8 smartphone cameras from different views simultaneously, as illustrated in Fig. 2. It covers 40 subjects in 10 kinds of scenes. To our best knowledge, it is the current largest multi-view 3D pose estimation dataset, with variable camera parameters and complex background environments. It is 215 \times of the famous outdoor dataset 3DPW [15]. From a practical perspective, it has several appealing strengths: **Firstly**, a large number of scenes introduce diversity in both backgrounds and lighting, enhancing the generalizability of models trained on FreeMan in real-world scenarios. This makes it particularly suitable for evaluating algorithmic performance in practical applications. **Secondly**, the distances between the 8 cameras and the actors are variant (*i.e.*, 2 to 5.5 meters) both within and among subjects, resulting in significant scale changes in human bodies across different sequences. **Thirdly**, although we employed mobile data collection devices, the annotation of the FreeMan dataset does not rely on

61 costly manual processes, thereby significantly enhancing the scalability of the dataset. **Lastly**, the
62 proposed FreeMan encompasses a wide range of pose estimation tasks, which include monocular
63 3D estimation, 2D-to-3D lifting, multi-view 3D estimation, and neural rendering of human subjects.
64 We present thorough evaluation baselines for the aforementioned tasks on FreeMan, highlighting the
65 inherent challenges of such a new benchmark.

66 In summary, this paper has made three contributions: (1) We have constructed a large-scale dataset
67 for 3D human pose estimation in uncontrolled environments. The impressive transferability of the
68 models trained on this dataset to real-world scenarios has been demonstrated. (2) We have showcased
69 a simple yet effective toolchain that enables the automatic generation of precise 3D annotations
70 from the collected data. (3) We provide comprehensive benchmarks for human pose estimation and
71 modeling on FreeMan, facilitating downstream applications. These baselines highlight potential
72 directions for future algorithmic enhancements.

73 2 Related Work

74 **Human Pose Datasets.** Human modeling is a significant task in computer vision. Existing datasets
75 predominantly rely on 2D and 3D keypoint annotations, with 3D keypoint datasets available in two
76 forms: monocular and multi-view. For 2D keypoint, there are some single-frame datasets such as
77 MPII [22] and COCO [23], which provide diverse images with 2D keypoints annotations, while
78 video datasets such as J-HMDB [24], Penn Action [25] and PoseTrack [26] provide 2D keypoints
79 with temporal information. In contrast, 3D keypoint datasets are often constructed in indoor scenes,
80 such as Human3.6M [12], CMU Panoptic [27], MPI-INF-3DHP [19], AIST++ [28] and HuMMan
81 [13] for multi-view. There also exists some outdoor datasets such as 3DPW [15] for monocular cases.
82 Details of these datasets are shown in Tab. 1. However, the majority of outdoor datasets such as
83 MPI-INF-3DHP, MuCo-3DHP, and 3DPW exhibit a limited variety of acquisition scenes, and the
84 datasets that involve fixed camera poses such as AIST++.

85 **3D Human Pose Estimation.** The present study categorizes the task of 3D pose estimation into three
86 distinct types, namely 2D-to-3D pose lifting, monocular 3D pose estimation, and multi-view 3D pose
87 estimation. In the 2D-to-3D pose lifting task, Martinez [29] proposed a simple baseline to regress
88 the 3d keypoints based on a convolutional neural network from 2D keypoints. While Pavllo[30],
89 Zheng[31] and Li[10]’s work better incorporated temporal information in this task. In monocular 3D
90 pose estimation task, HMR[16], SPIN[32] takes a single RGB image as input to perform 3D huna pose
91 estimation, which is often used as baselines for comparison with other algorithms, such as PARE[9],
92 SPEC[33] and HybrIK[34]. Additionally, multi-view methods are proposed to accommodate potential
93 body parts overlapping in monocular view. Iqbal’s [35] and MCSS [36] adopt weak supervision to
94 reduce the dependence on the 3D annotated pose, while Wandt et al. [37] and Kocabas et al. [38]
95 turned to self-supervise fashion to deal with multi-view data.

96 **Neural Rendering of Human Subjects.** With the development of NeRF[39] in dynamic scene
97 rendering, people also focus on the dynamic rendering of humans. Compared to dynamic scenes, the
98 non-rigid property of humans has more challenges. The prior knowledge of body movements can
99 provide a good prior for rendering, and many methods use SMPL[40] as a prior for body rendering.
100 Most methods reconstruct human bodies through multi-view videos[41–43], while recent works have
101 also employed single-view videos, such as HumanNeRF[44], FlexNeRF[45], YOTO[46].

102 3 FreeMan Dataset

103 FreeMan is a large-scale multi-view dataset in the wild, offering precise 3D pose annotations. It
104 comprises 11M frames from 1000 sessions, featuring 40 subjects across 10 types of scenes. The
105 dataset includes 10M frames recorded at 30FPS and an additional 1M frames at 60FPS. Next, we
106 highlight the diversity of FreeMan, from various camera settings and scenario selections

107 **Scenarios.** We design 10 types of real-world scenes, comprising 4 indoor and 6 outdoor scenes,
108 for our data collection. Fig. 3 illustrates the scene diversity of our FreeMan. The blue section

