

000
001
002
003
004
005
006
007
008
009
010
011054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Self-supervised 3D Human Pose Estimation from a Single Image

Anonymous CVPR submission

Paper ID 6

Abstract

We propose a new self-supervised method for predicting 3D human body pose from a single image. The prediction network is trained from a dataset of unlabelled images depicting people in typical poses and a set of unpaired 2D poses. By minimising the need for annotated data, the method has the potential for rapid application to pose estimation of other articulated structures (e.g. animals). The self-supervision comes from an earlier idea exploiting consistency between predicted pose under 3D rotation. Our method is a substantial advance on state-of-the-art self-supervised methods in training a mapping directly from images, without limb articulation constraints or any 3D empirical pose prior. We compare performance with state-of-the-art self-supervised methods using benchmark datasets that provide images and ground-truth 3D pose (Human3.6M, MPI-INF-3DHP). Despite the reduced requirement for annotated data, we show that the method outperforms on Human3.6M and matches performance on MPI-INF-3DHP. Qualitative results on a dataset of human hands show the potential for rapidly learning to predict 3D pose for articulated structures other than the human body.

1. Introduction

Estimating 3D pose for articulated objects is a long-standing problem. Its foundations arise from the early days of computer vision with model-based approaches representing the human body as an articulated structure of parts [10, 26]. Interest in estimating 3D human pose grew within the computer vision community motivated by the many real-world applications, for example, pedestrian detection [16], human-computer interaction [42], video surveillance [43], and sports analysis [30].

Initial work on estimating 3D pose addressed this problem by extracting a set of hand-crafted features, for example, segmentation masks [1]. Other early approaches, such as exemplar-based methods, use extensive datasets of 3D poses (commonly constructed from motion capture data) to search for the optimal 3D pose given its 2D projection

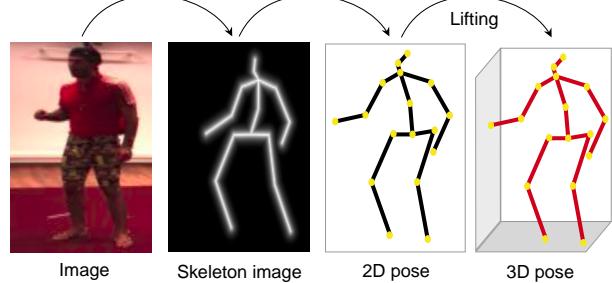


Figure 1. 3D pose estimation pipeline. Our approach jointly learns to estimate 3D pose from an image via intermediate representations of 2D pose. The pipeline is embedded within a larger network for end-to-end training.

[2, 9, 15]. Subsequently, these approaches were surpassed in performance by deep learning methods, initially using supervised learning to regress from images to joint positions [37] or heatmaps [27, 36]. While annotations for 2D joint positions in the image plane are relatively easy to obtain, getting ground truth 3D joint positions from images alone is not straightforward.

Because of the limited availability of 2D and 3D pose annotations, an immense amount of available data in the form of images remains unexploited for training in pose estimation. Although there are many annotated datasets for human pose estimation, the situation is very different for non-human articulated structures such as animals. Recently, unsupervised and self-supervised methods for 3D human pose estimation have made progress using this unlabelled data (at least lacking 3D annotations) and have demonstrated that it is possible to learn to estimate 3D poses from 2D poses relying on rotational consistency [3, 5, 41], multi-view setup [13, 31], and 3D structure constraints e.g. joint angles between limbs [18]. However, what will happen if we do not assume the availability of paired 2D poses? Is it still possible to train a model to predict 3D poses? What are the minimum assumptions to make this possible?

We have proposed a method which, for the first time, learns to map between images and 3D pose without requiring 3D pose annotations or paired 2D pose annotations.

108 Training only needs a set of unlabelled images depicting
109 people in different poses and an unrelated set of 2D human
110 poses. The motivation is twofold: (1) there is the potential
111 for exceeding current levels of performance by training on
112 massive unlabelled datasets, and (2) the method could, in
113 principle, be applied to articulated structures (e.g. animals)
114 where little or no 2D/3D annotated data is available. In our
115 proposed method, we learn a 2D pose predictor and a 3D
116 ‘lifting’ function to produce 3D joint positions from unla-
117 belled images (summarised in [Figure 1](#)) in an end-to-end
118 learning framework.
119

120 Our method simultaneously learns 2D and 3D pose re-
121 presentations in a largely unsupervised fashion, requiring
122 only an empirical prior on unpaired 2D poses. We demon-
123 strate its effectiveness on Human3.6M [\[11\]](#) and MPI-INF-
124 3DHP [\[23\]](#) datasets, two of the most popular benchmarks
125 for human pose estimation. We also show the method’s
126 adaptability to other articulated structures using a synthetic
127 dataset of human hands [\[33\]](#). In experiments, the approach
128 outperforms state-of-the-art self-supervised methods that
129 estimate 3D pose from images and require higher super-
130 vision in training. Overall, our method has the following
131 advantages:
132

- It does not assume any 3D pose annotations or paired 2D pose annotations.
- It holds the potential for quickly adapting to 3D pose prediction for other articulated structures (e.g. animals and jointed inanimate objects).

2. Related Work

140 Our method broadly relates to prior work that estimates
141 3D human pose directly from images, and mainly to self-
142 supervised deep learning methods. However, it also draws
143 inspiration from earlier work on the estimation of 3D pose
144 from 2D pose. Therefore, we review both perspectives, re-
145 gardless of the degree of supervision required for training.
146

3D pose from 2D pose

147 A range of methods take as input 2D poses and lift them
148 to 3D space. Frequently, the 2D poses come from an off-
149 the-shelf 2D pose estimator, or they are simply annotations
150 for a given dataset. Early techniques for estimating 3D
151 poses from 2D joint positions rely on classical classification
152 algorithms and physical constraints. For example, given the
153 joint connectivity and bone lengths, [\[19\]](#) use binary decision
154 trees to estimate the two possible states of a joint with
155 respect to its parent. Other methods implement the near-
156 est neighbour algorithm with large datasets of 3D poses to
157 search for the most likely 3D representation of a particular
158 2D pose [\[2, 9, 15\]](#).

