# Metric-Scale Truncation-Robust Heatmaps for 3D Human Pose Estimation

István Sárándi[1], Timm Linder[2], Kai O. Arras[2] and Bastian Leibe[1]

[1]Computer Vision Group, Visual Computing Institute, RWTH Aachen University, Germany
[2]Robert Bosch GmbH, Corporate Research, Renningen, Germany

{sarandi,leibe}@vision.rwth-aachen.de, {timm.linder,kaioliver.arras}@de.bosch.com

*Abstract*— Heatmap representations have formed the basis of 2D human pose estimation systems for many years, but their generalizations for 3D pose have only recently been considered. This includes 2.5D volumetric heatmaps, whose X and Y axes correspond to image space and the Z axis to metric depth around the subject. To obtain metric-scale predictions, these methods must include a separate, explicit post-processing step to resolve scale ambiguity. Further, they cannot encode body joint positions outside of the image boundaries, leading to incomplete pose estimates in case of image truncation. We address these limitations by proposing metric-scale truncation-robust (*MeTRo*) volumetric heatmaps, whose dimensions are defined in metric 3D space near the subject, instead of being aligned with image space. We train a fully-convolutional network to estimate such heatmaps from monocular RGB in an end-to-end manner. This reinterpretation of the heatmap dimensions allows us to estimate complete metric-scale poses without test-time knowledge of the focal length or person distance and without relying on anthropometric heuristics in post-processing. Furthermore, as the image space is decoupled from the heatmap space, the network can learn to reason about joints beyond the image boundary. Using ResNet-50 without any additional learned layers, we obtain state-of-the-art results on the Human3.6M and MPI-INF-3DHP benchmarks. As our method is simple and fast, it can become a useful component for real-time top-down multi-person pose estimation systems. We make our code publicly available to facilitate further research.[1]

## I. INTRODUCTION

Human pose estimation from camera input is a long-standing problem in computer vision with a wide range of applications including human-robot interaction [56], virtual reality [1], medicine [3], [43] and commerce [29]. Since the adoption of deep convolutional neural networks (CNN), and especially heatmap representations, we have witnessed rapid progress in pose estimation research [30], [52], [17]. Recently, deep CNNs have been successfully applied to the monocular 3D human pose estimation task as well [25], [26], [55], [21], [32]. Here a person's anatomical landmarks are sought in 3D space, *i.e.*, in millimeters, instead of pixels. These advances tie into one of the major themes of computer vision research, reconstructing 3D structure from images. Such tasks are especially challenging due to inherent geometric ambiguities. One class of ambiguities arise because different 3D articulations may share the same 2D projection. Another ambiguity is between the size of an object and its distance, since small objects near the camera look the same as large ones far away.
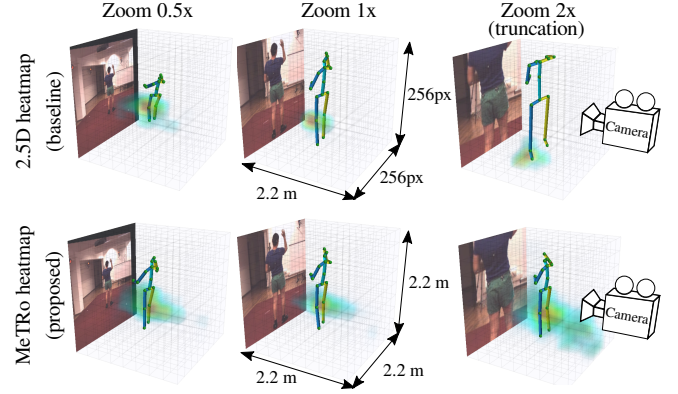


Fig. 1. By defining heatmaps in the 3D metric space around the person (*bottom row*) we can directly predict scale-correct and complete poses. This is in contrast to prior work (*top row*) that defines the X and Y heatmap axes in image space and requires further post-processing to obtain a metric-scale skeleton. The three columns show how zooming affects the heatmap representation (a knee heatmap is shown along with the soft-argmax decoded skeleton). Notice that our heatmap-space representation is largely invariant to image scaling and estimates a complete pose even under body-truncation at the image boundaries.

There is no clear consensus yet about the most effective way to represent and tackle these problems. One promising line of approaches extend 2D joint heatmaps with a depth axis, resulting in a 2.5D volumetric representation [35], [46], [15], [23]. Finding heatmap maxima gives the estimated pixel coordinates and root-relative depths per joint (a 2.5D pose). While these estimates can be highly accurate, the 2.5D representation does not address the challenging ambiguity between scale (person size) and distance. Indeed, to bridge the gap between a 2.5D and a 3D pose, one needs to perform scale recovery as a separate post-processing step. Multiple explicit anthropometric heuristics have been proposed as scale cues, *e.g.* bone length priors [35] or a skeleton length prior [44], computed by averaging over the training poses. However, these simple heuristics have difficulties when the experimental subjects have diverse heights. A further limitation is that 2.5D formulations are constrained to the estimation of joints that lie within the image boundaries. This can be problematic in practical applications with noisy bounding box detectors. While one could use an additional module to estimate missing joints, it is preferable to learn the complete skeleton estimation in a single unified stage.

