# $x$R-EgoPose: Egocentric 3D Human Pose from an HMD Camera

Denis Tome[1,2], Patrick Peluse[2], Lourdes Agapito[1] and Hernan Badino[2]
[1]University College London    [2]Facebook Reality Lab

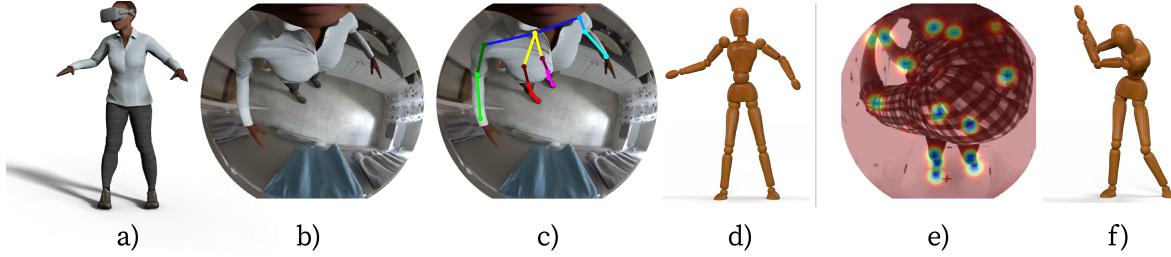{D.Tome, L.Agapito}@cs.ucl.ac.uk    {patrick.peluse, hernan.badino}@fb.com

Figure 1: **Left**: Our $x$R-EgoPose Dataset setup: (a) external camera viewpoint showing a synthetic character wearing the headset; (b) example of photorealistic image rendered from the egocentric camera perspective; (c) 2D and (d) 3D poses estimated with our algorithm. **Right**: results on real images; (e) real image acquired with our HMD-mounted camera with predicted 2D heatmaps; (f) estimated 3D pose, showing good generalization to real images.

## Abstract

*We present a new solution to egocentric 3D body pose estimation from monocular images captured from a downward looking fish-eye camera installed on the rim of a head mounted virtual reality device. This unusual viewpoint, just 2 cm. away from the user's face, leads to images with unique visual appearance, characterized by severe self-occlusions and strong perspective distortions that result in a drastic difference in resolution between lower and upper body. Our contribution is two-fold. Firstly, we propose a new encoder-decoder architecture with a novel dual branch decoder designed specifically to account for the varying uncertainty in the 2D joint locations. Our quantitative evaluation, both on synthetic and real-world datasets, shows that our strategy leads to substantial improvements in accuracy over state of the art egocentric pose estimation approaches. Our second contribution is a new large-scale photorealistic synthetic dataset – $x$R-EgoPose – offering 383K frames of high quality renderings of people with a diversity of skin tones, body shapes, clothing, in a variety of backgrounds and lighting conditions, performing a range of actions. Our experiments show that the high variability in our new synthetic training corpus leads to good generalization to real world footage and to state of the art results on real world datasets with ground truth. Moreover, an evaluation on the Human3.6M benchmark shows that the performance of our method is on par with top performing approaches on the more classic problem of 3D human pose from a third person viewpoint.*

## 1. Introduction

The advent of $x$R technologies (such as AR, VR, and MR) have led to a wide variety of applications in areas such as entertainment, communication, medicine, CAD design, art, and workspace productivity. These technologies mainly focus on immersing the user into a virtual space by the use of a head mounted display (HMD) which renders the environment from the very specific viewpoint of the user. However, current solutions have been focusing so far on the video and audio aspects of the user's perceptual system, leaving a gap in the touch and proprioception senses. Partial solutions to the proprioception problem have been limited to hands whose positions are tracked and rendered in real time by the use of controller devices. The 3D pose of the rest of the body can be inferred from inverse kinematics of the head and hand poses [16], but this often results in inaccurate estimates of the body configuration with a large loss of signal which impedes compelling social interaction [14] and even leads to motion sickness [36].

In this paper we present a new approach for full-body 3D human pose estimation *from a monocular camera installed on a HMD*. In our solution, the camera is mounted on the rim of a HMD looking down, effectively just 2cm. away from an average size nose. With this unique camera viewpoint, most of the lower body appears self-occluded (see right images of Fig. 2). In addition, the strong perspective distortion, due to the fish-eye lens and the camera being so close to the face, results in a drastic difference in resolu-

tion between the upper and lower body (see Fig. 3). Consequently, estimating 2D or 3D pose from images captured from this first person viewpoint is considerably more challenging than from the more standard external perspective and, therefore, even state of the art approaches to human pose estimation [42] underperform on our input data.

Our work tackles the two main challenges described above: *(i)* given the unique visual appearance of our input images and the scarcity of training data for the specific scenario of a HMD mounted camera, we have created a new large scale photorealistic synthetic dataset for training with both 2D and 3D annotations; and *(ii)* to tackle the challenging problem of self-occlusions and difference in resolution between lower and upper body we have proposed a new architecture that accounts for the uncertainty in the estimation of the 2D location of the body joints.

More specifically, our solution adopts a two-step approach. Instead of regressing directly the 3D pose from input images, we first train a model to extract the 2D heatmaps of the body joints and then regress the 3D pose via an autoencoder with a dual branch decoder. While one branch is trained to regress 3D pose from the encoding, the other reconstructs input 2D heatmaps. In this way, the latent vector is enforced to encode the uncertainty in the 2D joint estimates. The auto-encoder helps to infer accurate joint poses for occluded body parts or those with high uncertainty. Both sub-steps are first trained independently and finally end-to-end as the resulting network is fully differentiable. The training is performed on real and synthetic data. The synthetic dataset was created with a large variety of body shapes, environments, and body motions, and will be made publicly available to promote progress in this area.

Our contributions can be summarized as:

- A new encoder-decoder network for egocentric full-body 3D pose estimation from monocular images captured from a camera-equipped VR headset (Sec. 5.2). Our quantitative evaluation on both synthetic and real-world benchmarks with ground truth 3D annotations shows that our approach outperforms previous state of the art [55]. Our ablation studies show that the introduction of our novel decoder branch, trained to reconstruct the 2D input heatmaps, is responsible for the drastic improvements in 3D pose estimation.

- We show that our new approach generalizes, without modifications, to the standard scenario of an external front facing camera. Our method is currently the second best performing after [46] on the Human3.6M benchmark.

- A new large-scale training corpus, composed of 383K frames, that will be made publicly available to promote progress in the area of egocentric human pose

capture (see Section 4). Our new dataset departs from the only other existing monocular egocentric dataset from a headmounted fish-eye camera [55] in its photorealistic quality (see Fig. 2), different viewpoint (since the images are rendered from a camera located on a VR HMD), and its high variability in characters, backgrounds and actions.