159 In the deep learning context, Martinez *et al.* [\[22\]](#) present
160 a fully supervised approach to predict 3D positions given
161

162 2D joint locations using a fully connected network with
163 residual blocks. This network structure has become pop-
164 ular, and subsequent unsupervised approaches [\[3, 5, 38, 41\]](#)
165 incorporate it within their processes. For instance, Drover
166 *et al.* [\[5\]](#) propose a weakly supervised approach that lifts a
167 2D pose to 3D and then evaluates its 2D projection through
168 a GAN loss. Later, Chen *et al.* [\[3\]](#) extended this work by
169 adding a symmetrical pipeline of consecutive transfor-
170 mations (lifting, rotation, and projection) of the estimated 3D
171 representation. This cycle of transformations exploits geo-
172 metric consistency and removes the dependency on any 3D
173 correspondences.
174

175 More recently, Wandt *et al.* [\[38\]](#) incorporate two funda-
176 mental elements to the model in [\[3\]](#) that increase the per-
177 formance of the 3D lifting process: the use of normalising
178 flow (NF) and a learned elevation angle for the 3D rotations.
179 Previous methods have successfully used normalising flow to
180 estimate 3D prior distributions given 3D human poses [\[40\]](#).
181 However, the method in [\[38\]](#) is the first to use normalising
182 flow to infer the probability of a reconstructed 3D pose from
183 a prior distribution of the 2D input.
184

185 Unlike these previous methods, we do not assume access
186 to ground truth 2D poses as input. Instead, our model takes
187 a single image and predicts the 2D pose from it, which is
188 then lifted to 3D. Overall, it estimates both the 2D and 3D
189 poses from the input image, removing the dependency on
190 paired 2D pose annotations or pre-trained 2D pose predictors.
191

3D pose from an image

192 Work under this category is more related to our ap-
193 proach. Typically, methods for estimating 3D pose from
194 single images break down the task into two steps. First,
195 the 2D joints are localised, and then the 3D pose is es-
196 timated from these 2D joint positions. Early deep-learning
197 implementations of this two-step process depend on hav-
198 ing access to 2D ground truth joint locations and 3D data
199 for supervising training [\[20, 29, 44\]](#). Furthermore, most of
200 those methods integrate pre-trained 2D pose estimators, e.g.
201 a stacked hourglass network [\[27\]](#) and AlphaPose [\[6\]](#), for
202 solving the joint localization step.
203

204 On the unsupervised side, many approaches incorporate
205 specific assumptions for the 2D and 3D joint configura-
206 tion or even add some small portion of actual 3D data to guide
207 the training. For instance, [\[35\]](#) shows a unified multi-stage
208 CNN architecture to estimate 2D and 3D joint locations
209 from single images. This approach relies on a probabilis-
210 tic 3D model of the human pose responsible for lifting the
211 2D representations. Other methods incorporate motion in-
212 formation; for example, Zhou *et al.* [\[45\]](#) use sequences of
213 images and their corresponding 2D pose to guide their 3D
214 pose estimation framework.
215

Like our work, Kundu *et al.* [\[17\]](#) propose a self-

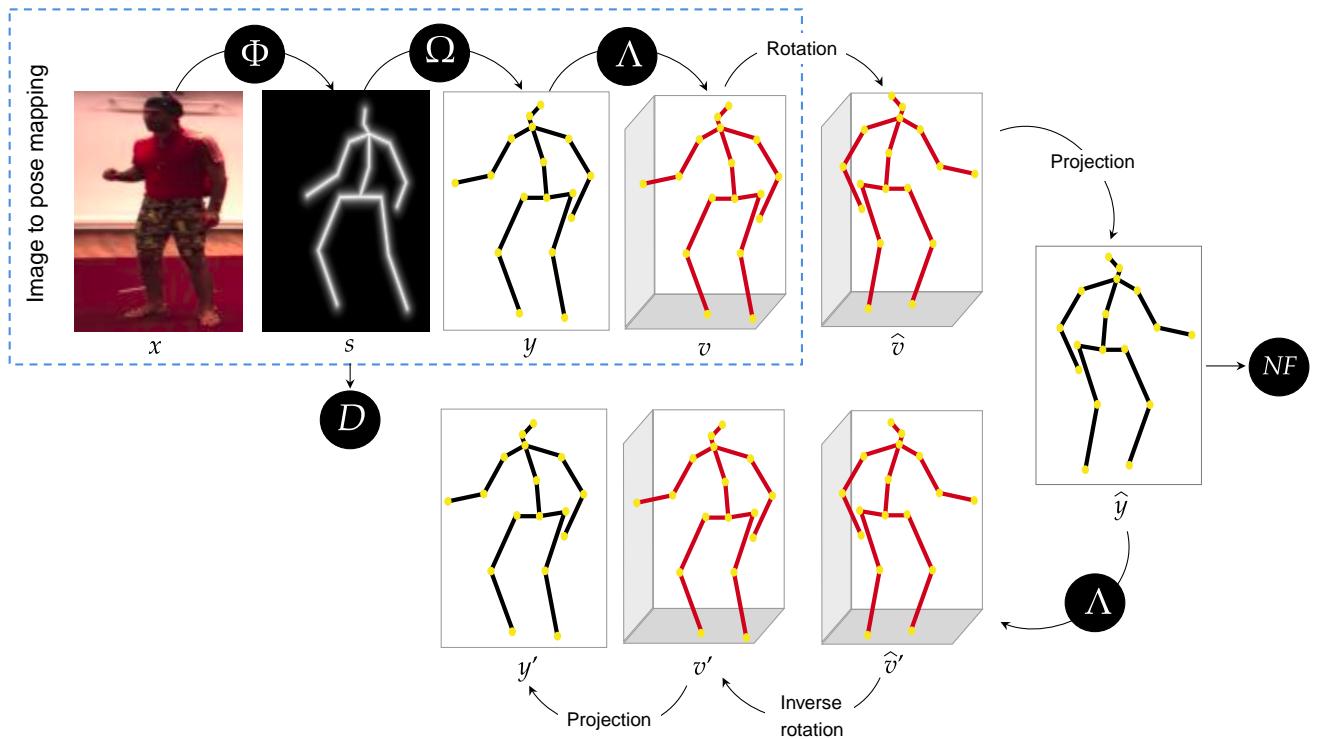


Figure 2. **Self-supervised architecture for estimating the 3D pose of a person.** Our method aims to map an image of a person x to its 3D pose v . To achieve this; we train the networks Φ , Ω , D , and Λ with help of a normalising flow network and a couple of rotations and projections. The pipeline starts with passing an input image x through Φ to obtain a skeleton image s , then Ω produces a set of 2D joint positions y given s . We then embed the mapping into a larger cyclic architecture, where the 2D pose y is lifted into 3D using Λ . The 3D pose v is then rotated, \hat{v} projected into 2D, \hat{y} is lifted back again into 3D through the same Λ , \hat{v}' inversely rotated, and finally v' projected back to the 2D pose y' which should resemble the original 2D pose y .

supervised architecture to learn 3D poses from unlabelled images. They incorporate three assumptions: human pose articulation constraints, a part-based 2D human puppet model, and unpaired 3D poses. Other approaches explore learning without direct supervision by producing synthetic multi-views of the same skeleton [31], accessing multi-view videos [25], or relying solely on 3D kinematic constraints [18]. In contrast, our method does not depend on multi-view images nor any 3D data; it requires only an unlabelled dataset of images depicting people in typical poses and an unpaired empirical prior for the 2D pose.