Our goal in this paper is to tackle the above limitations in a simple and efficient manner, while keeping the structural advantages of fully-convolutional heatmap estimation, as opposed to numerical coordinate regression. To this end,

---

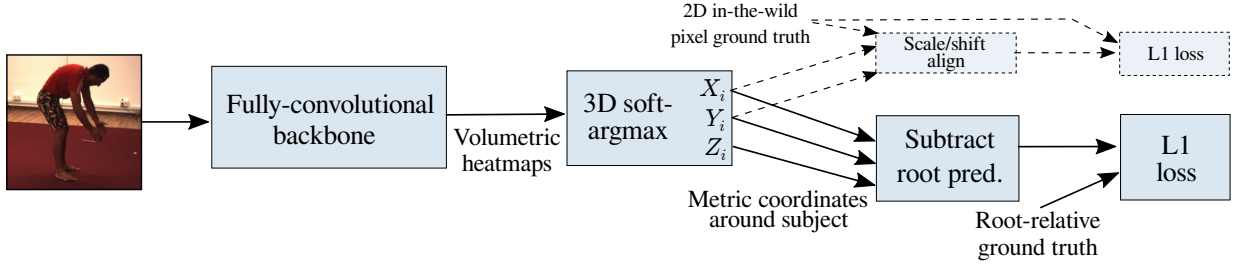[1]https://vision.rwth-aachen.de/metro-pose3d

Fig. 2. Overview of the method. We predict volumetric heatmaps using an off-the-shelf fully-convolutional backbone. Applying soft-argmax on these heatmaps and scaling by an image-independent constant factor yields joint coordinates in metric space up to translation. We minimize the root-relative $L^1$ loss. Focusing on simplicity, no learnable parameters are introduced outside the standard backbone. Note that reasoning about truncated body parts, scale-recovery and back-projection also happen implicitly within the backbone. Weak supervision from in-the-wild 2D-labeled data is incorporated by aligning the metric prediction to the 2D ground truth by scaling and translation and computing the $L^1$ loss (dashed arrows and boxes).

we propose training a fully-convolutional network to output metric-scale truncation-robust (MeTRo) heatmaps as illustrated in Fig. 1. All dimensions of these heatmaps are defined to have a fixed metric extent in meters. This is an unconventional task definition for fully-convolutional networks (FCN). FCNs are predominantly applied for pixel-wise prediction tasks, such as semantic segmentation, where the input and output are pixel-to-pixel aligned, or at least are in the same coordinate frame. In our proposed approach, the input pixel positions and the output metric positions only satisfy a looser form of spatial correspondence. Nevertheless, we show that somewhat surprisingly, such a mapping can still be learned effectively by a standard modern FCN backbone.

While explicit prior knowledge of problem structure is known to be beneficial, it is still an open question how much geometric computation needs to be performed explicitly and how much can be learned by deep networks from data. By skipping the 2.5D stage, we train the backbone FCN to implicitly reason about out-of-image joints, discover scale cues and learn the geometric perspective back-projection in an end-to-end manner. Our MeTRo heatmap representation can naturally encode body parts lying outside the image, since the prediction volume's bounds do not correspond to the image bounds. As there is no need to design an explicit scale recovery step, the pipeline becomes simpler and requires neither the focal length nor the root joint distance to be known at test time.

Recent approaches have achieved good generalization performance to in-the-wild images by using abundant and diverse images with 2D pose labels in the training procedure besides 3D data [55], [46], [23]. Applying such weak supervision is challenging in our representation, since the network does not make any pixel-based predictions, its outputs are directly on a metric scale. We tackle this by proposing a scale and translation invariant loss computation method for 2D-annotated examples using an alignment layer. Combined with the recently introduced differentiable soft-argmax [18], [24], [46], [31] layer, our method becomes end-to-end learned all the way from image to final 3D metric-scale prediction as shown in Fig. 2. Soft-argmax also allows rapid training with low-resolution heatmaps and using dense prediction with smaller strides at test time for higher quality results, without the need for a decoder module. Here we find that the

details of the striding mechanism are crucial and propose a "centered striding" method that distributes the output neuron receptive fields evenly over the image. Experimentally, our MeTRo heatmap estimation achieves state-of-the-art results on the two largest 3D pose benchmarks, Human3.6M and MPI-INF-3DHP. To isolate the effect of the representation, we perform direct comparisons with 2.5D heatmap learning using bone-length-based scale recovery [35], under otherwise equal training conditions. We find that scale cues can indeed be learned implicitly in this fashion and MeTRo outperforms the baseline on most test sequences.

## II. RELATED WORK

3D human pose estimation has had a long research history starting with hand-crafted features and part-based models [40]. Similar to other computer vision problems, the transition to deep convolutional networks has led to a dramatic performance increase in this task as well [48], [27], [26], [45], [25], [28], [46].

### A. Deep 3D Human Pose Estimation

Much of the inspiration in recent 3D pose estimator design has come from lessons learned in 2D pose research. DeepPose, the first neural method for 2D pose estimation [49] directly regressed 2D body joint coordinates on the RGB input via convolutional and fully-connected layers. Later top-performing methods have transitioned to predicting body joint heatmaps by fully-convolutional networks (*e.g.*, [30]) as an intermediate representation. These heatmaps are spatially discretized arrays (one for each joint), in which higher values indicate higher confidence that the particular joint is located at the corresponding position.