## 2. Related Work

We describe related work on monocular (single-camera) marker-less 3D human pose estimation focusing on two distinct capture setups: *outside-in* approaches where an external camera viewpoint is used to capture one or more subjects from a distance – the most commonly used setup; and *first person* or egocentric systems where a head-mounted camera observes the own body of the user. While our paper focuses on the second scenario, we build on recent advances in CNN-based methods for human 3D pose estimation. We also describe approaches that incorporate wearable sensors for first person human pose estimation.

**Monocular 3D Pose Estimation from an External Camera Viewpoint:** the advent of convolutional neural networks and the availability of large 2D and 3D training datasets [18, 3] has recently allowed fast progress in monocular 3D pose estimation from RGB images captured from external cameras. Two main trends have emerged: *(i)* fully supervised regression of 3D joint locations directly from images [22, 31, 47, 58, 32, 27] and *(ii)* pipeline approaches that decouple the problem into the tasks of 2D joint detection followed by 3D lifting [26, 29, 35, 1, 59, 60, 4, 43]. Progress in fully supervised approaches and their ability to generalize has been severely affected by the limited availability of 3D pose annotations for in-the-wild images. This has led to significant efforts in creating photorealistic synthetic datasets [39, 51] aided by the recent availability of parametric dense 3D models of the human body learned from body scans [24]. On the other hand, the appeal of two-step decoupled approaches comes from two main advantages: the availability of high-quality off-the-shelf 2D joint detectors [53, 30, 34, 6] that only require easy-to-harvest 2D annotations, and the possibility of training the 3D lifting step using 3D mocap datasets and their ground truth projections without the need for 3D annotations for images. Even simple architectures have been shown to solve this task with a low error rate [26]. Recent advances are due to combining the 2D and 3D tasks into a joint estimation [41, 42] and using weakly [54, 48, 50, 9, 33] or self-supervised losses [49, 38] or mixing 2D and 3D data for training [46].

**First Person 3D Human Pose Estimation:** while capturing users from an egocentric camera perspective for activity recognition has received significant attention in recent
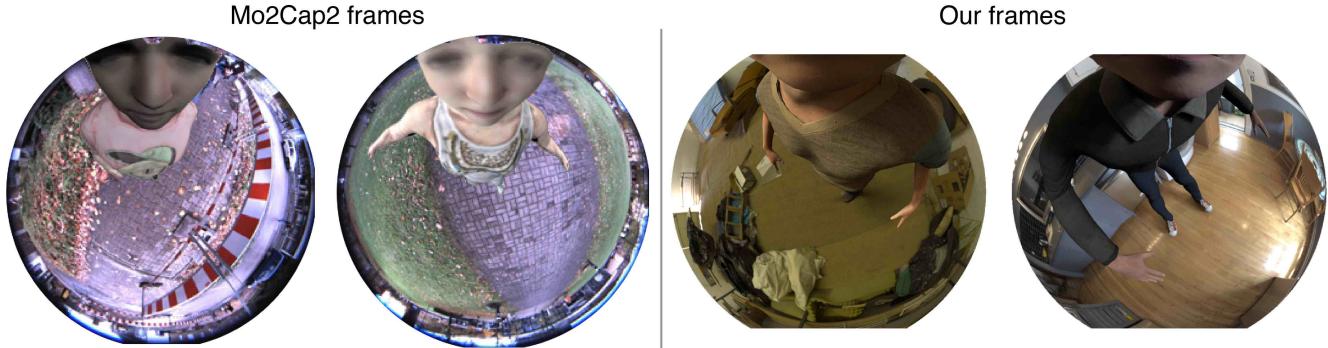
Mo2Cap2 frames          Our frames

Figure 2: Example images from our $x$R-EgoPose Dataset compared with the competitor Mo2Cap2 dataset [55]. The quality of our frames is far superior than the randomly sampled frames from mo2cap2, where the characters suffer color matching with respect to the background light conditions.

years [11, 25, 5], most methods detect, at most, only upper body motion (hands, arms or torso). Capturing full 3D body motion from head-mounted cameras is considerably more challenging. Some head-mounted capture systems are based on RGB-D input and reconstruct mostly hand, arm and torso motions [40, 57]. Jiang and Grauman [20] reconstruct full body pose from footage taken from a camera worn on the chest by estimating egomotion from the observed scene, but their estimates lack accuracy and have high uncertainty. A step towards dealing with large parts of the body not being observable was proposed in [2] but for external camera viewpoints. Rhodin *et al.* [37] pioneered the first approach towards full-body capture from a helmet-mounted stereo fish-eye camera pair. The cameras were placed around 25 cm away from the user's head, using telescopic sticks, which resulted in a fairly cumbersome setup for the user but with the benefit of capturing large field of view images where most of the body was in view. Monocular head-mounted systems for full-body pose estimation have more recently been demonstrated by Xu *et al.* [55] (who propose a real-time compact setup mounted on a baseball cap) although in this case the egocentric camera is placed a few centimeters further from the user's forehead than in our proposed approach. Our approach substantially outperforms Xu *et al.*'s method [55] by 20% or more on both indoor and outdoor sequences from their real world evaluation dataset.

**3D Pose Estimation from Wearable Devices:** Inertial Measurement Units (IMUs) worn by the subject provide a camera-free alternative solution to first person human pose estimation. However, such systems are intrusive and complex to calibrate. While reducing the number of sensors leads to a less invasive configuration [52] recovering accurate human pose from sparse sensor readings becomes a more challenging task. An alternative approach, introduced by Shiratori *et al.* [44] consists of a multi-camera structure-from-motion (SFM) approach using 16 limb-mounted cameras. Still very intrusive, this approach suffers from motion blur, automatic white balancing, rolling shutter effects and

motion in the scene, making it impractical in realistic scenarios.

## 3. Challenges in Egocentric Pose Estimation

Fig. 3 provides a visualization of the unique visual appearance of our HMD egocentric setup — the top row shows which body parts would become self-occluded from an egocentric viewpoint and dark green indicates highest and bright red lowest pixel resolution. This unusual visual appearance calls both for a new approach and a new training corpus. Our paper tackles both. Our new neural network architecture is specifically designed to encode the difference in uncertainty between upper and lower body joints caused by the varying resolution, extreme perspective effects and self-occlusions. On the other hand, our new large-scale synthetic training set — $x$R-EgoPose — contains 383K images rendered from a novel viewpoint: a fisheye camera mounted on a VR display. It has quite superior levels of photorealism in contrast with the only other existing monocular egocentric dataset [55] (see Fig. 2 for a side to side comparison), and large variability in the data. To enable quantitative evaluations on real world images, we contribute $x$R-EgoPose$^R$, a smaller scale real-world dataset acquired with a lightweight setup – a real fish-eye camera mounted on a VR display – with ground truth 3D pose annotations. Our extensive experimental evaluations show that our new approach outperforms the current state of the art in monocular egocentric 3D pose estimation [55] both on synthetic and real-world datasets.