3. Method

Our proposed 3D pose estimation model consists of a pipeline of three networks Φ , Ω , Λ mapping from full body images to 3D pose. This is shown in the upper-left blue dotted box in Figure 2. The pipeline consists of:

- a Convolutional Neural Network (CNN) Φ mapping from an input image x to an intermediate skeleton image s ,

- a second CNN Ω mapping from s to a 2D pose representation y , and
- a fully connected network Λ lifting the 2D pose y to the required 3D pose v .

The 3D pose is represented as an articulated structure of 3D line segments corresponding to the parts of the body (e.g., torso, head, upper arm, foot).

We train the three networks together by incorporating into a larger network (Figure 2) and optimise end-to-end. This larger network is constructed to incorporate a loop of transformations of the 3D pose. The degree of geometric consistency around the loop contributes to a loss function and provides self-supervision of the training. The training starts with a dataset of images depicting people in different poses. We also assume we have a (normally unrelated) dataset of typical 2D poses from which we obtain skeleton images using a differentiable rendering function κ . These will be used in a GAN framework D to help ensure the generated skeleton images are realistic. In the following sections, we provide more details about the components of our model.

324

Image to 3D pose mapping

The image-to-pose mapping is the composition of networks Φ, Ω, Λ to map an image x showing a person to its 3D pose representation v . The first part of the mapping is a CNN Φ , which maps from the image x to a skeleton image $s = \Phi(x)$ showing the person as a stick figure. Our network Φ adopts a similar architecture to the autoencoder in [14] but without the decoder stage. After training, the emergent skeleton in s aligns with the person in x as expected.

Then, the network Ω maps the skeleton image s to a 2D pose representation $y = \Omega(\Phi(x))$. Informally, Ω learns to extract 2D joint positions (x_i, y_i) from the skeleton image. Finally, Λ is a neural network that lifts the 2D pose to the required pose v in 3D. In particular, $\Lambda(y)$ estimates the depth $z_i = d_i + \Delta$ for each pair of (x_i, y_i) joint positions in the input y , where Δ is a constant depth. Then, the 3D position of joint v_i in the 3D pose v is given by

$$v_i = (x_i \cdot z_i, y_i \cdot z_i, z_i) \quad (1)$$

where z_i is forced to be larger than one, to prevent ambiguity from negative depths. In line with previous works [3, 38, 41], Δ is fixed to 10.

Our lifting network Λ is based on the work of [3, 22] and extended following [38]. In this context, our extended version not only outputs the depth z_i for each joint position (x_i, y_i) in the input, it also produces a value for the elevation angle α . This angle will be used when performing the rotations of the 3D pose v . In particular, we use α to fix the elevation angle of the vertical axis to the ground-plane about which the rotation is performed.

Skeleton images and discriminator

We encourage the training network to generate realistic skeleton images with the help of an empirical prior of 2D poses. Note that these 2D poses are unpaired, i.e., they are not annotations of the training images.

The 2D poses from our empirical prior are first rendered as skeleton images using the renderer proposed by [14]. Let C be a set of connected joint pairs (i, j) , e an image pixel location, and u a set of (x, y) 2D coordinates of body joint positions. The skeleton image rendering function is given by:

$$\kappa(u)_e = \exp\left(-\gamma \min_{(i,j) \in C, r \in [0,1]} \|e - ru_i - (1-r)u_j\|^2\right) \quad (2)$$

Informally κ works by defining a distance field from the line segments linking joints and applies an exponential fall-off to create the image.

Following [14], we use a discriminator network D that uses the prior skeleton images to encourage the predicted skeleton images to represent plausible poses. The task of D

is to determine whether or not a skeleton image $s = \Phi(x)$ looks like an authentic skeleton image such as those in the prior $w = \kappa(u)$. Formally, the goal is to learn $D(s) \in [0, 1]$ to match the reference distribution of skeleton images $p(w)$ and the distribution of predicted skeleton images $q(s)$. An adversarial loss [8] compares the unpaired samples w and the predictions s :

$$\mathcal{L}_D = \mathbb{E}_w(\log(D(w))) + \mathbb{E}_s(\log(1 - D(s))) \quad (3)$$

Random rotations and projections

A fundamental component of our model is the lifting process which allows learning an accurate 3D pose v from the estimated 2D input y . To provide self-supervision of the lifting function and ultimately the whole end-to-end network, we emulate a second virtual view of the 3D pose v by randomly rotating it $\hat{v} = R * v$. Previous work [3] has selected a rotation matrix R by uniformly sampling azimuth and elevation angles from a fixed distribution. Recently, [38] demonstrates that learning the distribution of elevation angles leads to better results. Thus, we follow their approach and use Λ to also predict the elevation angle for the rotation matrix. The rotation around the azimuth axis R_a is sampled from a uniform distribution $[-\pi, \pi]$.

In line with [38], we also predict the dataset's normal distribution of elevation angles R_e by calculating a batch-wise mean μ_e and standard deviation σ_e . We sample from the normal distribution $\mathcal{N}(\mu_e, \sigma_e)$ to rotate the pose in the elevation direction R_e . Then, the complete rotation matrix R is given by $R = R_e^T R_a R_e$.

After rotating the 3D pose, we project \hat{v} through a perspective projection. Then, the same lifting network $\Lambda(\hat{y})$ produces another 3D pose \hat{v}' which is then rotated back to the original view. The final 3D pose v' is projected to 2D using the same perspective projection. This loop of transformations of the 3D pose provides a self-supervised consistency loss. In this context, we assume that if the lifting network Λ accurately estimates the depth for the 2D input y , then the 3D poses \hat{v} and \hat{v}' should be similar. The same principle applies to y and the final 2D projection y' . This gives the following two components of the loss function:

$$\mathcal{L}_{3d} = \|\hat{v}' - \hat{v}\|^2 \quad (4)$$

$$\mathcal{L}_{2d} = \|y' - y\|^2 \quad (5)$$

Additionally, the 3D poses v and v' should be similar. However, instead of comparing with an L_2 loss, we follow [38, 41] and measure the change in the difference in 3D pose between two samples from a batch at corresponding stages in the network.

$$\mathcal{L}_{def} = \|(v'^{(j)} - v'^{(k)}) - (v^{(j)} - v^{(k)})\|^2 \quad (6)$$

Similar to Wandt *et al.* [38], we do not assume samples are from the same video sequence; the samples j and k may come from different sequences and subjects.