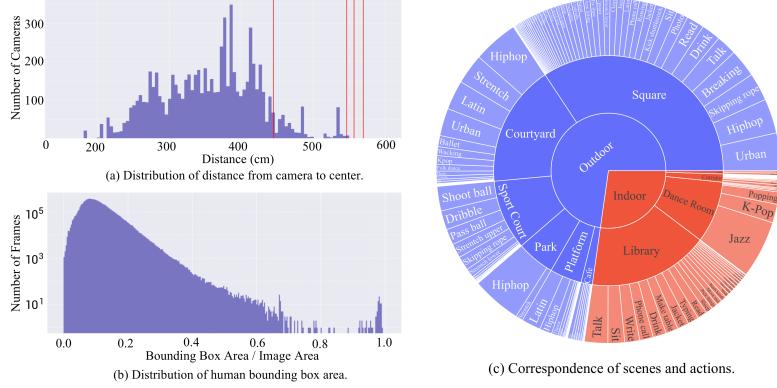


Figure 3: (a) Distribution of distance from the camera to the center of the system, indicated by translation along the z-axis in camera parameters. Four vertical red lines represent the distance of 4 cameras in Human3.6M [12]. (b) Distribution of human bounding box areas. The horizontal axis represents the ratio of the bounding box area over the image area. The vertical axis is in log scale. (c) Correspondence of scenes and actions. Areas of blocks represent the scale of the respective frame number. The outmost circle shows actions and the circle in the middle present 10 type of scenes in our dataset. Zoom in 10× for the best view.

represents the data collected outdoors, while the red sector refers to frames captured in indoor scenes. Specifically, there are 2.76 million frames captured indoors and 8.45 million frames captured outdoors. Additionally, there are different frame numbers collected under varying lighting conditions, with 1 million frames captured at night and 7.45 million frames captured during the day outdoors. Moreover, the central block of the circle denotes different scenarios, while the blocks on the outermost circle refer to actions. The areas of the blocks are proportional to the corresponding frame numbers. For more details about the distribution of scenarios, please refer to the supplementary material.

Action Set. Following the popular action recognition dataset NTU-RGBD120[47], we compose our action set with several common actions corresponding to scenes in daily life, *e.g.*, drinking and talking in a cafe and reading in the library. Also, subjects interact with real objects to make actions as close to real life as possible. As shown in the topmost row of Fig. 4, interaction with objects brings complicated occlusions, making our data more challenging. For outdoor scenarios, we set the data collection field as large as possible to help subjects perform actions with little restriction.

Camera Poses. Cameras in previous 3D human pose datasets [13, 12, 18] are fixed during data collection, resulting that only a few camera poses being included. As shown in Fig. 2, Our cameras are attached to lightweight tripods and are newly placed from time to time, and translation from the center of the system to camera d , which is the physical distance between the camera and the system center, can vary from 2m to 5.5m. Fig. 3 (a) shows the distribution of d and the corresponding number of cameras. Most cameras are located around 4 meters far away from the system center. Besides, we show the distribution of the human bounding box area in Fig. 3 (b), in a unit of ratio to the whole image area, to demonstrate the variation of human size. With variations in camera translation and human actions, the area of human bounding boxes varies from 0.01 to 0.7 of the whole image area.

Subjects. There are 40 subjects participating in the construction of FreeMan and recruitment is completely based on voluntary. All of them are well-informed and signed the agreement to make the data public for research purposes only.

4 Data Acquisition & Annotation Pipeline

Overview. To collect a large-scale dataset from real-world environments, we developed a comprehensive toolchain, as shown in Fig. 5. Unlike previous toolchains used in controlled or idealized conditions, we carefully accounted for potential challenges in outdoor settings, including calibration and synchronization errors. To overcome these issues, we incorporated automated correction procedures to ensure efficient data collection.



Figure 4: The diverse frames in FreeMan. The topmost two rows illustrate a range of indoor and outdoor scenes, highlighting the diversity of scene contexts, lighting conditions, and subjects. The third row exhibits frames from different cameras, showcasing the variance across views. The final row illustrates the temporal variation of human poses from a consistent viewpoint, emphasizing the dynamism of motion capture.

140 4.1 Hardware Setup

141 **Cameras.** We collect FreeMan via 8 Mi11 phones [48] indexed from 1 to 8 as our data collection
 142 devices. **Note 8 collection of one action as one session, which corresponds to RGB sequences from**
 143 **8 views.** and each phone is attached to a tripod to keep stable within each data collection session.

144 As shown in Fig. 2, all devices are positioned in a circle around a human at a height of approximately
 145 1.6 meters above the ground, with the distance from the camera to the system center varying from 2 to
 146 5.5 meters. which is similar to real-life usage scenarios. Each smartphone captures RGB sequences
 147 using the main camera at 1920×1080 resolution and 30/60 FPS. During the data collection process,
 148 actors perform actions facing the device with odd-numbered indices. As shown in Fig. 5(a), the only
 149 requirement beyond devices is to connect all devices to the same server stably.

150 **Camera Calibration.** At the beginning of each session, we calculate the intrinsic and extrinsic camera
 151 parameters in two steps with a chessboard tiled at the center of the system, following the standard
 152 implementation in OpenCV[49, 50]. We attach cameras to tripods and collect tiled chessboards from
 153 each view. In this phase, chessboard frames from all views are sent to the server and the server detects
 154 corner points, calculates camera extrinsic parameters of each device, and sends parameters back to
 155 cameras. Cameras are allowed for the next step only if extrinsic parameters are received.

156 **Device Synchronization.** Previous works [13, 12, 18] have synchronized devices using wired
 157 interfaces in a laboratory. However, the complexity of the entire system coupled with the difficulty
 158 in deploying it in real-world environments, has prompted us to consider alternative methods. To
 159 address issues related to usability and device constraints, we connect all devices wirelessly to a single
 160 server and developed an Android app that utilizes the Network Time Protocol (NTP) [51] to calculate
 161 the time difference between each device and the server’s clock. During the capture process, time
 162 information is stored locally on each device as a timecode, while the server records the synchronized
 163 capture interval for all devices. The starting frame is determined by matching the timecode to
 164 the frame closest to the server’s clock time. This approach ensures device synchronization, with
 165 synchronization errors limited to smaller than a single frame during our testing. At frame rates of
 166 30FPS and 60FPS, this corresponds to 33ms and 16ms, respectively.