## 4. $x$R-EgoPose Synthetic Dataset

The design of the dataset focuses on scalability, with augmentation of characters, environments, and lighting conditions. A rendered scene is generated from a random selection of characters, environments, lighting rigs, and animation actions. The animations are obtained from mocap data. A small random displacement is added to the positioning of the camera on the headset to simulate the typical
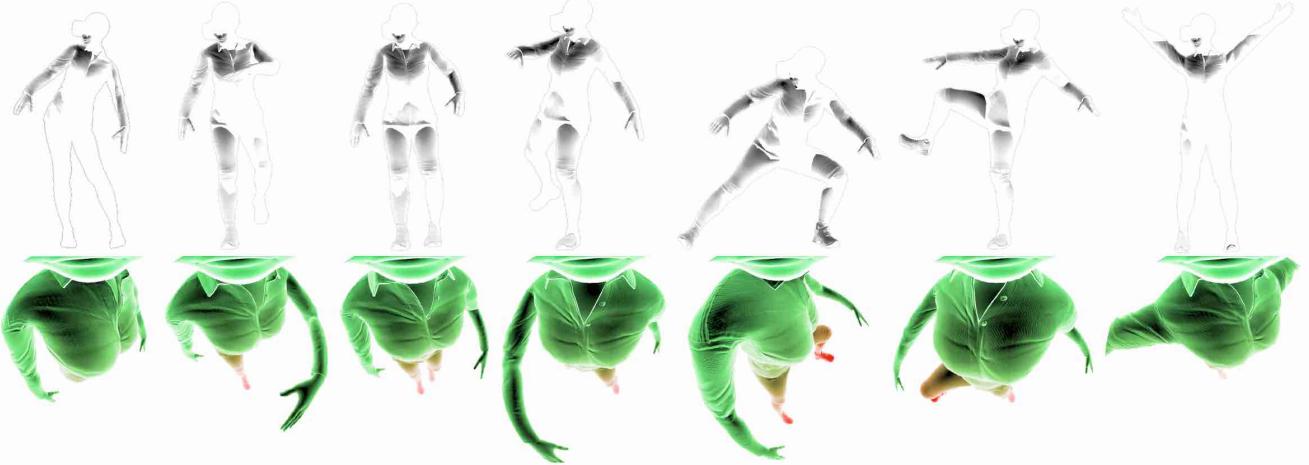
Figure 3: Visualization of different poses with the same synthetic actor. **Top:** poses rendered from an external camera viewpoint. White represents occluded areas of the body. poses rendered from the egocentric camera viewpoint. The color gradient indicates the density of image pixels for each area of the body: *green* indicates higher pixel density, whereas *red* indicates lower density. This figure illustrates the most important challenges faced in egocentric human pose estimation: severe self-occlusions, extreme perspective effects and drastically lower pixel density for the lower body.

variation of the pose of the headset with respect to the head when worn by the user.

**Characters**: To improve the diversity of body types, from a single character we generate additional *skinny short*, *skinny tall*, *full short*, and *full tall* versions The height distribution of each version varies from 155 cm. to 189 cm.

**Skin**: color tones include *white* (Caucasian, freckles or Albino), *light-skinned European*, *dark-skinned European* (darker Caucasian, European mix), *Mediterranean or olive* (Mediterranean, Asian, Hispanic, Native American), *dark brown* (Afro-American, Middle Eastern), and *black* (Afro-American, African, Middle Eastern). Additionally, we built random skin tone parameters into the shaders of each character used with the scene generator.

**Clothing**: Clothing types include Athletic Pants, Jeans, Shorts, Dress Pants, Skirts, Jackets, T-Shirts, Long Sleeves, and Tank Tops. Shoes include Sandals, Boots, Dress Shoes, Athletic Shoes, Crocs. Each type is rendered with different texture and colors.

**Actions**: the type of actions are listed in Table 1.

**Images**: the images have a resolution of $1024 \times 1024$ pixels and 16-bit color depth. For training and testing, we downsample the color depth to 8 bit. The frame rate is 30 fps. *RGB*, *depth*, *normals*, *body segmentation*, and *pixel world position* images are generated for each frame, with the option for exposure control for augmentation of lighting. Metadata is provided for each frame including 3D joint positions, height of the character, environment, camera pose, body segmentation, and animation rig.

**Render quality**: Maximizing the photorealism of the synthetic dataset was our top priority. Therefore, we animated the characters in Maya using actual mocap data [17], and used a standardized physically based rendering setup with

V-Ray. The characters were created with global custom shader settings applied across clothing, skin, and lighting of environments for all rendered scenes.

## 4.1. Training, Test, and Validation Sets

The dataset has a total size of 383K frames, with 23 male and 23 female characters, divided into three sets: *Train-set:* 252K frames; *Test-set:* 115K frames; and *Validation-set:* 16K frames. The gender distribution is: *Train-set:* 13M/11F, *Test-set:* 7M/5F and *Validation-set:* 3M/3F. Table 1 provides a detailed description of the partitioning of the dataset according to the different actions.

| Action | N. Frames | Size Train | Size Test |
|---|---|---|---|
| Gaming | 24019 | 11153 | 4684 |
| Gesticulating | 21411 | 9866 | 4206 |
| Greeting | 8966 | 4188 | 1739 |
| Lower Stretching | 82541 | 66165 | 43491 |
| Patting | 9615 | 4404 | 1898 |
| Reacting | 26629 | 12599 | 5104 |
| Talking | 13685 | 6215 | 2723 |
| Upper Stretching | 162193 | 114446 | 46468 |
| Walking | 34989 | 24603 | 9971 |

Table 1: Total number of frames per action and their distribution between train and test sets. Everything else not mentioned is validation data.

## 5. Architecture

Our proposed architecture, shown in Fig. 4, is a two step approach with two modules. The first module detects 2D heatmaps of the locations of the body joints in image space using a ResNet [13] architecture. The second module takes the 2D heatmaps as inputs and regresses the 3D coordinates
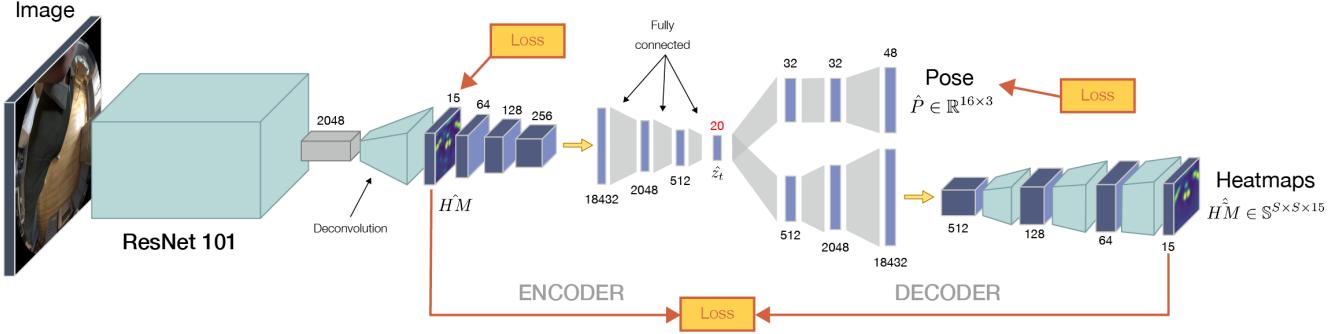
Figure 4: Our novel two-step architecture for egocentric 3D human pose estimation has two modules: *a)* the 2D heatmap estimator, based on ResNet101 [13] as the core architecture; *b)* the 3D lifting module takes 2D heatmaps as input and is based on our novel dual branch auto-encoder.

of the body joints using a novel dual branch auto-encoder.