Empirical prior on 2D pose

Like Wandt *et al.* [38], we use a normalizing flow to provide a prior over 2D pose. A normalising flow transforms a simple distribution (e.g. a normal distribution) into a complex distribution in such a way that the density of a sample under this complex distribution can be easily computed.

Let $Z \in \mathbb{R}^N$ be a normal distribution and g an invertible function $g(z) = \bar{y}$ with $\bar{y} \in \mathbb{R}^N$ as a projection of the 2D human pose vector \hat{y} in a PCA subspace. By a change of variables, the probability density function for \bar{y} is given by

$$p_Y(\bar{y}) = p_Z(f(\bar{y})) \left| \det \left(\frac{\delta f}{\delta \bar{y}} \right) \right| \quad (7)$$

where f is the inverse of g and $\frac{\delta f}{\delta \bar{y}}$ is the Jacobian of f . Following the normalising flow implementation in [38], we represent f as a neural network [4] and optimise over a dataset of 2D poses with negative log likelihood loss:

$$\mathcal{L}_{NF} = -\log(p_Y(\bar{y})) \quad (8)$$

Additional losses

We compute a loss from the mapping from skeleton images to 2D pose $y = \Omega(s)$. We use the same loss as [14], but without pretraining Ω , i.e., we learn this mapping simultaneously with all the other networks. \mathcal{L}_Ω is given by

$$\mathcal{L}_\Omega = \|(\Omega(\kappa(u)) - u)\|^2 + \lambda \|(\kappa(y) - s)\|^2 \quad (9)$$

where u represents a 2D pose from the unpaired prior, s is the predicted skeleton image, and λ is a balancing coefficient set to 0.1. The function κ is the skeleton image renderer defined in Equation 2.

Based on the proven effectiveness of incorporating relative bone lengths into pose estimation methods [21, 28, 38], we add this to impose a soft constraint when estimating the 3D pose. Following the formulation in [38], we calculate the relative bone lengths b_n for the n -th bone divided by the mean of all bones of a given pose v . We use a pre-calculated relative bone length \bar{b}_n as the mean of a Gaussian prior. Then, the negative log-likelihood of the bone lengths defines a loss function \mathcal{L}_{bl} ,

$$\mathcal{L}_{bl} = -\log \left(\prod_{n=1}^N \mathcal{N}(b_n | \bar{b}_n, \sigma_b) \right) \quad (10)$$

where N is the number of bones defined by the connectivity between joints. Note that this is a soft constraint that allows variation in the relative bone lengths between individuals.

3.1. Training

We train Φ, Ω, D , and Λ from scratch. Only the normalising flow is independently pre-trained, as indicated in [38]. The complete loss function for training our model has seven components expressed as \mathcal{L}_D (Equation 3), \mathcal{L}_Ω (Equation 9), \mathcal{L}_{2d} (Equation 5), \mathcal{L}_{3d} (Equation 4), \mathcal{L}_{def} (Equation 6), \mathcal{L}_{NF} (Equation 8), and \mathcal{L}_{bl} (Equation 10). For convenience in ablation studies, we group three of these loss terms and represent them as \mathcal{L}_{base}

$$\mathcal{L}_{base} = \mathcal{L}_{2d} + \mathcal{L}_{3d} + \mathcal{L}_{def} \quad (11)$$

Thus, the final composite loss function is defined as:

$$\mathcal{L} = \mathcal{L}_D + \mathcal{L}_\Omega + \mathcal{L}_{base} + \mathcal{L}_{NF} + \mathcal{L}_{bl} \quad (12)$$

During testing we only keep the pipeline composed of the trained Φ , Ω , and Λ networks illustrated in the upper-left box in Figure 2. Please see the supplementary section for a more detailed description of the networks and training.

4. Experiments

4.1. Datasets

Human3.6M: Human3.6M [12] is a widely used large-scale pose dataset consisting of videos of 11 subjects doing 17 activities against a static background. There are 3.6M frames depicting the human body and corresponding 3D body poses. In line with the standard protocol II on Human3.6M [12], we use frames from subjects S1, S5, S6, S7, and S8 for training. For testing, we use frames from subjects S9 and S11. We pre-processed the video data by cropping the human body on each frame and removing the background, using the bounding boxes and segmentation masks provided in the dataset.

MPI-INF-3DHP: MPI-INF-3DHP [24] is a human pose dataset containing 3D annotations. Unlike Human3.6M, this dataset incorporates studio and outdoor recordings. The dataset comprises eight subjects with two video sequences for each, doing different activities, e.g. walking, sitting, exercising, and reaching. We train our model with the images from the train split and evaluate with the given test set.

HandDB: HandDB [33] is a dataset of images showing human hands under different scenarios. For our experiments, we use part of the subset of hands generated from synthetic data, which contains 2D annotations for 21 key points: four for each of the five fingers and one for the wrist. For training and testing, we select two sequences (synth2 and synth3) of images from the four included in this subset and split them 80/20, respectively.

540
 541
 542
 543
 544
 545
 546
 547
 548
 549
 550
 551
 552
 553
 554
 555
 556
4.2. Metrics
 557
 558
 559
 560
 561
 562
 563
 564
 565
 Following previous methods [3, 5, 22, 25, 31], we use the standard *Protocol II* to evaluate the Human3.6M dataset [12]. *Protocol II* performs a rigid alignment between the predicted 3D pose and the 3D ground truth via the Procrustes method [7]. Then, it calculates the Mean Per Joint Position Error (MPJPE), which takes the average Euclidean distance between the ground-truth joint positions and the corresponding estimated positions across all 17 joints [12]. For simplicity, we refer to this metric as P-MPJPE (for Procrustes-MPJPE). For evaluation of the MPI-INF-3DHP dataset, we show the Percentage of Correct Keypoints (PCK), which represents the percentage of estimated joint positions within a fixed distance of 150mm from their respective ground truth. We also report the corresponding area under the curve (AUC).