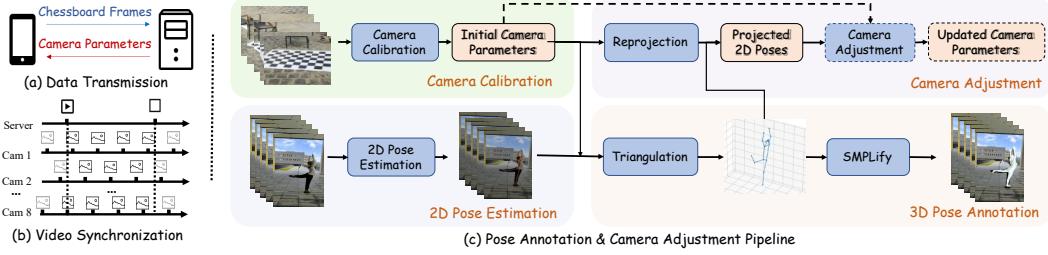


Figure 5: The illustration of data collection and annotation toolchain: (a) depicts the transmission of signals between cameras and servers for camera calibration, where chessboard frames are sent to the server, and camera parameters are returned. (b) demonstrates the synchronization process among devices. (c) showcases the pipeline for annotation and camera adjustment.



Figure 6: The examples of human pose annotations are presented as follows. In the left part, the first row displays 2D keypoints directly generated by HRNet-w48 [53], while the second row illustrated the re-projected 2D poses. When HRNet fails to deal with occlusions, our toolchain automatically corrects annotations by incorporating multi-view priors. The right part showcases the SMPL annotation examples for each view in our dataset.

167 4.2 Pose Annotation.

168 After camera calibration, subjects enter the capture area to perform actions. Once videos are collected,
 169 we use a state-of-the-art detector YOLOX[52] to detect human locations and HRNet-w48[53] to detect
 170 2D keypoints of 8 views $K_{2D} \in \mathbb{R}^{8 \times 17 \times 2}$ in COCO[23] format. To eliminate the effect of potential
 171 wrong keypoints output, keypoint predictions with confidence lower than a threshold ϕ are removed.
 172 Then remaining 2D keypoints are used for triangulation to get 3D human pose $K_{3D} \in \mathbb{R}^{17 \times 3}$ with
 173 pre-computed camera parameters. Here, we set ϕ to be 0.5. Furthermore, we optimize K_{3D} with
 174 smoothness constraints and bone length constraints introduced in HuMMAN[13] resulting in optimized
 175 3D pose $\tilde{K}_{3D} \in \mathbb{R}^{17 \times 3}$. Then we fit a standard SMPL[40] model to the estimated 3D skeleton by
 176 SMPLify [54] to produce a rough mesh annotation. After that, we project 3D keypoints to 2D image
 177 planes of each view using corresponding camera parameters. With regularization in triangulation and
 178 optimization along the temporal axis, the re-projected 2D poses \tilde{K}_{2D} is a refined version compared
 179 with K_{2D} . Comparison between original K_{2D} and \tilde{K}_{2D} are shown in the left part of Fig. 6.

180 **Camera Adjustment.** Finally, for cameras that failed in calibration accidentally, we adjust camera
 181 parameters through pose computation[55] if needed. If more than two cameras failed in calibration,
 182 the session is discarded directly. Otherwise, we ignore the uncalibrated views first, then triangulate
 183 the remaining 2D poses into 3D poses. We can obtain the camera parameters of uncalibrated cameras
 184 by minimizing the re-projection error of 2D-3D point correspondences.

185 4.3 Annotation Quality Assessment

186 To demonstrate the effectiveness of our toolchain, we test the toolchain on Human3.6M[12]. We
 187 select 3 different actions of each subject in the training set, which covers 10% sequences of the whole
 188 training set and all kinds of actions. Performance is assessed by Euclidean distance between estimated
 189 2D poses and ground truth 2D poses in units of pixels. The error results in 9.9 pixels for images with
 190 the size of 1000×1000 , demonstrating that our toolchain can generate accurate annotations.

191 To validate the effectiveness of camera adjustment for failure cases in calibration, we randomly mask
 192 a camera for all sessions and use the other 3 views to get 3D poses, then compute the parameters
 193 of the missing one via pose computation. Then 3D poses are projected to the corrupted view by
 194 estimated camera parameters, mean error between resulting 2D poses and ground truth is 7.42 pixels.

195 5 Benchmarks

196 We have constructed four benchmarks utilizing images and annotations derived from our dataset.
 197 The data is subdivided based on subjects, allocating 18 subjects for training, 7 for validation, and
 198 15 for testing purposes. This partitioning results in three subsets composed of $5.87M$, $700K$, and
 199 $3.69M$ frames, respectively. For each benchmark, subject lists of each subset are shared, and only
 200 views selected from the session vary for each task. **(1) Monocular 3D Human Pose Estimation**
 201 (**HPE**). This task involves taking an RGB image or sequence from a monocular view as input and
 202 predicting 3D coordinates relative to the camera. For this benchmark, we randomly select one view
 203 from each session to construct the dataset. The performance of algorithms is measured by widely
 204 used Mean Per Joint Position Error (MPJPE)[12] and Procrustes analysis MPJPE (PA-MPJPE)[40].
 205 **(2) 2D-to-3D Lifting.** Given that 2D human poses can be predicted using existing 2D keypoint
 206 detectors [56–58, 53], the primary goal of this task is to effectively elevate these 2D poses into the 3D
 207 space within the camera coordinate system. The evaluation metrics are the same as HPE. **(3) Multi-**
 208 **View 3D Human Pose Estimation.** Due to occlusions in motion capture, monocular methods often
 209 encounter difficulties. Estimating the 3D human pose from multiple views presents a natural solution
 210 to overcome this problem. For this task, models are provided with images or videos from multiple
 211 views, along with corresponding camera parameters. The objective is to predict the 3D coordinates of
 212 human joints in the same world coordinate system as the cameras. The performance is measured by
 213 MPJPE and average precision (AP) following previous work[59]. **(4) Neural Rendering of Human**
 214 **Subjects.** The free-viewpoint rendering of humans is a significant issue in human modeling. With
 215 the rise in popularity of neural radiance fields (NeRF) [39] for the novel view rendering task, several
 216 methods, including HumanNeRF [44], have emerged. These methods utilize monocular human
 217 motion videos as input to synthesize novel views of dynamic humans through NeRF. The widely used
 218 metrics of prediction are PSNR, SSIM[60] and LPIPS[61].