One of the most important advantage of this pipeline approach is that 2D and 3D modules can be trained independently according to the available training data. For instance, if a sufficiently large corpus of images with 3D annotations is unavailable, the 3D lifting module can be trained instead using 3D mocap data and projected heatmaps without the need of paired images. Once the two modules are pre-trained the entire architecture can be fine-tuned end-to-end since it is fully differentiable. A further advantage of this architecture is that the second branch is only needed at training time (see Sec. 5.2) and can be removed at test time, guaranteeing the same performance and a faster execution.

### 5.1. 2D Pose Detection

Given an RGB image $\mathbf{I} \in \mathbb{R}^{368 \times 368 \times 3}$ as input, the 2D pose detector infers 2D poses, represented as a set of heatmaps $\mathbf{HM} \in \mathbb{R}^{47 \times 47 \times 15}$, one for each of the body joints. For this task we have used a standard *ResNet 101* [13] architecture, where the last average pooling and fully connected layers have been replaced by a deconvolutional layer, with kernel size = 3 and stride = 2. The weights have been randomly initialized using Xavier initialization [12]. The model was trained using normalized input images, obtained by subtracting the mean value and dividing by the standard deviation, and using the mean square error of the difference between the ground truth heatmaps and the predicted ones as the loss:

$$L_{2D} = \text{mse}(\mathbf{HM}, \widehat{\mathbf{HM}}) \qquad (1)$$

We also trained alternative 2D pose detectors including the CPM [53] and the Stacked Hourglass Network [30] resulting in comparable performance at a higher computational cost.

### 5.2. 2D-to-3D Mapping

The 3D pose module takes as input the 15 heatmaps computed by the previous module and outputs the final 3D pose $\mathbf{P} \in \mathbb{R}^{16 \times 3}$. Note that the number of output 3D joints is 16 since we include the head, whose position cannot be estimated in the 2D images, as the person is wearing a headset,

but can be regressed in 3D. In most pipeline approaches the *3D lifting* module typically takes as input the 2D coordinates of the detected joints. Instead, similarly to [33], our approach regresses the 3D pose from heatmaps, not just 2D locations. The main advantage is that these carry important information about the uncertainty of the 2D pose estimates.

The main novelty of our architecture (see Fig. 4), is that we ensure that this uncertainty information is not lost. While the encoder takes as input a set of heatmaps, and encodes them into the embedding $\hat{\mathbf{z}}$, the decoder has two branches – one to regress the 3D pose from $\hat{\mathbf{z}}$ and another to reconstruct the input heatmaps. The purpose of this branch is to force the latent vector to encode the probability density function of the estimated 2D heatmaps.

The overall loss function for the auto-encoder becomes

$$\begin{aligned} L_{\text{AE}} = \; & \lambda_p(||\mathbf{P} - \hat{\mathbf{P}}||^2 + R(\mathbf{P}, \hat{\mathbf{P}})) + \\ & \lambda_{hm}||\widehat{\mathbf{HM}} - \widetilde{\mathbf{HM}}||^2 \end{aligned} \qquad (2)$$

$\mathbf{P}$ the ground truth; $\widetilde{\mathbf{HM}}$ is the set of heatmaps regressed by the decoder from the latent space and $\widehat{\mathbf{HM}}$ are the heatmaps regressed by ResNet (see Sec. 5.1). Finally $R$ is the loss over the 3D poses $R(\mathbf{P}, \hat{\mathbf{P}}) = \lambda_\theta \theta(\mathbf{P}, \hat{\mathbf{P}}) + \lambda_L L(\mathbf{P}, \hat{\mathbf{P}})$ with

$$\theta(\mathbf{P}, \hat{\mathbf{P}}) = \sum_l^L \frac{\mathbf{P}_l \cdot \hat{\mathbf{P}}_l}{||\mathbf{P}|| * ||\hat{\mathbf{P}}_l||} \quad L(\mathbf{P}, \hat{\mathbf{P}}) = \sum_l^L ||\mathbf{P}_l - \hat{\mathbf{P}}_l||$$

corresponding to the cosine-similarity error and the limb-length error, with $\mathbf{P}_l \in \mathbb{R}^3$ the $l^{th}$ limb of the pose. An important advantage of this loss is that the model can be trained on a mix of 3D and 2D datasets simultaneously: if an image sample only has 2D annotations then $\lambda_p = 0$, such that only the heatmaps are contributing to the loss. In Section 6.2 we show how having a larger corpus of 2D annotations can be leveraged to improve final 3D body pose estimates.

### 5.3. Training Details

The model has been trained on the entire training set for 3 epochs, with a learning rate of $1e - 3$ using batch

| Approach | Evaluation error (mm) | Gaming | Gesticulating | Greeting | Lower Stretching | Patting | Reacting | Talking | Upper Stretching | Walking | All |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Martinez [26] | Upper body | 58.5 | 66.7 | 54.8 | 70.0 | 59.3 | 77.8 | 54.1 | 89.7 | 74.1 | 79.4 |
| | Lower body | 160.7 | 144.1 | 183.7 | 181.7 | 126.7 | 161.2 | 168.1 | 159.4 | 186.9 | 164.8 |
| | Average | 109.6 | 105.4 | 119.3 | 125.8 | 93.0 | 119.7 | 111.1 | 124.5 | 130.5 | 122.1 |
| **Ours - single-branch** | Upper body | 114.4 | 106.7 | 99.3 | 90.9 | 99.1 | 147.5 | 95.1 | 119.0 | 104.3 | 112.5 |
| | Lower body | 162.2 | 110.2 | 101.2 | 175.6 | 136.6 | 203.6 | 91.9 | 139.9 | 159.0 | 148.3 |
| | Average | 138.3 | 108.5 | 100.3 | 133.3 | 117.8 | 175.6 | 93.5 | 129.0 | 131.9 | 130.4 |
| **Ours - dual-branch** | Upper body | 48.8 | 50.0 | 43.0 | 36.8 | 48.6 | 56.4 | 42.8 | 49.3 | 43.2 | **50.5** |
| | Lower body | 65.1 | 50.4 | 46.1 | 65.2 | 70.2 | 65.2 | 45.0 | 58.8 | 72.2 | **65.9** |
| | Average | 56.0 | 50.2 | 44.6 | 51.1 | 59.4 | 60.8 | 43.9 | 53.9 | 57.7 | **58.2** |