One line of 3D pose research builds on top of 2D heatmaps and infers the 3D pose from them by exemplar-matching [4], regression [25] or probabilistic inference [48]. One downside of such approaches is that the image content only indirectly influences the 3D estimation, as it acts on the result of the 2D estimation stage. Furthermore, 2D-to-3D lifting is performed in a numerical coordinate representation, which does not benefit from the built-in convolutional structure of CNNs.

Nibali *et al.* [32] predict three marginal heatmaps per body joint, for the XY, XZ and YZ planes, respectively. Pavlakos *et al.* have proposed extending 2D heatmaps with a root-relative

metric depth axis [35]. One can obtain the 2D pixel positions and root-relative depths of the joints by finding maxima in the heatmaps.

One downside of heatmap representations has been the requirement of a dense output, which can become especially costly in 3D. The recently proposed soft-argmax [18], [24], [31] *a.k.a.* integral regression [46] method greatly alleviates this problem. As opposed to hard-argmax, which simply finds the location of the highest heatmap activation, soft-argmax is computed as the weighted average of all voxel grid coordinates, using softmaxed heatmap activations as the weights. For example, a low resolution heatmap can encode a joint position lying halfway between two bin centers by outputting 0.5 for both bins. By virtue of being differentiable unlike hard-argmax, it also obviates the need for explicit heatmap-level supervision (*e.g.*, voxel-wise cross-entropy). Instead, the loss can be computed (and its gradients backpropagated) from the coordinates yielded by soft-argmax.

Besides 2D heatmaps, Mehta *et al.* estimate three further output channels per joint, the so-called *location maps* [26]. These are read out at the position of the corresponding heatmap's peak to obtain the X, Y and Z coordinates on a metric scale. Note how in this approach the final 3D joint coordinates are generated in the form of activation *values* (of the location maps at the heatmap peaks), as opposed to high-activation *locations*. We can thus think of it a conceptual hybrid of direct numerical coordinate regression and heatmap estimation. A downside of this method is that it requires high-resolution location maps and cannot benefit from the soft-argmax approach.

*B. Scale Ambiguity*

It is well-known that projecting a 3D world onto a 2D image plane results in ambiguity between size and distance (depth). However, the end goal for 3D scene understanding and 3D human pose estimation in particular is a metric-space output at the true scale. The ambiguity can only be resolved using semantic scale cues, *i.e.* prior knowledge of the usual size of humans and other objects appearing in the scene. Unfortunately, not all papers describe how this step is performed. Some authors report their results assuming a known focal length and known ground-truth root joint distance [32], [46], [47], [6] and leave their estimation as a separate task. A simple anthropometric approach is used by Pavlakos *et al.* Given 2D pixel positions and root relative depth estimates from volumetric heatmaps, they optimize the absolute person distance such that the back-projected skeleton's bone lengths match the average over the training set in a least squares sense [35]. A detailed description of this convex optimization problem is given in [36]. We use this scale recovery approach as our main baseline comparison throughout the paper. Sun *et al.* employ a similar idea, but use the overall skeleton length and a weak perspective model instead [44]. Some recent works have shown that direct regression of person height from an image is a challenging task [10], [7]. Véges *et al.* make use of a monocular depth

prediction network pretrained on various indoor and outdoor datasets to help with absolute person distance estimation [50].

*C. Truncated Pose Estimation*

Single-person 3D human pose estimation benchmarks, such as Human3.6M [13], [14], assume that the input is a tight crop around a whole person. In practical applications, however, we need to obtain the bounding box using imperfect person detectors, which may result in body truncation. Performance under truncation has not been studied extensively in the literature. Vosoughi *et al.* created randomly truncated crops from Human3.6M images, and showed that current methods perform poorly on truncated person images, even when only considering the present (within-boundary) joints [51]. They tackled the problem using direct numerical coordinate regression, similar to early 2D pose estimation methods [49]. In this paper, we show that our approach performs significantly better on the truncated task. Other methods, such as LCR-Net [39], can also produce out-of-image predictions, but this aspect has not been explicitly evaluated by its authors.

## III. APPROACH

Given an input RGB image crop $I \in \mathbb{R}^{w \times h \times 3}$ depicting a person, we aim to predict a (root-relative) 3D skeleton, consisting of $J$ joint coordinates $\{(X_j, Y_j, Z_j)^T\}_{j=1}^{J}$ at metric scale (*i.e.* in millimeters).

*A. Metric-Space Volumetric Heatmap Representation*

First, we apply an off-the-shelf fully-convolutional backbone with effective stride $s$ to produce $d \cdot J$ spatial output channels, where $d$ is the number of discretization bins along the depth axis of the prediction volume.

We then split the resulting array along the channel axis into $J$ volumes, each of shape $(w/s) \times (h/s) \times d$. 3D spatial softmax is applied over each of them, resulting in volumetric heatmap activations $V^{(j)} \in \mathbb{R}^{(w/s) \times (h/s) \times d}$. The 3D joint coordinates are then decoded using the soft-argmax technique with *fixed* scaling factors:

$$\begin{bmatrix} X_j \\ Y_j \\ Z_j \end{bmatrix} = \sum_{p,q,r} V_{p,q,r}^{(j)} \begin{bmatrix} p \cdot s/w \cdot W \\ q \cdot s/h \cdot H \\ r/d \cdot D \end{bmatrix}, \quad (1)$$