219 6 Experiments

220 We experiment with the four benchmarks. In human 3D pose estimation tasks, we conduct several
 221 transfer tests with other standard datasets to evaluate the effectiveness and transferability of our
 222 proposed FreeMan dataset. Existing similar datasets, Human3.6M[12] & HuMMAn[13], are used for
 223 comparison. Since *HuMMAn* only releases 1% of data, we only involve it in monocular 3D human
 224 pose estimation and 2D-to-3D lifting. In the neural rendering of human subjects tasks, we train the
 225 model from one of the 8 views and test from the rest of the views in selected sessions.

226 6.1 Monocular 3D Human Pose Estimation

Method	HMR			PARE		
	Train	Supervision	Test	MPJPE	PA	MPJPE
Human3.6M	2D+3D KPTs	3DPW	279.92	133.13	118.54	81.22
HuMMAn	2D+3D KPTs	3DPW	407.57	192.75	110.99	63.11
HuMMAn	2D KPTs+SMPL	3DPW	475.73	184.15	114.20	66.19
HuMMAn	2D+3D KPTs+SMPL	3DPW	437.52	203.17	114.33	72.12
FreeMan	2D+3D KPTs	3DPW	157.46	87.93 ^{↑33.95%}	118.31	68.72 ^{↑15.39%}
FreeMan	2D KPTs+SMPL	3DPW	151.85	88.85 ^{↑51.75%}	94.27	60.39 ^{↑8.76%}
FreeMan	2D+3D KPTs+SMPL	3DPW	159.31	91.33 ^{↑55.04%}	98.33	64.51 ^{↑10.55%}

Table 2: Monocular 3D HPE performance of HMR[16] and PARE[9] trained on different dataset for monocular Human Pose Estimation. PA stands for PA-MPJPE and both metrics are in unit of millimeters. The lower metrics is, the better performance model obtains. All released part of HuMMAn is used for training. [↑] refers to the improvement relative to HuMMAn and [↑] refers to the improvent relative to Human3.6M.

227 **Implementation details.** For the Human3.6M [12] and HuMMAn [13] datasets, all views in their
 228 training set are utilized. In contrast, for FreeMan, we randomly sample a single view from sessions

Train	Test	AP@25mm (%) ↑	AP@50mm (%) ↑	AP@75mm (%) ↑	AP@100mm (%) ↑	Recall@500mm (%) ↑	MPJPE@500mm (mm) ↓
Human3.6M	Human3.6M	32.32	97.47	98.61	98.99	100.00	25.95
Human3.6M	FreeMan	0.00	0.00	0.00	0.00	0.06	89.85
Human3.6M	FreeMan (w/ GT Root)	0.00	1.27	11.44	21.40	96.20	154.41
FreeMan	FreeMan	43.38	88.77	97.73	99.12	99.97	26.61
FreeMan	Human3.6M	0.00	5.77	82.85	92.62	96.68	62.37
FreeMan	Human3.6M (w/ GT Root)	0.00	6.60	87.91	95.38	100.00	58.30

Table 3: Multi-View 3D Pose Estimation results of VoxelPose[59]. Ground truth root position (GT Root) is not used if not specified. Recall@500mm shows the percentage that falls within the threshold, and the MPJPE@500mm indicates the average MPJPE values within the threshold. Rows highlighted shows the best setting in cross-domain test.

in the training split, resulting that the frame numbers of all three datasets being 312K, 253K, and 233K, respectively. To enhance efficiency, videos of all three datasets are downsampled to 10FPS, following the implementation of MPMpose [62]. We trained HMR[16] and PARE[9] models on different datasets using configurations open-sourced by [63]. Please refer to Appendix for more.

Results. We perform testing on the test set of 3DPW [15]. The performance of the models trained on different datasets, with varying types of supervision, are reported in Tab. 2. Notably, the HMR models trained on FreeMan exhibit significantly better performance on the 3DPW test set compared to those trained on Human3.6M and HuMMAN. The PA-MPJPE scores are 42.59mm and 112.63mm, respectively, indicating that our dataset demonstrates superior generalizability compared to the others. The same results are obtained with PARE, further confirming that our FreeMan outperforms even in more advanced algorithms. This can be attributed to the diversity of scene contexts and human actions present in our dataset, which provides better transferability in real-world scenarios.

6.2 Multi-View 3D Human Pose Estimation

Implementation Details. We conduct in-domain and cross-domain tests between Human3.6M and FreeMan to evaluate the effectiveness and generalization ability. We conduct the experiments of multi-view 3D human pose estimation with VoxelPose[64], which locates the human root first and then regresses 3D joint location accordingly. COCO-format poses in FreeMan are transformed to match that in Human3.6M via interpolation. We trained VoxelPose[64] on the two datasets following official implementation. For Human3.6M, bounding box annotations are from [65] and its validation set is used for the test. For FreeMan, we only use videos of odd-indexed views from the training set.

Results. Results of all experiments are reported in Tab. 3. For in-domain testing, the model trained on FreeMan achieves MPJPE@500mm of 26.61mm on test set consisting of *odd-indexed* views. For cross-domain testing, the model trained on FreeMan achieves Recall@500mm of 96.68% while MPJPE@500mm is 62.37mm on Human3.6M validation set. However, the model trained on the Human3.6M dataset fails to locate human on FreeMan test set. To get rid of the effects of root location, we input the ground truth root locations to model directly. With this setting, the model trained on Human3.6M obtains MPJPE@500mm of 154.41mm on FreeMan test set, while the model trained on FreeMan can obtain MPJPE@500mm of 58.30mm on Human3.6M validation set. Results show that the model trained on FreeMan has a much better generalization ability, while that on Human3.6M struggles in transfer testing.