Table 2: Quantitative evaluation with Martinez *et al.* [26], a state-of-the-art approach developed for front-facing cameras. Both upper and lower body reconstructions are shown as well. A comparison with our own architecture using a single-branch decoder is also included. Note how the competing approach fails consistently across different actions in lower body reconstructions. This experiment emphasizes how, even a state-of-the-art 3D lifting method developed for external cameras fails on this challenging task.

normalization on a mini-batch of size 16. The deconvolutional layer used to identify the heatmaps from the features computed by *ResNet* has kernel size $= 3$ and stride $= 2$. The convolutional and deconvolutional layers of the encoder have kernel size $= 4$ and stride $= 2$. Finally, all the layers of the encoder use leakly ReLU as activation function with $0.2$ leakiness. The $\lambda$ weights used in the loss function were identified through grid search and set to $\lambda_{hm} = 10^{-3}$, $\lambda_p = 10^{-1}$, $\lambda_\theta = -10^{-2}$ and $\lambda_L = 0.5$ . The model has been trained from scratch with Xavier weight initializer.

## 6. Quantitative Evaluation

We evaluate the proposed approach quantitatively on a variety of egocentric 3D human pose datasets: *(i)* the test-set of $x$R-EgoPose, our synthetic corpus, *(ii)* the test-set of $x$R-EgoPose$^R$, our smaller scale real-world dataset acquired with a real fish-eye camera mounted on a VR display and with ground truth 3D poses, and *(iii)* the Mo2Cap2 test-set [55] which includes 2.7K frames of real images with ground truth 3D poses of two people captured in indoor and outdoor scenes.

In addition we evaluate quantitatively on the Human3.6M dataset to show that our architecture generalizes well without any modifications to the case of an external camera viewpoint.

**Evaluation protocol**: Unless otherwise mentioned, we report the Mean Per Joint Position Error - MPJPE:

$$E(\mathbf{P}, \hat{\mathbf{P}}) = \frac{1}{N_f} \frac{1}{N_j} \sum_{f=1}^{N_f} \sum_{j=1}^{N_j} ||\mathbf{P}_j^{(f)} - \hat{\mathbf{P}}_j^{(f)}||_2 \quad (3)$$

where $\mathbf{P}_j^{(f)}$ and $\hat{\mathbf{P}}_j^{(f)}$ are the 3D points of the ground truth and predicted pose at frame $f$ for joint $j$, out of $N_f$ number of frames and $N_j$ number of joints.

### 6.1. Evaluation on our Egocentric Synthetic Dataset

**Evaluation on $x$R-EgoPose test-set**: Firstly, we evaluate our approach on the test-set of our synthetic $x$R-EgoPose

dataset. It was not possible to establish a comparison with state of the art monocular egocentric human pose estimation methods such as Mo2Cap2 [55] given that their code has not been made publicly available. Instead we compare with Martinez *et al.* [26], a recent state of the art method for a traditional external camera viewpoint. For a fair comparison the training-set of our $x$R-EgoPose dataset has been used to re-train this model. In this way we can directly compare the performance of the 2D to 3D modules.

Table 2 reports the MPJPE (Eq. 3) for both methods showing that our approach (Ours-dual-branch) outperforms that by Martinez *et al.* by 36.4% in the upper body reconstruction, 60% in the lower body reconstruction, and 52.3% overall, showing a considerable improvement.

**Effect of the second decoder branch:** Table 2 also reports an ablation study to compare the performance of two versions of our approach: with (Ours-dual-branch) and without (Ours-single-branch) the second branch for the decoder which reconstructs the heatmaps $\hat{H\!M}$ from the encoding **z**. The overall average error of the single branch encoder is $130.4$ mm, far from the $58.2$ mm error achieved by our novel dual-branch architecture.

**Reconstruction errors per joint type**: Table 4 reports a decomposition of the reconstruction error into different individual joint types. The highest errors are in the hands (due to hard occlusions when they go outside of the field of view) and feet (due to self-occlusions and low resolution).

### 6.2. Evaluation on Egocentric Real Datasets

**Comparison with Mo2Cap2 [55]**: We compare the results of our approach with those given by our direct competitor, Mo2Cap2, on their real world test set, including both indoor and outdoor sequences. To guarantee a fair comparison, the authors of [55] provided us the heatmaps from their 2D joint estimator. In this way, both 3D reconstruction networks use the same input. Table 5 reports the MPJPE errors for both methods. Our dual-branch approach substantially outperforms Mo2Cap2 [55] in both indoor and outdoor scenarios.

| **Protocol #1** | Chen [7] | Hossain [15]* | Dabral [8]* | Tome [48] | Moreno [29] | Kanazawa [21] | Zhou [61] | Jahangiri [19] | Mehta [27] | Martinez [26] | Fang [10] | Sun [45] | Sun [46] | **Ours** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Errors (mm)** | 114.2 | 51.9 | 52.1 | 88.4 | 87.3 | 88.0 | 79.9 | 77.6 | 72.9 | 62.9 | 60.4 | 59.1 | **49.6** | 53.4 |
| **Protocol #2** | Yasin [56] | Hossain [15]* | Dabral [8]* | Rogez [39] | Chen [7] | Moreno [29] | Tome [48] | Zhou [61] | Martinez [26] | Kanazawa [21] | Sun [45] | Fang [10] | Sun [46] | **Ours** |
| **Errors (mm)** | 108.3 | 42.0 | 36.3 | 88.1 | 82.7 | 76.5 | 70.7 | 55.3 | 47.7 | 58.8 | 48.3 | 45.7 | **40.6** | 45.24 |

Table 3: Comparison with other state-of-the-art approaches on the Human3.6M dataset (front-facing cameras). Approaches with * make use of temporal information.

Note that the dataset provided by the stereo egocentric system EgoCap [37] cannot be directly used for comparison, due to the hugely different camera position relative to the head (their stereo cameras are 25cm. from the head).

**Evaluation on $x$R-EgoPose$^\mathbf{R}$**: The $\sim$ 10K frames of our small scale real-world data set were captured from a fisheye camera mounted on a VR HMD worn by three different actors wearing different clothes, and performing 6 different actions. The ground truth 3D poses were acquired using a custom mocap system. The network was trained on our synthetic corpus ($x$R-EgoPose) and fine-tuned using the data from two of the actors. The test set contained data from the unseen third actor. Examples of the input views and the reconstructed poses are shown in Fig. 6). The MPJPE [18] errors (Eq. 3) are shown in Table 6. These results show good generalization of the model (trained mostly on synthetic data) to real images.