where the $p, q, r$ are 0-based integer indices into the volumetric heatmap array and $W, H, D$ are the fixed metric width, height and depth extents of the full prediction volume. We set these extents as 2.2 meters in our work, which allows capturing people of usual height even in a stretched out pose. The final root-relative prediction is obtained by subtracting the predicted root coordinates from all joint positions. Supervision is applied on these root-relative coordinates. Crucially, the position of the root joint prediction within the volume is not explicitly prescribed for the network, the gradients are backpropagated through the root-joint-subtraction operation. No camera calibration-based back-projection, nor bone or skeleton size-based rescaling is needed. The network is trained to perform these operations within the backbone.
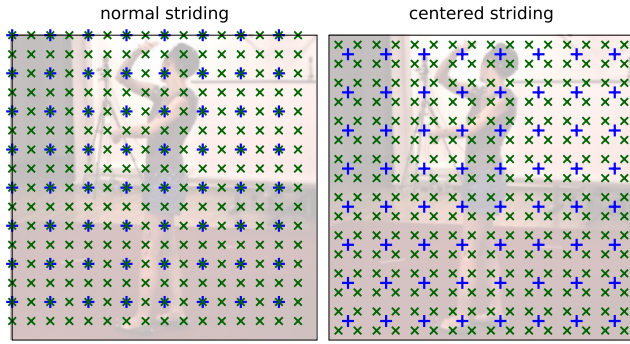
Fig. 3. Receptive field centers of the output neurons in a strided FCN on a 256x256 px input image (+: stride 32, ×: stride 16). *Left:* Normal striding logic, where the top left result is kept per 2x2 block. Consequently, the receptive field centers are not symmetrically distributed and dense prediction introduces bias. *Right:* We use centered striding by reversing the stride logic in the last strided layer (*i.e.*, bottom right result taken, instead of top left). This way the receptive fields are symmetrically distributed over the image and dense prediction at test-time introduces new bins in a proportional manner around each training-time bin.

## B. Architecture

In contrast to prior work that employs decoders with upsampling layers and multiple refinement stages with intermediate supervision, we show that the task can be tackled in a significantly simpler fashion. Indeed, we apply the widely used ResNet-50 [12] architecture to predict spatial heatmaps, without any additional learnable layers, such as transposed convolutions. ResNet-50 has an effective stride of 32, resulting in heatmaps of spatial size $8 \times 8$ from the input image of size $256 \times 256$ during training. The depth of the volume is set to 8.

## C. Centered Striding for Dense Prediction

At test time we apply the trained network with an effective stride of 4, to obtain heatmaps with spatial size 64, which is the same size as in [46] and [35]. This is called dense prediction and is commonly used in image segmentation [5]. In this technique, striding is removed from a given number of convolutional layers and the dilation rate of subsequent convolutions is increased correspondingly. To avoid a mismatch between the distribution of the heatmap neuron receptive field centers between training and test time, we apply a slight modification to the striding logic. The first column of Fig. 3 shows the usual case of a 256x256 input image processed with a training stride 32 (+) and test stride 16 (×). Clearly, the coverage changes significantly between training and test and is not symmetric over the image. This is because each convolutional layer with stride 2 returns the *top left* output for each 2x2 block. To tackle the issue, we propose *centered striding* (second column in Fig. 3), where the last strided convolutional layer of the backbone is "reversed", such that it outputs the *bottom right* result per each 2x2 block. The result is a more evenly distributed coverage over the image, without changing the resolution of either the input or the output. This benefit is evaluated in Section V.

## D. Scale and Translation Invariant Loss for 2D Supervision

Similar to recent approaches [55], [46], [23], we train simultaneously on 3D-labeled data from motion capture studios and 2D-labeled, in-the-wild data from the MPII dataset [2], to incorporate more appearance variability in the training process. Only the arm and leg joints are used from MPII, since we found these to be the most consistently labeled across datasets. Half of each mini-batch is filled with examples of either kind. Supervision via 2D labels is straightforward when using 2.5D heatmaps, as the $X$ and $Y$ heatmap axes correspond to the space in which the 2D labels are defined. However, since our prediction volume is defined on a metric scale and is not aligned with image space, we propose a 2D loss computation method that is invariant to prediction scale and translation. To this end, we first orthographically project the predicted skeleton onto the image plane by discarding the Z coordinate. Then we align the projected prediction to the 2D pixel-scale ground truth by translation and uniform scaling to the least-squares optimal fit before computing the loss. This alignment layer is differentiable and gradients can be backpropagated through it, in a similar manner to batch normalization layers. We note that a similar scale-invariant loss has been used by Rhodin *et al.* to enforce multi-view consistency of 3D poses [38].

## E. Estimation of Truncated Poses

Our metric-space heatmap representation decouples the image boundary from the heatmap boundary. This enables the prediction of joint locations outside the image frame without additional design effort, the network is simply trained to output complete poses at a metric scale, regardless of how the input image is scaled or cropped. To evaluate this aspect, we follow Vosoughi *et al.* [51] by randomly cropping H36M inputs, keeping at least 1/4 of the area of the person bounding square. Examples of such crops are in the second row of Fig. 4. We consider two scenarios. In the first one, the above described sampling of truncated crops is only performed at test time. In the second case, such crops are used for training as well.