6.3 2D-to-3D Pose Lifting

Implementation Details. We employ four methods, either CNN-based or Transformer-based, including SimpleBaseline[29], VideoPose3D[30], PoseFormer[31] and MHFormer[10] in this task, and all the methods follow the corresponding official implementation. The results of VideoPose3D and PoseFormer can be found in Appendix. For training set in this task, we select one view from every session and down-sample the videos to 15FPS, resulting in the frame number to be 350K, which is similar to the amount of released part of HuMMAN (253K) and much smaller than Human3.6M (1500K). In order to verify the generalization across datasets, we unify the test set to be AIST++[28]. To verify the effect of the dataset scale, we also train our model on the whole training set.

Results. The experimental results are shown in Tab. 4. Results of the in-domain test on FreeMan are provided as a baseline for future work. For in-domain testing, the MPJPE of SimpleBaseline trained on FreeMan, 79.22mm, is larger than that on HuMMAN[13] (78.5mm) and Human3.6M[12]

271 (53.4mm), demonstrating that FreeMan is a more challenging benchmark. Besides, all the methods
 272 trained on FreeMan tend to generalize better than that on HuMMan and Human3.6M when testing on
 273 AIST++ under the same setting. Although the scale of FreeMan training set is of a similar magnitude
 274 as HuMMan’s, which is much smaller than Human3.6M’s, models trained on FreeMan outperform
 275 models trained on the other two by a large margin. Furthermore, when the training set is expanded to
 276 all frames in training split, FreeMan can further boost models to achieve better performance, proving
 277 that our large-scale data helps to improve model performance.

278 6.4 Neural Rendering of Human Subjects.

279 **Implementation Details.** We employ 10 scenes captured by FreeMan to train HumanNeRF[44], a
 280 deep neural network that aims to achieve high-quality human-centric new synthesis. To obtain human
 281 body segmentation annotations, we utilize the SAM (Segment Anything)[66] algorithm using our
 282 bounding boxes as prompts. Throughout the training step, we randomly select one view for each
 283 session and render the rest 7 view as novel views for testing. We then calculate metrics including
 284 PSNR, SSIM, and LPIPS, to evaluate the performance of the model. Please refer to Appendix for
 285 results of data at 60FPS.

286 **Results.** The reconstruction results in 10 scenes are shown in Tab. 5. The best reconstruction achieves
 287 a high PSNR of 30.11 which indicates FreeMan contains contents that the model can learn and fit
 288 very well. While the performance varies, the lowest PSNR of 23.86 shows FreeMan also contains
 289 contents that are outside of model’s learning scope and challenging. Additionally, the results in 10
 290 scenes including both familiar contents that the model can handle well and more challenging new
 291 contents demonstrates the diversity of our dataset.

Algorithm	Train	Test	MPJPE (mm)	PA (mm)
SimpleBaseline	FreeMan	FreeMan	90.53	54.17
	FreeMan [†]	FreeMan	79.22	49.11
	Human3.6M	AIST++	212.57	138.98
	HuMMan	AIST++	255.5	116.86
	FreeMan	AIST++	156.96	105.85 ^{+10.30%}
	FreeMan [†]	AIST++	126.23	88.07 ^{+24.64%}
MHFormer	FreeMan	FreeMan	93.00	63.50
	FreeMan [†]	FreeMan	77.06	53.38
	Human3.6M	AIST++	171.19	133.37
	HuMMan	AIST++	188.73	101.52
	FreeMan	AIST++	132.99	88.79 ^{+12.54%}
	FreeMan [†]	AIST++	124.34	79.22 ^{+21.97%}

Table 4: Performance of methods with different training and testing datasets in 2D-to-3D Pose Lifting. PA stands for PA-MPJPE. [†] refer to experiments with the whole training set of FreeMan. Smaller MPJPE and PA-MPJPE indicate better performance. Highlighted rows show training on our dataset achieves the best performance in the transfer test. ↑ refers to the improvement relative to HuMMan.

Scene	PSNR↑	SSIM↑	LPIPS*↓
Square	25.98	0.9501	58.38
Corridor	24.57	0.9340	81.39
Sports Port	26.33	0.9662	30.09
Park	23.86	0.9439	73.61
Courtyard	28.56	0.9630	53.99
Dance Room	30.11	0.9658	43.34
Library	29.41	0.9665	31.53
Platform	26.79	0.9439	70.01
Lobby	25.41	0.9387	78.80
Cafe	27.32	0.9644	37.88

Table 5: Neural rendering results by using HumanNeRF[44] on 10 scenes of FreeMan. Note that $LPIPS^* = LPIPS \times 10^3$.

292 7 Conclusion

293 We present FreeMan, a novel large-scale multi-view 3D pose estimation dataset with comprehensive
 294 tasks and 3D human pose annotations. We elaborately develop a simple yet effective annotation
 295 pipeline to automatically annotate frame-level 3D landmarks and precise 3D human motions at a
 296 much lower cost. We establish benchmarks covering multiple tasks in human modeling, including
 297 monocular and multi-view 3D human pose estimation, 2D-to-3D pose lifting, and neural rendering of
 298 human subjects. Extensive experimental results demonstrate the strengths of the proposed FreeMan.

299 **Limitation and future work.** While we have validated the superiority of our dataset on numerous
 300 in-domain and cross-domain cases, there is still a wide range of potential applications that are worth
 301 exploring. Additionally, since the data is automatically annotated, more evaluation schemes for the
 302 accuracy of annotations and understanding their potential impact on downstream applications remain
 303 open issues. We consider these aspects as future work for further investigation. As a large-scale
 304 human motion dataset, our FreeMan addresses the existing gap between the current datasets and
 305 real-scene applications, and we are optimistic that it will catalyze the development of algorithms
 306 designed to model and sense human behavior in real-world scenes.