557 558 **4.3. Results**

559
 560 Using our trained model with the Human3.6M dataset,
 561 we predict 3D pose (consisting of 17 joint positions) for
 562 every frame in all videos of subjects S9 and S11 ($\sim 584k$
 563 frames) and calculate the average P-MPJPE. **Table 1** com-
 564 pares our method with the state-of-the-art 3D pose estima-
 565 tion methods.

| 566 Assumptions | 567 Method | 568 P-MPJPE(\downarrow) |
|----------------------------------|-----------------------------------|---|
| 3D pose from 2D landmarks | | |
| 569 3D poses | 570 Martinez <i>et al.</i> [22] | 571 52.1 |
| 572 2D poses | 573 Chen <i>et al.</i> [3] | 574 68.0 |
| 575 2D poses | 576 Drover <i>et al.</i> [5] | 577 64.6 |
| 578 2D poses | 579 Yu <i>et al.</i> [41] | 580 52.3 |
| 581 2D poses | 582 Wandt <i>et al.</i> [38] | 583 36.7 |
| 3D pose from images | | |
| 584 3D Poses | 585 Chen <i>et al.</i> [2] | 586 114.2 |
| 587 3D Poses | 588 Mitra <i>et al.</i> [25] | 589 72.5 |
| 590 MV + P3D | 591 Rhodin <i>et al.</i> [32] | 592 128.6 |
| 593 MV + P3D | 594 Rhodin <i>et al.</i> [31] | 595 98.2 |
| 596 MV + 2D | 597 Wandt <i>et al.</i> [39] | 598 53.0 |
| 599 Unpaired 3D | 600 Kundu <i>et al.</i> [17] | 601 99.2 |
| 602 3D Struct. | 603 Kundu <i>et al.</i> [18] | 604 89.4 |
| 605 Unpaired 2D | 606 Ours | 607 96.7 |

584 **Table 1. P-MPJPE (in mm's) for all activities in Human3.6M.**
 585 MV = Multi view, Unpaired 3D = Unpaired 3D poses, 3D Struct. =
 586 3D body structure constraints, Unpaired 2D = Unpaired 2D poses,
 587 P3D = Partial 3D poses (i.e., one portion of the 3D pose annotations
 588 available).

589
 590 We include supervised [2, 25], semi-supervised [31, 32,
 591 39], and self-supervised [17, 18] approaches that estimate
 592 the 3D pose from images. For a more comprehensive com-
 593 parison, we also consider supervised [22] and unsupervised

594
 595 [3, 5, 34, 38, 41] methods that estimate 3D pose from 2D
 596 landmarks. The performance of our method exceeds that of
 597 some methods that rely on 3D supervision [2], multi-view
 598 images [31, 32] or priors on 3D data [17].

599
 600 To demonstrate our model’s generalisation performance,
 601 we evaluate using the MPI-INF-3DHP test dataset under
 602 different settings. First, we trained the model using Hu-
 603 man3.6M. Second, we train our model with images from the
 604 MPI-INF-3DHP training set. Finally, we train with both im-
 605 ages from the MPI-INF-3DHP training data and the training
 606 set of Human3.6M. For the first experiment, the set of 2D
 607 poses used for the normalising flow prior on 2D pose and
 608 the derived empirical prior on skeletons comes from Hu-
 609 man3.6M, and for the rest the empirical prior comes from
 610 2D poses on MPI-INF-3DHP. **Table 2** presents the PCK and
 611 AUC scores for the different evaluation conditions.

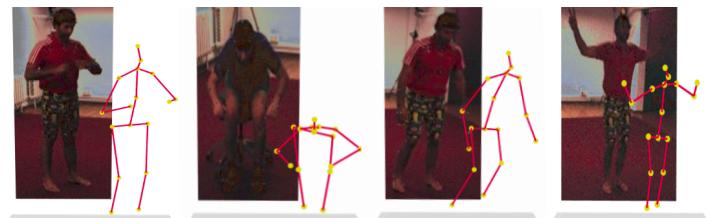
| | Method | PCK(\uparrow) | AUC(\uparrow) |
|-------------------|--------------------------------|-------------------|-------------------|
| 613 Unpaired 3D | 614 Kundu <i>et al.</i> [17] | 615 83.2 | 616 58.7 |
| 616 3D Struct. | 617 Kundu <i>et al.</i> [18] | 618 79.2 | 619 43.4 |
| 620 Unpaired 2D | 621 Ours ⁻ | 622 58.7 | 623 24.3 |
| 624 Unpaired 2D | 625 Ours [*] | 626 69.6 | 627 32.8 |
| 628 Unpaired 2D | 629 Ours ⁺ | 630 75.3 | 631 40.0 |

632 **Table 2. Evaluation results on MPI-INF-3DHP dataset.** First
 633 column shows the main assumption for each method, where Un-
 634 paired 3D = Unpaired 3D poses, 3D Struct. = 3D body struc-
 635 ture constraints, and Unpaired 2D = Unpaired 2D poses. Ours⁻
 636 indicates the model trained with Human3.6M and tested with MPI-
 637 INF-3DHP. Ours^{*} represents the model trained with images from
 638 MPI-INF-3DHP. Ours⁺ indicates that the MPI-INF-3DHP train set
 639 has been extended with images from Human3.6M.

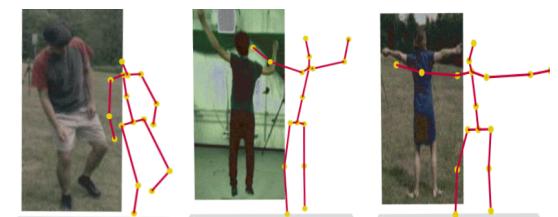
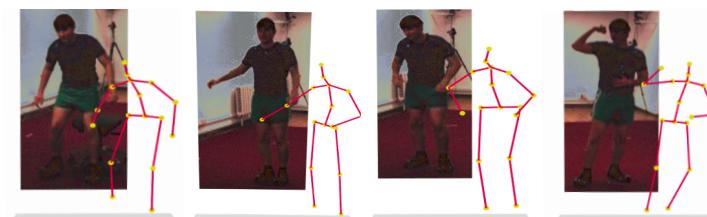
640 **Figure 3** shows qualitative results on images from Hu-
 641 man3.6M and MPI-INF-3DHP datasets. It includes a ran-
 642 dom selection of predicted 3D poses and their correspond-
 643 ing input image.

644 To demonstrate the adaptability of the method, we ap-
 645 plied it to estimate hand pose using part of the synthetic
 646 subset from [33]. This required a different structure for the
 647 target 3D pose. For training and testing the model, we se-
 648 lect sequences of images showing hands under similar con-
 649 ditions (synth2 and synth3). We augment the training set
 650 offline by making two rotated versions of each image (45°
 651 and 90°). We use half of the 2D annotations provided with
 652 the dataset to build the prior on 2D hand poses. The train-
 653 ing set does not include the images corresponding to those
 654 annotations. With the trained model, we estimate 3D hand
 655 poses consisting of 21 key points (representing hand joint
 656 positions). Since the synthetic subset of HandDB does not
 657 contain 3D annotations (just 2D), we only show qualitative
 658 results in **Figure 4**.