## 6.3. Evaluation on Front-facing Cameras

**Comparison on Human3.6M dataset**: We show that our proposed approach is not specific for the egocentric case, but also provides excellent results in the more standard case of front-facing cameras. For this evaluation, we chose the Human3.6M dataset [18]. We used two evaluation protocols. *Protocol 1* has five subjects (S1, S5, S6, S7, S8) used in training, with subjects (S9, S11) used for evaluation. The MPJPE error is computed on every 64th frame. *Protocol 2* contains six subjects (S1, S5, S6, S7, S8, S9) used for training, and the evaluation is performed on every 64th frame of Subject 11 (Procrustes aligned MPJPE is used for evaluation). The results are shown in Table 3, from where it

| Joint | Error (mm) | Joint | Error (mm) |
|---|---|---|---|
| Left Leg | 34.33 | Right Leg | 33.85 |
| Left Knee | 62.57 | Right Knee | 61.36 |
| Left Foot | 70.08 | Right Foot | 68.17 |
| Left Toe | 76.43 | Right Toe | 71.94 |
| Neck | 6.57 | Head | 23.20 |
| Left Arm | 31.36 | Right Arm | 31.45 |
| Left Elbow | 60.89 | Right Elbow | 50.13 |
| Left Hand | 90.43 | Right Hand | 78.28 |

Table 4: Average reconstruction error per joint using Eq. 3, evaluated on the entire test-set (see Sec. 4).
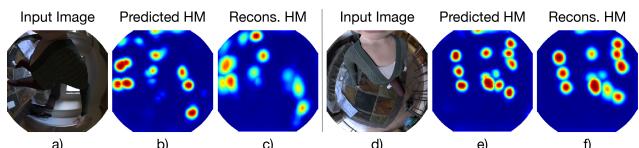


Figure 5: Examples of reconstructed heatmaps generated by the latent vector **z**. They reproduce the correct uncertainty of the predicted 2D joint positions.

can be seen that our approach is on par with state-of-the-art methods, scoring second overall within the non-temporal methods.

## 6.4. Mixing 2D and 3D Ground Truth Datasets

An important advantage of our architecture is that the model can be trained on a mix of 3D and 2D datasets simultaneously: if an image sample only has 2D annotations but no 3D ground truth labels, the sample can still be used and only the heatmaps will contribute to the loss. We evaluated the effect of adding additional images with 2D but no 3D labels on both scenarios: egocentric and front-facing cameras. In the egocentric case we created two subsets of the $x$R-EgoPose test-set. The first subset contained 50% of all the available image samples with both 3D and 2D labels. The second contained 100% of the image samples with 2D labels, but only 50% of the 3D labels. Effectively the second subset contained twice the number of images with 2D annotations only. Table 7a compares the results between the two subsets, from where it can be seen that the final 3D pose estimate benefits from additional 2D annotations. Equivalent behavior is seen on the Human3.6M dataset. Figure 7b shows the improvements in reconstruction error when additional 2D annotations from COCO [23] and MPII [3] are used.

## 6.5. Encoding Uncertainty in the Latent Space

Figure 5 demonstrates the ability of our approach to encode the uncertainty of the input 2D heatmaps in the latent vector. Examples of input 2D heatmaps and those reconstructed by the second branch of the decoder are shown for comparison.

| INDOOR | walking | sitting | crawling | crouching | boxing | dancing | stretching | waving | total (mm) |
|---|---|---|---|---|---|---|---|---|---|
| 3DV'17 [27] | 48.76 | 101.22 | 118.96 | 94.93 | 57.34 | 60.96 | 111.36 | 64.50 | 76.28 |
| VCNet [28] | 65.28 | 129.59 | 133.08 | 120.39 | 78.43 | 82.46 | 153.17 | 83.91 | 97. 85 |
| Xu [55] | 38.41 | 70.94 | 94.31 | 81.90 | 48.55 | 55.19 | 99.34 | 60.92 | 61.40 |
| **Ours** | **38.39** | **61.59** | **69.53** | **51.14** | **37.67** | **42.10** | **58.32** | **44.77** | **48.16** |

| OUTDOOR | walking | sitting | crawling | crouching | boxing | dancing | stretching | waving | total (mm) |
|---|---|---|---|---|---|---|---|---|---|
| 3DV'17 [27] | 68.67 | 114.87 | 113.23 | 118.55 | 95.29 | 72.99 | 114.48 | 72.41 | 94.46 |
| VCNet [28] | 84.43 | 167.87 | 138.39 | 154.54 | 108.36 | 85.01 | 160.57 | 96.22 | 113.75 |
| Xu [55] | 63.10 | **85.48** | 96.63 | 92.88 | 96.01 | 68.35 | 123.56 | 61.42 | 80.64 |
| **Our** | **43.60** | 85.91 | **83.06** | **69.23** | **69.32** | **45.40** | **76.68** | **51.38** | **60.19** |

Table 5: Quantitative evaluation on Mo2Cap2 dataset [55], both indoor and outdoor test-sets. Our approach outperforms all competitors by more than **21.6%** (13.24 mm) on indoor data and more than **25.4%** (20.45 mm) on outdoor data.
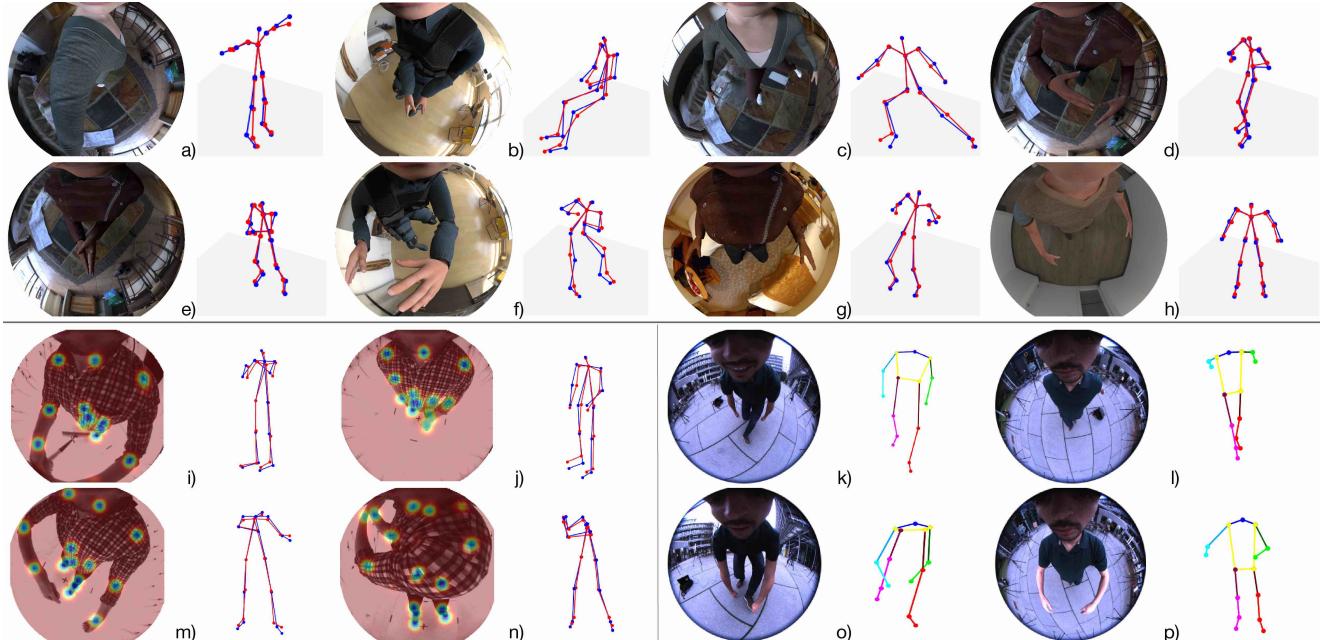