307 **References**

- 308 [1] Chung-Yi Weng, Brian Curless, Pratul P. Srinivasan, Jonathan T. Barron, and Ira Kemelmacher-
309 Shlizerman. Humannerf: Free-viewpoint rendering of moving people from monocular video. In
310 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
311 pages 16210–16220, June 2022.
- 312 [2] Jae Shin Yoon, Lingjie Liu, Vladislav Golyanik, Kripasindhu Sarkar, Hyun Soo Park, and Chris-
313 tian Theobalt. Pose-guided human animation from a single image in the wild. In *Proceedings of*
314 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15039–15048,
315 2021.
- 316 [3] Thomas Waltemate, Dominik Gall, Daniel Roth, Mario Botsch, and Marc Erich Latoschik.
317 The impact of avatar personalization and immersion on virtual body ownership, presence, and
318 emotional response. *IEEE Transactions on Visualization and Computer Graphics*, 24(4):1643–
319 1652, 2018.
- 320 [4] Shivam Grover, Kshitij Sidana, and Vanita Jain. Pipeline for 3d reconstruction of the human
321 body from ar/vr headset mounted egocentric cameras. *arXiv preprint arXiv:2111.05409*, 2021.
- 322 [5] Lik-Hang Lee, Tristan Braud, Pengyuan Zhou, Lin Wang, Dianlei Xu, Zijun Lin, Abhishek
323 Kumar, Carlos Bermejo, and Pan Hui. All one needs to know about metaverse: A complete
324 survey on technological singularity, virtual ecosystem, and research agenda. *arXiv preprint*
325 *arXiv:2110.05352*, 2021.
- 326 [6] Jae Shin Yoon. Metaverse in the wild: Modeling, adapting, and rendering of 3d human avatars
327 from a single camera. 2022.
- 328 [7] Doina Popescu Ljungholm. Metaverse-based 3d visual modeling, virtual reality training experi-
329 ences, and wearable biological measuring devices in immersive workplaces. *Psychosociological*
330 *Issues in Human Resource Management*, 10(1), 2022.
- 331 [8] Abdelfetah Bentout, Mustapha Aouache, Abderraouf Maoudj, and Isma Akli. Human–robot
332 interaction in industrial collaborative robotics: a literature review of the decade 2008–2017.
333 *Advanced Robotics*, 33(15-16):764–799, 2019.
- 334 [9] Muhammed Kocabas, Chun-Hao P. Huang, Otmar Hilliges, and Michael J. Black. PARE: part
335 attention regressor for 3d human body estimation. *CoRR*, abs/2104.08527, 2021.
- 336 [10] Wenhao Li, Hong Liu, Hao Tang, Pichao Wang, and Luc Van Gool. Mhformer: Multi-hypothesis
337 transformer for 3d human pose estimation. *CoRR*, abs/2111.12707, 2021.
- 338 [11] Tao Wang, Jianfeng Zhang, Yujun Cai, Shuicheng Yan, and Jiashi Feng. Direct multi-view
339 multi-person 3d human pose estimation. *Advances in Neural Information Processing Systems*,
340 2021.
- 341 [12] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large
342 scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE*
343 *Transactions on Pattern Analysis and Machine Intelligence*, 36(7):1325–1339, jul 2014.
- 344 [13] Zhongang Cai, Daxuan Ren, Ailing Zeng, Zhengyu Lin, Tao Yu, Wenjia Wang, Xiangyu Fan,
345 Yang Gao, Yifan Yu, Liang Pan, Fangzhou Hong, Mingyuan Zhang, Chen Change Loy, Lei
346 Yang, and Ziwei Liu. Humman: Multi-modal 4d human dataset for versatile sensing and
347 modeling. October 2022.
- 348 [14] Dushyant Mehta, Oleksandr Sotnychenko, Franziska Mueller, Weipeng Xu, Srinath Sridhar,
349 Gerard Pons-Moll, and Christian Theobalt. Single-shot multi-person 3d pose estimation from
350 monocular rgb. In *3D Vision (3DV), 2018 Sixth International Conference on*. IEEE, sep 2018.