## F. Training Details

*1) Loss:* Prior work has shown that the $L^1$ loss is preferable in soft-argmax-based pose estimation [46]. To balance the losses computed on 3D and 2D examples, we use a fixed weighting factor tuned on a separate validation set of Human3.6M, yielding the overall loss as $\mathcal{L} = \mathcal{L}_{3D}^1 + \lambda \mathcal{L}_{2D}^1$.

*2) Training Schedule:* We initialize the network with ImageNet-pretrained weights and use the Adam optimizer with weight decay [20] and a batch size of 64. We decay the learning rate exponentially by an overall factor of 100, in two parts: from $10^{-4}$ to $3.33 \times 10^{-5}$ over 25 epochs and from $3.33 \times 10^{-6}$ to $10^{-6}$ in 2 final cooldown epochs.

*3) Randomness:* As usual in deep learning, several sources of randomness influence the exact results of an experiment: random weight initialization, data shuffling, data augmentation and hardware-level non-determinism of execution order. We
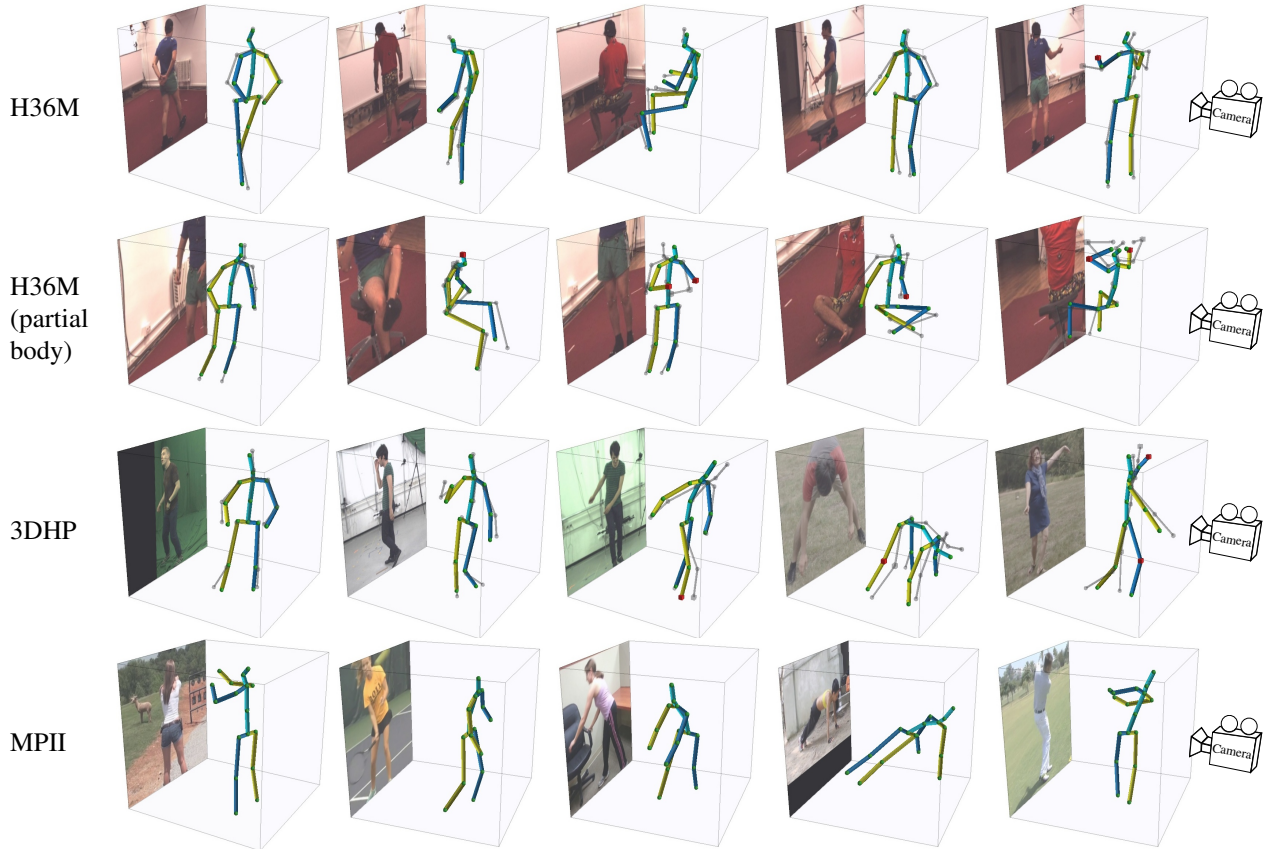
Fig. 4. Qualitative results of our method on different datasets. Predictions are shown in color, ground truth in gray (except for MPII, where it is unavailable). Green spheres mark predictions within 150 mm of the ground truth, red cubes beyond that threshold. *Best viewed in color.*

control these (except the last) by consistently seeding the random number generators. To distinguish random fluctuations from algorithmic differences, we repeat our main experiments with 5 different seeds and report the mean and standard deviation of the evaluation metrics.

## IV. DATASETS AND PREPROCESSING

We conduct experiments on the largest 3D pose estimation benchmarks: Human3.6M (shortened as H36M) [13], [14] and MPI-INF-3DHP (3DHP) [27].

H36M [13], [14] was captured with 4 cameras in a motion capture studio. Two evaluation protocols have been established over the years. In Protocol 1, the training subjects are 1, 5, 6, 7, 8, while 9 and 11 are used for testing. Prediction and ground-truth are aligned at the root joint, but no Procrustes alignment is performed. In Protocol 2, subjects 1, 5, 6, 7, 8, 9 are used in training and 11 in evaluation, with Procrustes alignment between prediction and ground truth. Every 64th frame is evaluated, as in prior work.