Figure 6: Qualitative results on synthetic and real images acquired with a camera physically mounted on a HMD: **(top)** 3D poses reconstructed from synthetic images. Blue are ground truth poses and red predictions; **(bottom)** reconstructed 3D predictions (in red) from **real images** captured in a mocap studio compared to ground truth poses (in blue), and reconstruction of images the wild from mo2cap2 [55] with poses shown using the same alignment for better visualization.

| Action | Error (mm) | Action | Error (mm) |
|---|---|---|---|
| Greeting | 51.78 | Upper Stretching | 61.09 |
| Talking | 47.46 | Throwing Arrow | 88.54 |
| Playing Golf | 68.74 | **Average** | **61.71** |
| Shooting | 52.64 | | |

Table 6: Average reconstruction error per joint using Eq. 3, evaluated on real data captured in a mocap studio.

| 3D | 2D | Error (mm) |
|---|---|---|
| 50% | 50% | 68.04 |
| 50% | 100% | 63.98 |

(a) $x$R-EgoPose

| Training dataset | Error (mm) |
|---|---|
| H36M | 67.9 |
| H36M + COCO + MPII | 53.4 |

(b) Human3.6M

Table 7: Having a larger corpus of 2D annotations can be leveraged to improve final 3D pose estimation

# 7. Conclusions

We have presented a solution to the problem of 3D body pose estimation from a monocular camera installed on a HMD. Our fully differentiable network estimates input images to heatmaps, and from heatmaps to 3D pose via a novel dual-branch auto-encoder which was fundamental for accurate results. We have also introduced the $x$R-EgoPose dataset, a new large scale photorealistic synthetic dataset that was essential for training and will be made publicly available to promote research in this exciting area. While our results are state-of-the-art, there are a few failures cases due to extreme occlusion and the inability of the system to measure hands when they are out of the field of view. Adding additional cameras to cover more field of view and enable multi-view sensing is the focus of our future work.

# References

[1] I. Akhter and M. J. Black. Pose-conditioned joint angle limits for 3d human pose reconstruction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1446–1455, 2015. 2

[2] M. Amer, S. V. Amer, and A. Maria. Deep 3d human pose estimation under partial body presence. In *Proceedings of the IEEE International Conference on Image Processing (ICIP)*, 2018. 3

[3] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014. 2, 7

[4] F. Bogo, A. Kanazawa, C. Lassner, P. Gehler, J. Romero, and M. J. Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In *European Conference on Computer Vision*, pages 561–578. Springer, 2016. 2

[5] C. Cao, Y. Zhang, Y. Wu, H. Lu, and J. Cheng. Egocentric gesture recognition using recurrent 3d convolutional neural networks with spatiotemporal transformer modules. *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017. 2

[6] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017. 2

[7] C.-H. Chen and D. 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. 7

[8] R. Dabral, A. Mundhada, U. Kusupati, S. Afaque, and A. Jain. Structure-aware and temporally coherent 3d human pose estimation. *arXiv preprint arXiv:1711.09250*, 2017. 7

[9] D. Drover, C.-H. Chen, A. Agrawal, A. Tyagi, and C. P. Huynh. Can 3d pose be learned from 2d projections alone? *arXiv preprint arXiv:1808.07182*, 2018. 2

[10] H.-S. Fang, Y. Xu, W. Wang, X. Liu, and S.-C. Zhu. Learning pose grammar to encode human body configuration for 3d pose estimation. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. 7

[11] A. Fathi, A. Farhadi, and J. M. Rehg. Understanding egocentric activities. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2011. 2

[12] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256, 2010. 5

[13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4, 5

[14] U. Hess, K. Kafetsios, H. Mauersberger, C. Blaison, and C. Kessler. Signal and noise in the perception of facial emotion expressions: From labs to life. *Pers Soc Psychol Bull*, 42(8), 2016. 1

[15] M. R. I. Hossain and J. J. Little. Exploiting temporal information for 3d human pose estimation. In *European Conference on Computer Vision*, pages 69–86. Springer, 2018. 7

[16] https://medium.com/@DeepMotionInc/how-to-make-3-point-tracked-full-body-avatars-in-vr-34b3f6709782. How to make 3 point tracked full-body avatars in vr, https://medium.com/@deepmotioninc/how-to-make-3-point-tracked-full-body-avatars-in-vr-34b3f6709782, last accessed on 2019-03-19. 1

[17] https://www.mixamo.com/. Animated 3d characters, last accessed on 2019-03-19. 4

[18] C. Ionescu, D. Papava, V. Olaru, and C. 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, 2014. 2, 7

[19] E. Jahangiri and A. L. Yuille. Generating multiple diverse hypotheses for human 3d pose consistent with 2d joint detections. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 805–814, 2017. 7

[20] H. Jiang and K. Grauman. Seeing invisible poses: Estimating 3d body pose from egocentric video. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3501–3509. IEEE, 2017. 3

[21] A. Kanazawa, M. J. Black, D. W. Jacobs, and J. Malik. End-to-end recovery of human shape and pose. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7122–7131, 2018. 7

[22] S. Li and A. B. Chan. 3d human pose estimation from monocular images with deep convolutional neural network. In *Asian Conference on Computer Vision*, pages 332–347. Springer, 2014. 2

[23] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 7

[24] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. Smpl: A skinned multi-person linear model. *ACM Transactions on Graphics (TOG)*, 34(6):248, 2015. 2

[25] M. Ma, H. Fan, and K. M. Kitani. Going deeper into first-person activity recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1894–1903, 2016. 2

[26] J. Martinez, R. Hossain, J. Romero, and J. J. Little. A simple yet effective baseline for 3d human pose estimation. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2017. 2, 6, 7

[27] D. Mehta, H. Rhodin, D. Casas, P. Fua, O. Sotnychenko, W. Xu, and C. Theobalt. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *3D Vision (3DV)*, 2017. 2, 7, 8

[28] D. Mehta, S. Sridhar, O. Sotnychenko, H. Rhodin, M. Shafiei, H.-P. Seidel, W. Xu, D. Casas, and C. Theobalt. Vnect: Real-time 3d human pose estimation with a single rgb camera. *ACM Transactions on Graphics (TOG)*, 36(4):44, 2017. 8