- 351 [15] Timo von Marcard, Roberto Henschel, Michael Black, Bodo Rosenhahn, and Gerard Pons-Moll.
 352 Recovering accurate 3d human pose in the wild using imus and a moving camera. In *European*
 353 *Conference on Computer Vision (ECCV)*, sep 2018.
- 354 [16] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery
 355 of human shape and pose. In *Computer Vision and Pattern Recognition (CVPR)*, 2018.
- 356 [17] Leonid Sigal, Alexandru O. Balan, and Michael J. Black. HumanEva: Synchronized video and
 357 motion capture dataset and baseline algorithm for evaluation of articulated human motion. *Int.*
 358 *J. Comput. Vis.*, 87(1-2):4–27, 2010.
- 359 [18] TomasSimon HanbyulJoo, HaoLiu XulongLi, LinGui LeiTan, and TimothyGodisart SeanBaner-
 360 jee. Panoptic studio: A massively multiview system for social interaction capture. *IEEE*
 361 *Transactions on Pattern Analysis and Machine Intelligence*, 41(1), 2019.
- 362 [19] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu,
 363 and Christian Theobalt. Monocular 3d human pose estimation in the wild using improved cnn
 364 supervision. In *3D Vision (3DV), 2017 Fifth International Conference on*. IEEE, 2017.
- 365 [20] Qi Fang, Qing Shuai, Junting Dong, Hujun Bao, and Xiaowei Zhou. Reconstructing 3d human
 366 pose by watching humans in the mirror. In *CVPR*, 2021.
- 367 [21] Rui long Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. Learn to dance with aist++:
 368 Music conditioned 3d dance generation, 2021.
- 369 [22] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose
 370 estimation: New benchmark and state of the art analysis. In *IEEE Conference on Computer*
 371 *Vision and Pattern Recognition (CVPR)*, June 2014.
- 372 [23] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James
 373 Hays, Pietro Perona, Deva Ramanan, Piotr Doll’ar, and C. Lawrence Zitnick. Microsoft COCO:
 374 common objects in context. *CoRR*, abs/1405.0312, 2014.
- 375 [24] H. Jhuang, J. Gall, S. Zuffi, C. Schmid, and M. J. Black. Towards understanding action
 376 recognition. In *International Conf. on Computer Vision (ICCV)*, pages 3192–3199, December
 377 2013.
- 378 [25] Weiyu Zhang, Menglong Zhu, and Konstantinos G. Derpanis. From actemes to action: A
 379 strongly-supervised representation for detailed action understanding. In *2013 IEEE International*
 380 *Conference on Computer Vision*, pages 2248–2255, 2013.
- 381 [26] Mykhaylo Andriluka, Umar Iqbal, Anton Milan, Eldar Insafutdinov, Leonid Pishchulin, Juergen
 382 Gall, and Bernt Schiele. Posetrack: A benchmark for human pose estimation and tracking.
 383 *CoRR*, abs/1710.10000, 2017.
- 384 [27] Hanbyul Joo, Hao Liu, Lei Tan, Lin Gui, Bart Nabbe, Iain Matthews, Takeo Kanade, Shohei
 385 Nobuhara, and Yaser Sheikh. Panoptic studio: A massively multiview system for social motion
 386 capture. In *ICCV*, 2015.
- 387 [28] Rui long Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. Learn to dance with aist++:
 388 Music conditioned 3d dance generation, 2021.
- 389 [29] Julieta Martinez, Rayat Hossain, Javier Romero, and James J. Little. A simple yet effective
 390 baseline for 3d human pose estimation. *CoRR*, abs/1705.03098, 2017.
- 391 [30] Dario Pavllo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3d human
 392 pose estimation in video with temporal convolutions and semi-supervised training. *CoRR*,
 393 abs/1811.11742, 2018.

- 394 [31] Ce Zheng, Sijie Zhu, Matías Mendieta, Taojiannan Yang, Chen Chen, and Zhengming Ding. 3d
 395 human pose estimation with spatial and temporal transformers. *CoRR*, abs/2103.10455, 2021.
- 396 [32] Nikos Kolotouros, Georgios Pavlakos, Michael J. Black, and Kostas Daniilidis. Learning to
 397 reconstruct 3d human pose and shape via model-fitting in the loop. In *Proceedings of the IEEE*
 398 *International Conference on Computer Vision*, 2019.
- 399 [33] Muhammed Kocabas, Chun-Hao P. Huang, Joachim Tesch, Lea Müller, Otmar Hilliges,
 400 and Michael J. Black. SPEC: seeing people in the wild with an estimated camera. *CoRR*,
 401 abs/2110.00620, 2021.
- 402 [34] Jiefeng Li, Chao Xu, Zhicun Chen, Siyuan Bian, Lixin Yang, and Cewu Lu. Hybrik: A hybrid
 403 analytical-neural inverse kinematics solution for 3d human pose and shape estimation. *CoRR*,
 404 abs/2011.14672, 2020.
- 405 [35] Umar Iqbal, Pavlo Molchanov, and Jan Kautz. Weakly-supervised 3d human pose learning
 406 via multi-view images in the wild. In *Proceedings of the IEEE/CVF Conference on Computer*
 407 *Vision and Pattern Recognition*, pages 5243–5252, 2020.
- 408 [36] Rahul Mitra, Nitesh B Gundavarapu, Abhishek Sharma, and Arjun Jain. Multiview-consistent
 409 semi-supervised learning for 3d human pose estimation. In *Proceedings of the ieee/cvf confer-*
 410 *ence on computer vision and pattern recognition*, pages 6907–6916, 2020.
- 411 [37] Bastian Wandt, Marco Rudolph, Petrissa Zell, Helge Rhodin, and Bodo Rosenhahn. Canonpose:
 412 Self-supervised monocular 3d human pose estimation in the wild. In *Proceedings of the*
 413 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13294–13304,
 414 2021.
- 415 [38] Muhammed Kocabas, Salih Karagoz, and Emre Akbas. Self-supervised learning of 3d human
 416 pose using multi-view geometry. In *Proceedings of the IEEE/CVF conference on computer*
 417 *vision and pattern recognition*, pages 1077–1086, 2019.
- 418 [39] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi,
 419 and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis.
 420 *Communications of the ACM*, 65(1):99–106, 2021.
- 421 [40] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black.
 422 SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*,
 423 34(6):248:1–248:16, October 2015.
- 424 [41] Lingjie Liu, Weipeng Xu, Marc Habermann, Michael Zollhöfer, Florian Bernard, Hyeongwoo
 425 Kim, Wenping Wang, and Christian Theobalt. Neural human video rendering by learning
 426 dynamic textures and rendering-to-video translation. *arXiv preprint arXiv:2001.04947*, 2020.
- 427 [42] Atsuhiro Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. Neural articulated radiance
 428 field. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages
 429 5762–5772, 2021.
- 430 [43] Chung-Yi Weng, Brian Curless, and Ira Kemelmacher-Shlizerman. Vid2actor: Free-viewpoint
 431 animatable person synthesis from video in the wild. *arXiv preprint arXiv:2012.12884*, 2020.
- 432 [44] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-
 433 Shlizerman. Humannerf: Free-viewpoint rendering of moving people from monocular video. In
 434 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
 435 16210–16220, 2022.
- 436 [45] Vinoj Jayasundara, Amit Agrawal, Nicolas Heron, Abhinav Shrivastava, and Larry S Davis.
 437 Flexnerf: Photorealistic free-viewpoint rendering of moving humans from sparse views. In
 438 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
 439 21118–21127, 2023.