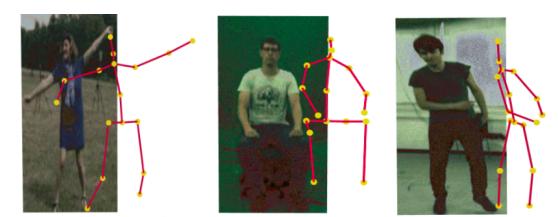
648

A. Results on Human3.6M dataset**B. Results on MPI-INF-3DHP dataset**

649



650



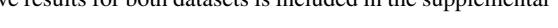
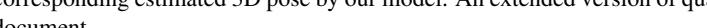
651



652



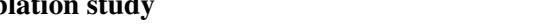
653



654



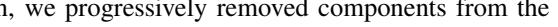
655



656



657



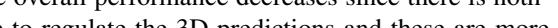
658



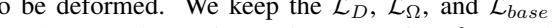
659



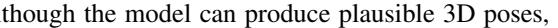
660



661



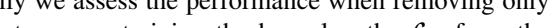
662



663



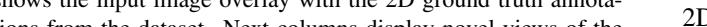
664



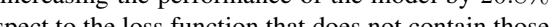
665



666



667



668



669



670



671



672



673

Figure 3. **Qualitative results on images from Human3.6M and MPI-INF-3DHP datasets.** Each figure contains the input image and its corresponding estimated 3D pose by our model. An extended version of qualitative results for both datasets is included in the supplemental document.

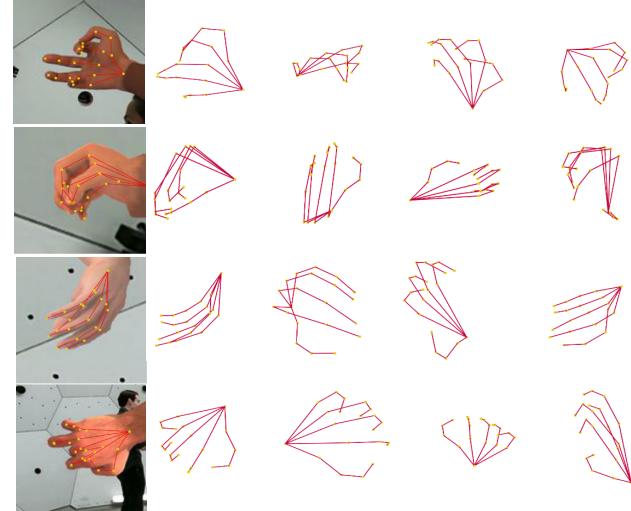


Figure 4. **Qualitative results on HandDB dataset.** First column shows the input image overlay with the 2D ground truth annotations from the dataset. Next columns display novel views of the 3D hand predictions.

4.4. Ablation study

To evaluate the effectiveness of the design of our loss function, we progressively removed components from the complete loss expressed by [Equation 12](#). First, we evaluate the model guided by the discriminator loss \mathcal{L}_D and \mathcal{L}_Ω . As expected, even when the 2D predictions are mostly accurate, the overall performance decreases since there is nothing else to regulate the 3D predictions and these are more likely to be deformed. We keep the \mathcal{L}_D , \mathcal{L}_Ω , and \mathcal{L}_{base} losses for the second experiment, i.e., removing \mathcal{L}_{NF} and \mathcal{L}_{bl} . Although the model can produce plausible 3D poses, the performance is still inferior.

Finally we assess the performance when removing only the loss term constraining the bone lengths \mathcal{L}_{bl} from the original loss formulation ([Equation 12](#)). Adding the combination of the loss terms for the normalising flow prior on 2D pose \mathcal{L}_{NF} and relative bone length \mathcal{L}_{bl} has proven to be useful, increasing the performance of the model by 20.8% with respect to the loss function that does not contain those terms. [Table 3](#) shows the results for each of the modifications to the loss function.

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

| 756 | Configuration | P-MPJPE(\downarrow) |
|-----|---|-------------------------|
| 757 | $\mathcal{L}_D + \mathcal{L}_\Omega$ | 139.3 |
| 758 | $\mathcal{L}_D + \mathcal{L}_\Omega + \mathcal{L}_{base}$ | 122.1 |
| 759 | $\mathcal{L}_D + \mathcal{L}_\Omega + \mathcal{L}_{base} + \mathcal{L}_{NF}$ | 112.3 |
| 760 | $\mathcal{L}_D + \mathcal{L}_\Omega + \mathcal{L}_{base} + \mathcal{L}_{NF} + \mathcal{L}_{bl}$ | 96.7 |
| 761 | | |

Table 3. **Ablation studies.** Experiments with different loss terms using Human3.6M dataset for training and testing.

5. Discussion

Our method outperforms self-supervised state-of-the-art approaches that estimate 3D pose from images and assume unpaired 3D data for supervision [17]. Also, it performs better than recent methods that rely on 3D supervision [2] or multi-view images [31, 32]. Moreover, its performance is similar to one method that assume 3D kinematic constraints [18]. We achieve superior performance to Kundu *et al.* [18] in 20% of the activities (*Discuss*, *Pose*, and *Wait*), and close scores for the rest.

Most failure cases of our model on the Human3.6M dataset appear for snapshots from activities such as *Sitting* and *Sitting Down*. We assume this occurs because of the self-occlusions and perspective ambiguity in these activities. However, according to the examples shown in Figure 5, the model can still produce plausible 3D poses for most cases, even if they do not exactly match their respective 3D ground truth. The high P-MPJPE comes from mismatches between the ‘joints’ representing the body’s extremities, e.g. hands and feet.

With the experiments using the MPI-INF-3DHP dataset, we demonstrate the cross-dataset generalisation of our method. We also show that it works well even if we train the NF network with a different dataset (Human3.6M). Moreover, when training with images from Human3.6M and MPI-INF-3DHP, the experimental results suggest that increasing the diversity of images in the training set could help to improve the overall performance. In principle, extending the dataset of images is relatively straightforward since 3D annotations are not required.

6. Conclusion

We demonstrate how to estimate 3D human pose with a training architecture requiring only images depicting people in different poses and an unpaired set of typical 2D poses. We demonstrate qualitatively that our approach holds the potential for rapidly learning about the pose of articulated structures other than the human body without the need to collect ground-truth 3D pose data, e.g. human hands.