3DHP [27] shows 8 training subjects in a green-screen studio. Test frames come from 3 scenes, each with 2 subjects: green-screen studio, studio without green screen, and outdoor. The latter two make this benchmark more challenging than H36M. In this dataset, the hip and pelvis joints are labeled closer to the legs than in MPII. We follow [55] and move these joints towards the neck by a fifth of the pelvis-neck vector before comparing with MPII-annotated skeletons for 2D loss computation. 3DHP provides two ground truth variants: usual metric-space poses and "universal" (height-normalized) ones. To analyze scale recovery performance, we use metric-scale evaluation, but to be comparable with prior work we also provide results with universal skeletons.

We downsample the videos from 50 to 10 fps. To further reduce redundancy, frames are only kept for training if at least one body joint moves at least 100 mm since the previous kept frame. For 3DHP, we train on images from chest-height cameras as [27], and only on examples where all joints are within the image.

For H36M examples we use the provided bounding boxes. The 3DHP dataset provides no boxes, we therefore generate them ourselves by combining the bounding box of the labeled joint positions and the most confident person detection from YOLOv3 [37]. For the 2D examples of MPII, we use the provided rough center positions and person sizes as the center and side length of the box, respectively.

We crop the image to the person's bounding square and resize it to $256 \times 256$ px. Perspective effects must be taken into account when centering the image on the subject as this induces an implicit rotation of the camera [27]. We compensate for this effect by transforming image and the target joint positions to match the rotated camera frame. The green-screen 3DHP sequences are gamma-adjusted with an exponent of 0.67.

We apply geometric augmentations (scaling, rotation, translation, horizontal flip) and color distortions (brightness,

TABLE I

COMPARISON ON H36M PROTOCOL 1, USING MEAN PER JOINT POSITION ERROR (MPJPE) WITHOUT PROCRUSTES ALIGNMENT. WE GIVE MEAN AND STANDARD DEVIATION OF THE OVERALL METRIC FOR 5 DIFFERENT RANDOM SEEDS. ALL METHODS USE EXTRA 2D POSE DATA IN TRAINING.

| | Dir. | Dis. | Eat | Gre. | Pho. | Pose | Pur. | Sit | SitD | Sm. | Pho. | Wait | Walk | WD | WT | Avg ↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Methods using ground-truth scale or depth information at test time* | | | | | | | | | | | | | | | | |
| Sun *et al.* [45] | 52.8 | 54.8 | 54.2 | 54.3 | 61.8 | 53.1 | 53.6 | 71.7 | 86.7 | 61.5 | 67.2 | 53.4 | 47.1 | 61.6 | 53.4 | 59.1 |
| Nibali *et al.* [32] | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | 57.0 |
| Luvizon *et al.* [23] | 51.5 | 53.4 | 49.0 | 52.5 | 53.9 | 50.3 | 54.4 | 63.6 | 73.5 | 55.3 | 61.9 | 50.1 | 46.0 | 60.2 | 51.0 | 55.1 |
| Sun *et al.* [46] | 47.5 | **47.7** | 49.5 | 50.2 | 51.4 | 43.8 | 46.4 | 58.9 | 65.7 | 49.4 | 55.8 | 47.8 | **38.9** | **49.0** | 43.8 | 49.6 |
| Chen *et al.* [6] | **45.3** | 49.8 | 46.1 | 49.6 | **48.2** | **41.7** | 47.4 | **53.1** | **55.2** | **48.0** | 57.7 | 45.6 | 40.8 | 52.4 | 45.2 | **48.4** |
| *Methods using no ground truth scale or depth information at test time* | | | | | | | | | | | | | | | | |
| Pavlakos *et al.* [35] | 67.4 | 72.0 | 66.7 | 69.1 | 72.0 | 77.0 | 65.0 | 68.3 | 83.7 | 96.5 | 71.7 | 65.8 | 74.9 | 59.1 | 63.2 | 71.9 |
| Zhou *et al.* [55] | 54.8 | 60.7 | 58.2 | 71.4 | 62.0 | 53.8 | 55.6 | 75.2 | 111.6 | 64.2 | 65.5 | 66.0 | 51.4 | 63.2 | 55.3 | 64.9 |
| Martinez *et al.* [25] | 51.8 | 56.2 | 58.1 | 59.0 | 69.5 | 55.2 | 58.1 | 74.0 | 94.6 | 62.3 | 78.4 | 59.1 | 49.5 | 65.1 | 52.4 | 62.9 |
| Fang *et al.* [9] | 50.1 | 54.3 | 57.0 | 57.1 | 66.6 | 53.4 | 55.7 | 72.8 | 88.6 | 60.3 | 73.3 | 57.7 | 47.5 | 62.7 | 50.6 | 60.4 |
| Yang *et al.* [53] | 51.5 | 58.9 | 50.4 | 57.0 | 62.1 | 49.8 | 52.7 | 69.2 | 85.2 | 57.4 | 65.4 | 58.4 | 43.6 | 60.1 | 47.7 | 58.6 |
| Pavlakos *et al.* [34] | 48.5 | 54.4 | 54.4 | 52.0 | 59.4 | 49.9 | 52.9 | 65.8 | 71.1 | 56.6 | 65.3 | 52.9 | 44.7 | 60.9 | 47.8 | 56.2 |
| Liu *et al.* [19] | 47.0 | 53.1 | 50.3 | 48.8 | 56.0 | 48.1 | 47.6 | 65.9 | 72.6 | 52.3 | 61.4 | 49.1 | **39.3** | 54.2 | **40.6** | 52.4 |
| 2.5D mean bone len. | 45.1 | 50.4 | 45.4 | **47.8** | **50.0** | 44.6 | 49.8 | 59.0 | 69.4 | 49.4 | 56.5 | 48.0 | 39.6 | **49.4** | 45.0 | 50.2±0.3 |
| MeTRo (proposed) | 46.3 | **48.3** | 43.3 | 48.2 | 50.2 | 45.1 | **46.1** | 56.2 | 66.8 | 49.3 | 54.5 | **46.7** | 40.1 | 49.6 | 46.2 | **49.3±0.7** |