[29] F. Moreno-Noguer. 3d human pose estimation from a single image via distance matrix regression. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1561–1570. IEEE, 2017. 2, 7

[30] A. Newell, K. Yang, and J. Deng. Stacked hourglass networks for human pose estimation. In *European Conference on Computer Vision*, pages 483–499. Springer, 2016. 2, 5

[31] S. Park, J. Hwang, and N. Kwak. 3d human pose estimation using convolutional neural networks with 2d pose information. In *European Conference on Computer Vision, Workshops*, pages 156–169. Springer, 2016. 2

[32] G. Pavlakos, X. Zhou, K. G. Derpanis, and K. Daniilidis. Coarse-to-fine volumetric prediction for single-image 3d human pose. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 1263–1272. IEEE, 2017. 2

[33] G. Pavlakos, L. Zhu, X. Zhou, and K. Daniilidis. Learning to estimate 3d human pose and shape from a single color image. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 5

[34] L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. V. Gehler, and B. Schiele. Deepcut: Joint subset partition and labeling for multi person pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4929–4937, 2016. 2

[35] V. Ramakrishna, T. Kanade, and Y. Sheikh. Reconstructing 3d human pose from 2d image landmarks. In *European Conference on Computer Vision*, pages 573–586. Springer, 2012. 2

[36] J. T. Reason and J. J. Brand. *Motion sickness*. Academic press, 1975. 1

[37] H. Rhodin, C. Richardt, D. Casas, E. Insafutdinov, M. Shafiei, H.-P. Seidel, B. Schiele, and C. Theobalt. Egocap: egocentric marker-less motion capture with two fisheye cameras. *ACM Transactions on Graphics (TOG)*, 35(6):162, 2016. 3, 7

[38] H. Rhodin, M. Salzmann, and P. Fua. Unsupervised geometry-aware representation for 3d human pose estimation. *arXiv preprint arXiv:1804.01110*, 2018. 2

[39] G. Rogez and C. Schmid. Mocap-guided data augmentation for 3d pose estimation in the wild. In *Advances in Neural Information Processing Systems*, pages 3108–3116, 2016. 2, 7

[40] G. Rogez, J. S. Supancic, and D. Ramanan. First-person pose recognition using egocentric workspaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4325–4333, 2015. 3

[41] G. Rogez, P. Weinzaepfel, and C. Schmid. Lcr-net: Localization-classification-regression for human pose. In *CVPR 2017-IEEE Conference on Computer Vision & Pattern Recognition*, 2017. 2

[42] G. Rogez, P. Weinzaepfel, and C. Schmid. LCR-Net++: Multi-person 2D and 3D Pose Detection in Natural Images. *CoRR*, abs/1803.00455v1, 2018. 2

[43] M. Sanzari, V. Ntouskos, and F. Pirri. Bayesian image based 3d pose estimation. In *European Conference on Computer Vision*, pages 566–582. Springer, 2016. 2

[44] T. Shiratori, H. S. Park, L. Sigal, Y. Sheikh, and J. K. Hodgins. Motion capture from body-mounted cameras. In *ACM Transactions on Graphics (TOG)*, volume 30, page 31. ACM, 2011. 3

[45] X. Sun, J. Shang, S. Liang, and Y. Wei. Compositional human pose regression. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2602–2611, 2017. 7

[46] X. Sun, B. Xiao, F. Wei, S. Liang, and Y. Wei. Integral human pose regression. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 529–545, 2018. 2, 7

[47] B. Tekin, I. Katircioglu, M. Salzmann, V. Lepetit, and P. Fua. Structured prediction of 3d human pose with deep neural networks. In *British Machine Vision Conference (BMVC)*, 2016. 2

[48] D. Tome, C. Russell, and L. Agapito. Lifting from the deep: Convolutional 3d pose estimation from a single image. *CVPR 2017 Proceedings*, pages 2500–2509, 2017. 2, 7

[49] H.-Y. Tung, H.-W. Tung, E. Yumer, and K. Fragkiadaki. Self-supervised learning of motion capture. In *Advances in Neural Information Processing Systems*, pages 5242–5252, 2017. 2

[50] H.-Y. F. Tung, A. Harley, W. Seto, and K. Fragkiadaki. Adversarial inversion: Inverse graphics with adversarial priors. *arXiv preprint arXiv:1705.11166*, 2017. 2

[51] G. Varol, J. Romero, X. Martin, N. Mahmood, M. J. Black, I. Laptev, and C. Schmid. Learning from synthetic humans. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2017)*, 2017. 2

[52] T. von Marcard, B. Rosenhahn, M. J. Black, and G. Pons-Moll. Sparse inertial poser: Automatic 3d human pose estimation from sparse imus. In *Computer Graphics Forum*, volume 36, pages 349–360. Wiley Online Library, 2017. 3

[53] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh. Convolutional pose machines. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4724–4732, 2016. 2, 5

[54] J. Wu, T. Xue, J. J. Lim, Y. Tian, J. B. Tenenbaum, A. Torralba, and W. T. Freeman. Single image 3d interpreter network. In *European Conference on Computer Vision*, pages 365–382. Springer, 2016. 2

[55] W. Xu, A. Chatterjee, M. Zollhoefer, H. Rhodin, P. Fua, H.-P. Seidel, and C. Theobalt. Mo$^2$Cap$^2$ : Real-time mobile 3d motion capture with a cap-mounted fisheye camera. *IEEE Transactions on Visualization and Computer Graphics*, pages 1–1, 2019. 2, 3, 6, 8

[56] H. Yasin, U. Iqbal, B. Kruger, A. Weber, and J. Gall. A dual-source approach for 3d pose estimation from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4948–4956, 2016. 7

[57] H. Yonemoto, K. Murasaki, T. Osawa, K. Sudo, J. Shimamura, and Y. Taniguchi. Egocentric articulated pose tracking for action recognition. In *International Conference on Machine Vision Applications (MVA)*, 2015. 3

[58] X. Zhou, X. Sun, W. Zhang, S. Liang, and Y. Wei. Deep kinematic pose regression. In *European Conference on Computer Vision*, pages 186–201. Springer, 2016. 2

[59] X. Zhou, M. Zhu, S. Leonardos, and K. Daniilidis. Sparse representation for 3d shape estimation: A convex relaxation approach. *IEEE transactions on pattern analysis and machine intelligence*, 39(8):1648–1661, 2017. 2

[60] X. Zhou, M. Zhu, S. Leonardos, K. G. Derpanis, and K. Daniilidis. Sparseness meets deepness: 3d human pose estimation from monocular video. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4966–4975, 2016. 2

[61] X. Zhou, M. Zhu, G. Pavlakos, S. Leonardos, K. G. Derpanis, and K. Daniilidis. Monocap: Monocular human motion capture using a cnn coupled with a geometric prior. *IEEE transactions on pattern analysis and machine intelligence*, 2018. 7