Overall, using human datasets the qualitative and quantitative results suggest that our method is comparable to other self-supervised state-of-the-art approaches that estimate 3D

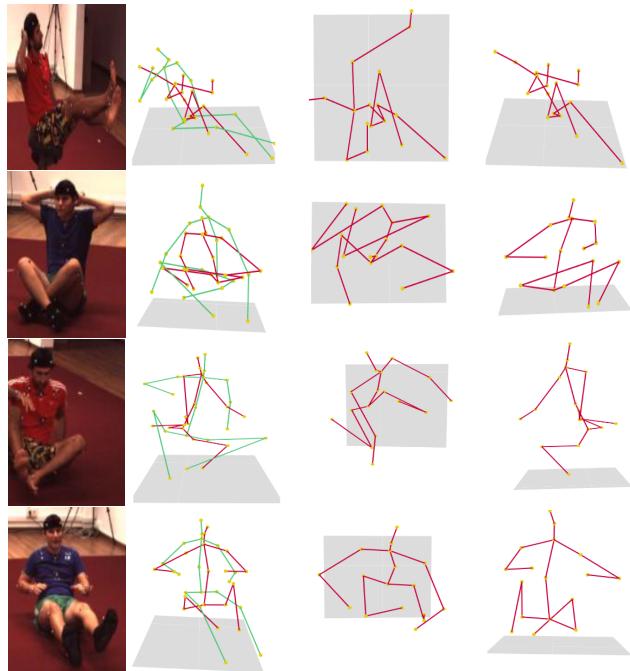


Figure 5. **Failure cases on Human3.6M.** 3D predictions with a P-MPJPE greater than 200mm. The first column shows the input images. The second column displays the predicted 3D pose (coloured in red) aligned with its respective ground truth (coloured in green). Following columns show different views of the predicted 3D pose.

pose from images despite requiring a less onerous dataset for training. Furthermore, it performs better than state-of-the-art methods that rely on multi-view images or 3D pose annotations for supervision. Prior work has demonstrated the value of using temporal information from image sequences and domain adaptation networks. Incorporating these into our approach would be a promising direction for future work. Finally, the way is open to apply the method to much larger datasets of unlabelled images to see whether performance continues to improve, and to apply the method to other articulated structures (e.g., mice, dogs and other animals), exploiting the relatively light requirement for self-supervision in the form of an empirical prior on 2D poses.

References

- [1] Ankur Agarwal and Bill Triggs. Recovering 3d human pose from monocular images. *IEEE transactions on pattern analysis and machine intelligence*, 28(1):44–58, 2005. 1
- [2] Ching-Hang Chen and Deva Ramanan. 3d human pose estimation= 2d pose estimation+ matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7035–7043, 2017. 1, 2, 6, 8
- [3] Ching-Hang Chen, Ambrish Tyagi, Amit Agrawal, Dylan Drover, Stefan Stojanov, and James M Rehg. Unsupervised 3d pose estimation with geometric self-supervision. In *Pro-*

- 864 *ceedings of the IEEE/CVF Conference on Computer Vision*
865 *and Pattern Recognition*, pages 5714–5724, 2019. 1, 2, 4, 6 918
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
- ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 5714–5724, 2019. 1, 2, 4, 6
- [4] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*, 2016. 5
- [5] Dylan Drover, Ching-Hang Chen, Amit Agrawal, Ambrish Tyagi, and Cong Phuoc Huynh. Can 3d pose be learned from 2d projections alone? In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018. 1, 2, 6
- [6] Hao-Shu Fang, Shuqin Xie, Yu-Wing Tai, and Cewu Lu. Rmpe: Regional multi-person pose estimation. In *Proceedings of the IEEE international conference on computer vision*, pages 2334–2343, 2017. 2
- [7] Colin Goodall. Procrustes methods in the statistical analysis of shape. *Journal of the Royal Statistical Society: Series B (Methodological)*, 53(2):285–321, 1991. 6
- [8] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. 4
- [9] Ankur Gupta, Julieta Martinez, James J Little, and Robert J Woodham. 3d pose from motion for cross-view action recognition via non-linear circulant temporal encoding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2601–2608, 2014. 1, 2
- [10] David Hogg. Model-based vision: a program to see a walking person. *Image and Vision computing*, 1(1):5–20, 1983. 1
- [11] 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. 2
- [12] 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, jul 2014. 5, 6
- [13] Umar Iqbal, Pavlo Molchanov, and Jan Kautz. Weakly-supervised 3d human pose learning via multi-view images in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5243–5252, 2020. 1
- [14] Tomas Jakab, Ankush Gupta, Hakan Bilen, and Andrea Vedaldi. Self-supervised learning of interpretable keypoints from unlabelled videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8787–8797, 2020. 4, 5
- [15] Hao Jiang. 3d human pose reconstruction using millions of exemplars. In *2010 20th International Conference on Pattern Recognition*, pages 1674–1677. IEEE, 2010. 1, 2
- [16] Wonhui Kim, Manikandasirram Srinivasan Ramanagopal, Charles Barto, Ming-Yuan Yu, Karl Rosaen, Nick Goumas, Ram Vasudevan, and Matthew Johnson-Roberson. Pedx:
- Benchmark dataset for metric 3-d pose estimation of pedestrians in complex urban intersections. *IEEE Robotics and Automation Letters*, 4(2):1940–1947, 2019. 1 919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