contrast, hue, saturation). Synthetic occlusion is added with 70% probability, half of which are rectangles with uniform white noise as in [54], half are segmented non-person objects from the Pascal VOC dataset [8] as in [41], [42]. On the 3DHP dataset we also apply background augmentation with 70% probability following [27], but no compositing for clothes and chair. The backgrounds are taken from the INRIA Holidays dataset [16] excluding person images. We do not use ensembling or test-time augmentation, all evaluation is done on a single crop.

We use the standard metrics from the literature. The main metric on 3DHP is the percentage of correct keypoints (PCK), *i.e.* the fraction of joints predicted within a certain distance of the ground truth (150 mm by convention). The AUC metric is the area under the PCK curve as the threshold ranges from 0 to 150 mm. The metric on H36M is mean per joint position error (MPJPE).

## V. RESULTS

We achieve state-of-the-art performance on H36M with 49.3 mm MPJPE in the scenario where no ground truth information (focal length, root joint distance) is allowed to be accessed at test-time (see Table I). This is only surpassed by Chen et al.'s [6] method (48.4), however they do use the ground truth root joint depth for back-projection at test-time and do not perform scale recovery. Similarly, Sun *et al.* [46] obtain comparable results (49.6), however they also access the ground-truth root joint depth at test time, for image cropping [47]. Besides simplifying the prediction pipeline and allowing for truncation-robust prediction (see below), our metric heatmap representation also performs better than the 2.5D baseline with bone-length-based scale recovery under the same conditions. On Protocol 2 (Table II), the benefit of our method is masked by the use of Procrustes alignment, which explicitly ignores the quality of scale recovery. It is therefore unsurprising that our method performs about equally

well as the 2.5D variant (within the standard deviation of repeated experiments).

On 3DHP, our method outperforms prior work by a large margin, including ones trained on more datasets as well (Table III). Both with universal (height-normalized) skeletons and true metric-scale ones, the MeTRo representation outperforms the baseline due to its better performance on indoor images, where scale cues such as the size of chairs and other objects in the motion capture room can be relied on. The outdoor scenes were recorded on an empty field with no useful scale cues and the explicit bone-length-based scale recovery performs better in that scenario. Qualitative results are in Fig. 4.

We analyze scale recovery in more detail in Table IV. The 2.5D baseline using mean training bone lengths performs worse on H36M and equivalently on 3DHP than the proposed approach. Interestingly, our MeTRo approach outperforms the 2.5D baseline on H36M even when the latter uses ground truth bone lengths for each test frame (51.9).

Table VII shows that training data augmentations improve performance by a large margin.

When tested on truncated crops, our method by far outperforms prior approaches (Table V). This is true even for our default training configuration, but performance improves substantially when training on truncated images as well. Qualitative examples are in the second row of Fig. 4.

### A. Speed-Accuracy Tradeoff

Given a bounding box crop, inference only requires a single forward pass of a standard backbone. Table VI shows that 511 crops can be processed per second on an RTX 2080 Ti desktop GPU when operating on batches of 8 crops at stride 32 (the time cost of performing the detection stage is not considered). Varying the heatmap resolution using dense prediction provides diminishing returns (Table VI), showing that soft-argmax can cope with heatmaps of very coarse resolution. This means our method is attractive for use in top-down multi-person pose estimation systems.

TABLE II

COMPARISONS ON HUMAN3.6M UNDER PROTOCOL 2 WITH PROCRUSTES ALIGNMENT TO THE GROUND TRUTH.

| | Nie [33] | Pavlakos [35] | Sun [45] | Martinez [25] | Sun [46] | Nibali [32] | Habibie [11] | Chen [6] | 2.5D baseline | MeTRo (proposed) |
|---|---|---|---|---|---|---|---|---|---|---|
| P-MPJPE | 79.5 | 51.9 | 48.3 | 47.7 | 40.6 | 40.4 | 49.2 | **33.7** | 34.5±0.4 | 34.7±0.5 |

TABLE III

COMPARISON ON MPI-INF-3DHP WITH PRIOR METHODS. *EVALUATED BEFORE A FEW ANNOTATIONS WERE CHANGED IN THE DATASET. DASHES (−) REFLECT A LACK OF PUBLISHED INFORMATION. SUPERSCRIPTS INDICATE THE TRAINING DATA (FIRST CHARACTERS OF 3DHP, H36M, MPII, LSP AND COCO). WE GIVE THE MEAN AND STANDARD DEVIATION FOR 5 RUNS WITH DIFFERENT RANDOM SEEDS.