- 440 [46] Jaehyeok Kim, Dongyoong Wee, and Dan Xu. You only train once: Multi-identity free-viewpoint
 441 neural human rendering from monocular videos. *arXiv preprint arXiv:2303.05835*, 2023.
- 442 [47] Jun Liu, Amir Shahroudy, Mauricio Perez, Gang Wang, Ling-Yu Duan, and Alex C Kot. Ntu
 443 rgb+ d 120: A large-scale benchmark for 3d human activity understanding. *IEEE transactions*
 444 *on pattern analysis and machine intelligence*, 42(10):2684–2701, 2019.
- 445 [48] Mi11. <https://www.mi.com/global/product/mi-11/>, 2022.
- 446 [49] G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000.
- 447 [50] Z. Zhang. A flexible new technique for camera calibration. *IEEE Transactions on Pattern*
 448 *Analysis and Machine Intelligence*, 22(11):1330–1334, 2000.
- 449 [51] D.L. Mills. Internet time synchronization: the network time protocol. *IEEE Transactions on*
 450 *Communications*, 39(10):1482–1493, 1991.
- 451 [52] Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun. Yolox: Exceeding yolo series in
 452 2021. *arXiv preprint arXiv:2107.08430*, 2021.
- 453 [53] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning
 454 for human pose estimation. In *CVPR*, 2019.
- 455 [54] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and
 456 Michael J. Black. Keep it SMPL: Automatic estimation of 3D human pose and shape from a
 457 single image. In *Computer Vision – ECCV 2016*, Lecture Notes in Computer Science. Springer
 458 International Publishing, October 2016.
- 459 [55] Eric Marchand, Hideaki Uchiyama, and Fabien Spindler. Pose estimation for augmented
 460 reality: A hands-on survey. *IEEE Transactions on Visualization and Computer Graphics*,
 461 22(12):2633–2651, 2016.
- 462 [56] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose
 463 estimation using part affinity fields. In *Proceedings of the IEEE conference on computer vision*
 464 *and pattern recognition*, pages 7291–7299, 2017.
- 465 [57] Yilun Chen, Zhicheng Wang, Yuxiang Peng, Zhiqiang Zhang, Gang Yu, and Jian Sun. Cascaded
 466 pyramid network for multi-person pose estimation. In *Proceedings of the IEEE conference on*
 467 *computer vision and pattern recognition*, pages 7103–7112, 2018.
- 468 [58] Hao-Shu Fang, Jiefeng Li, Hongyang Tang, Chao Xu, Haoyi Zhu, Yuliang Xiu, Yong-Lu Li,
 469 and Cewu Lu. Alphapose: Whole-body regional multi-person pose estimation and tracking in
 470 real-time. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- 471 [59] Hanyue Tu, Chunyu Wang, and Wenjun Zeng. Voxelpose: Towards multi-camera 3d human
 472 pose estimation in wild environment. In *European Conference on Computer Vision (ECCV)*,
 473 2020.
- 474 [60] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and
 475 Jonathan T Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown
 476 illumination. *ACM Transactions on Graphics (TOG)*, 40(6):1–18, 2021.
- 477 [61] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unrea-
 478 sonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE*
 479 *conference on computer vision and pattern recognition*, pages 586–595, 2018.
- 480 [62] MMpose Contributors. Openmmlab pose estimation toolbox and benchmark. <https://github.com/open-mmlab/mmPose>, 2020.

- 482 [63] Hui En Pang, Zhongang Cai, Lei Yang, Tianwei Zhang, and Ziwei Liu. Benchmarking and
483 analyzing 3d human pose and shape estimation beyond algorithms. In *Thirty-sixth Conference*
484 *on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.
- 485 [64] Hanyue Tu, Chunyu Wang, and Wenjun Zeng. Voxelpose: Towards multi-camera 3d human
486 pose estimation in wild environment. In *European Conference on Computer Vision*, pages
487 197–212. Springer, 2020.
- 488 [65] Haibo Qiu, Chunyu Wang, Jingdong Wang, Naiyan Wang, and Wenjun Zeng. Cross view fusion
489 for 3d human pose estimation. In *Proceedings of the IEEE/CVF International Conference on*
490 *Computer Vision*, pages 4342–4351, 2019.
- 491 [66] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson,
492 Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv*
493 preprint arXiv:2304.02643, 2023.

494 **Checklist**

- 495 1. For all authors...
 - 496 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
497 contributions and scope? **[Yes]**
 - 498 (b) Did you describe the limitations of your work? **[Yes]**
 - 499 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
 - 500 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
501 them? **[Yes]**
- 502 2. If you are including theoretical results...
 - 503 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
 - 504 (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 505 3. If you ran experiments (e.g. for benchmarks)...
 - 506 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
507 mental results (either in the supplemental material or as a URL)? **[Yes]**
 - 508 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
509 were chosen)? **[Yes]**
 - 510 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
511 ments multiple times)? **[No]**
 - 512 (d) Did you include the total amount of computing and the type of resources used (e.g.,
513 type of GPUs, internal cluster, or cloud provider)? **[Yes]**
- 514 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 515 (a) If your work uses existing assets, did you cite the creators? **[Yes]**
 - 516 (b) Did you mention the license of the assets? **[Yes]**
 - 517 (c) Did you include any new assets either in the supplemental material or as a URL? **[No]**
 - 518 (d) Did you discuss whether and how consent was obtained from people whose data you're
519 using/curating? **[Yes]**
 - 520 (e) Did you discuss whether the data you are using/curating contains personally identifiable
521 information or offensive content? **[N/A]**
- 522 5. If you used crowdsourcing or conducted research with human subjects...
 - 523 (a) Did you include the full text of instructions given to participants and screenshots, if
524 applicable? **[N/A]**
 - 525 (b) Did you describe any potential participant risks, with links to Institutional Review
526 Board (IRB) approvals, if applicable? **[N/A]**
 - 527 (c) Did you include the estimated hourly wage paid to participants and the total amount
528 spent on participant compensation? **[N/A]**