- 972 [29] Sungheon Park, Jihye Hwang, and Nojun Kwak. 3d human
973 pose estimation using convolutional neural networks with 2d
974 pose information. In *European Conference on Computer Vision*,
975 pages 156–169. Springer, 2016. 2
- 976 [30] Konstantinos Rematas, Ira Kemelmacher-Shlizerman, Brian
977 Curless, and Steve Seitz. Soccer on your tabletop. In *Proceedings of the IEEE Conference on Computer Vision and
978 Pattern Recognition*, pages 4738–4747, 2018. 1
- 979 [31] Helge Rhodin, Mathieu Salzmann, and Pascal Fua. Unsu-
980 pervised geometry-aware representation for 3d human pose
981 estimation. In *Proceedings of the European Conference on
982 Computer Vision (ECCV)*, pages 750–767, 2018. 1, 3, 6, 8
- 983 [32] Helge Rhodin, Jörg Spörrli, Isinsu Katircioglu, Victor Con-
984 stantin, Frédéric Meyer, Erich Müller, Mathieu Salzmann,
985 and Pascal Fua. Learning monocular 3d human pose estima-
986 tion from multi-view images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
987 pages 8437–8446, 2018. 6, 8
- 988 [33] Tomas Simon, Hanbyul Joo, Iain Matthews, and Yaser
989 Sheikh. Hand keypoint detection in single images using mul-
990 tiview bootstrapping. In *Proceedings of the IEEE conference on
991 Computer Vision and Pattern Recognition*, pages 1145–
992 1153, 2017. 2, 5, 6
- 993 [34] Cheng Sun, Diego Thomas, and Hiroshi Kawasaki. Unsu-
994 pervised 3d human pose estimation in multi-view-multi-pose
995 video. In *2020 25th International Conference on Pattern
996 Recognition (ICPR)*, pages 5959–5964. IEEE, 2021. 6
- 997 [35] Denis Tome, Chris Russell, and Lourdes Agapito. Lift-
998 ing from the deep: Convolutional 3d pose estimation from
999 a single image. In *Proceedings of the IEEE Conference on
1000 Computer Vision and Pattern Recognition*, pages 2500–
1001 2509, 2017. 2
- 1002 [36] Jonathan Tompson, Ross Goroshin, Arjun Jain, Yann LeCun,
1003 and Christoph Bregler. Efficient object localization using
1004 convolutional networks. In *Proceedings of the IEEE con-
1005 ference on computer vision and pattern recognition*, pages
1006 648–656, 2015. 1
- 1007 [37] Alexander Toshev and Christian Szegedy. Deeppose: Human
1008 pose estimation via deep neural networks. In *Proceedings of
1009 the IEEE conference on computer vision and pattern recog-
1010 nition*, pages 1653–1660, 2014. 1
- 1011 [38] Bastian Wandt, James J Little, and Helge Rhodin. Elepose:
1012 Unsupervised 3d human pose estimation by predicting cam-
1013 era elevation and learning normalizing flows on 2d poses.
1014 *arXiv preprint arXiv:2112.07088*, 2021. 2, 4, 5, 6
- 1015 [39] Bastian Wandt, Marco Rudolph, Petrissa Zell, Helge Rhodin,
1016 and Bodo Rosenhahn. Canonpose: Self-supervised monocular
1017 3d human pose estimation in the wild. In *Proceedings of
1018 the IEEE/CVF Conference on Computer Vision and Pattern
1019 Recognition*, pages 13294–13304, 2021. 6
- 1020 [40] Tom Wehrbein, Marco Rudolph, Bodo Rosenhahn, and Bas-
1021 tian Wandt. Probabilistic monocular 3d human pose es-
1022 timation with normalizing flows. In *Proceedings of the
1023 IEEE/CVF international conference on computer vision*,
1024 pages 11199–11208, 2021. 2
- 1025 [41] Zhenbo Yu, Bingbing Ni, Jingwei Xu, Junjie Wang, Cheng-
1026 long Zhao, and Wenjun Zhang. Towards alleviating the mod-
1027eling ambiguity of unsupervised monocular 3d human pose
1028 estimation. In *Proceedings of the IEEE/CVF International
1029 Conference on Computer Vision*, pages 8651–8660, 2021. 1, 2, 4, 6
- 1030 [42] Zhengyou Zhang. Microsoft kinect sensor and its effect.
1031 *IEEE multimedia*, 19(2):4–10, 2012. 1
- 1032 [43] Liang Zheng, Yujia Huang, Huchuan Lu, and Yi Yang. Pose-
1033 invariant embedding for deep person re-identification. *IEEE
1034 Transactions on Image Processing*, 28(9):4500–4509, 2019.
1035 [44] Xingyi Zhou, Qixing Huang, Xiao Sun, Xiangyang Xue, and
1036 Yichen Wei. Towards 3d human pose estimation in the wild:
1037 a weakly-supervised approach. In *Proceedings of the IEEE
1038 International Conference on Computer Vision*, pages 398–
1039 407, 2017. 2
- 1040 [45] Xiaowei Zhou, Menglong Zhu, Spyridon Leonardos, Kon-
1041 stantinos G Derpanis, and Kostas Daniilidis. Sparseness
1042 meets deepness: 3d human pose estimation from monocular
1043 video. In *Proceedings of the IEEE conference on computer
1044 vision and pattern recognition*, pages 4966–4975, 2016. 2
- 1045 [46] Mingming Jiang, Mingming Tang, and Shihui Wang. Multi-
1046 view 3d human pose estimation via multi-task learning. In
1047 *Proceedings of the IEEE/CVF International Conference on
1048 Computer Vision*, pages 1049–1055, 2019. 1
- 1048 [47] Mingming Jiang, Mingming Tang, and Shihui Wang. Multi-
1049 view 3d human pose estimation via multi-task learning. In
1050 *Proceedings of the IEEE/CVF International Conference on
1051 Computer Vision*, pages 1056–1062, 2019. 1
- 1051 [48] Mingming Jiang, Mingming Tang, and Shihui Wang. Multi-
1052 view 3d human pose estimation via multi-task learning. In
1053 *Proceedings of the IEEE/CVF International Conference on
1054 Computer Vision*, pages 1057–1063, 2019. 1
- 1055 [49] Mingming Jiang, Mingming Tang, and Shihui Wang. Multi-
1056 view 3d human pose estimation via multi-task learning. In
1057 *Proceedings of the IEEE/CVF International Conference on
1058 Computer Vision*, pages 1059–1065, 2019. 1
- 1059 [50] Mingming Jiang, Mingming Tang, and Shihui Wang. Multi-
1060 view 3d human pose estimation via multi-task learning. In
1061 *Proceedings of the IEEE/CVF International Conference on
1062 Computer Vision*, pages 1064–1070, 2019. 1
- 1063 [51] Mingming Jiang, Mingming Tang, and Shihui Wang. Multi-
1064 view 3d human pose estimation via multi-task learning. In
1065 *Proceedings of the IEEE/CVF International Conference on
1066 Computer Vision*, pages 1066–1072, 2019. 1
- 1067 [52] Mingming Jiang, Mingming Tang, and Shihui Wang. Multi-
1068 view 3d human pose estimation via multi-task learning. In
1069 *Proceedings of the IEEE/CVF International Conference on
1070 Computer Vision*, pages 1068–1074, 2019. 1
- 1071 [53] Mingming Jiang, Mingming Tang, and Shihui Wang. Multi-
1072 view 3d human pose estimation via multi-task learning. In
1073 *Proceedings of the IEEE/CVF International Conference on
1074 Computer Vision*, pages 1075–1077, 2019. 1
- 1075 [54] Mingming Jiang, Mingming Tang, and Shihui Wang. Multi-
1076 view 3d human pose estimation via multi-task learning. In
1077 *Proceedings of the IEEE/CVF International Conference on
1078 Computer Vision*, pages 1078–1079, 2019. 1