| | Stand/ Walk | Exer-cise | Sit on Chair | Cro./ Reach | On Floor | Sport | Misc. | Green Screen | No Gr.Sc. | Out-door | Total | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | PCK↑ | AUC↑ | MPJPE↓ |
| | | | | | | PCK | | | | | | | |
| *Universal, height-normalized skeletons (simplified scale recovery task)* | | | | | | | | | | | | | |
| Rogez et al. [39]* | 70.5 | 56.3 | 58.5 | 69.4 | 39.6 | 57.7 | 57.6 | – | – | – | 59.7 | 27.6 | 158.4 |
| Zhou et al.[H+M] [55]* | 85.4 | 71.0 | 60.7 | 71.4 | 37.8 | 70.9 | 74.4 | 71.7 | 64.7 | 72.7 | 69.2 | 32.5 | 137.1 |
| Mehta et al.[3+M+L+H] [26]* | 87.7 | 77.4 | 74.7 | 72.9 | 51.3 | 83.3 | 80.1 | – | – | – | 76.6 | 40.4 | 124.7 |
| Mehta et al.[3+M+L+H] [27]* | 86.6 | 75.3 | 74.8 | 73.7 | 52.2 | 82.1 | 77.5 | 84.6 | 72.4 | 69.7 | 75.7 | 39.3 | 117.6 |
| Mehta et al.[3+M+L+C] [28]* | 83.8 | 75.0 | 77.8 | 77.5 | 55.1 | 80.4 | 72.5 | – | – | – | 75.2 | 37.8 | 122.2 |
| Luo et al.[3+M+H] [21], [22] | 95.5 | 82.3 | 89.9 | 84.6 | 66.5 | 92.0 | 93.0 | – | – | – | 84.3 | 47.5 | 84.5 |
| Nibali et al.[3+M] [32] | – | – | – | – | – | – | – | – | – | – | 87.6 | 48.8 | 87.6 |
| 2.5D mean bone len.[3+M] | **95.9** | 91.9 | 88.6 | **92.8** | **77.2** | 95.1 | 92.9 | 93.1 | 90.5 | **89.1** | 91.2±0.1 | 57.0±0.3 | 72.2±0.7 |
| MeTRo (proposed)[3+M] | **95.9** | **93.2** | **91.6** | 92.7 | 76.4 | **95.9** | **93.1** | **94.4** | **91.8** | 87.9 | **91.8**±0.3 | **60.3**±0.5 | **67.6**±1.3 |
| *Metric-scale skeletons (full scale recovery task)* | | | | | | | | | | | | | |
| 2.5D mean bone len.[3+M] | 94.1 | 90.5 | 84.2 | **93.3** | **75.8** | 93.8 | **92.2** | 89.5 | 89.2 | **90.5** | 89.6±0.7 | 52.1±1.2 | **80.6**±2.1 |
| MeTRo (proposed)[3+M] | **95.0** | **90.6** | **88.7** | 90.0 | 72.0 | 93.7 | 91.6 | **91.3** | **89.4** | 87.0 | 89.6±0.5 | **52.6**±0.6 | 81.1±1.2 |

TABLE IV

COMPARISON WITH BASELINE METHODS OF SCALE RECOVERY, WITH OR WITHOUT ACCESS TO GROUND TRUTH INFORMATION. FOR 3DHP THE METRIC-SCALE (NON-UNIVERSAL) SKELETONS ARE USED HERE.

| | Uses test ground truth? | H36M MPJPE↓ | 3DHP (non-univ.) PCK↑ |
|---|---|---|---|
| 2.5D GT root depth | yes | **49.0** | **90.8** |
| 2.5D GT bone len. | yes | 51.9 | 90.3 |
| 2.5D mean train bone len. | no | 50.2 | **89.6** |
| MeTRo (proposed) | no | **49.3** | **89.6** |

TABLE V

MPJPE SCORES ON H36M UNDER TRUNCATION, EVALUATING ALL OR ONLY THE PRESENT JOINTS. *=TRAINING WAS NOT PERFORMED WITH TRUNCATED CROPS. OTHER METHODS' RESULTS ARE FROM [51].

| | Mehta*[26] | Zhou*[55] | Vosoughi [51] | **MeTRo*** | **MeTRo** |
|---|---|---|---|---|---|
| All joints | 396.4 | 400.5 | 185.0 | 124.7 | **77.8** |
| Present joints | 338.0 | 332.5 | 173.6 | 76.8 | **59.8** |

TABLE VI

TEST SPEED (CROPS PER SECOND, FPS) AND ACCURACY (MPJPE) TRADEOFF WITH THE TWO STRIDING VARIANTS FROM FIG. 3.

| | Striding variant | Test stride | | | |
|---|---|---|---|---|---|
| | | 32 | 16 | 8 | 4 |
| MPJPE | normal strides | 53.1 | 52.5 | 52.7 | 52.9 |
| | center-aligned | 50.9 | 50.2 | 50.0 | **49.3** |
| Speed (crop per sec.) | no batching | 160 | 150 | 105 | 38 |
| | batch size 8 | **511** | 475 | 292 | 92 |

TABLE VII

AUGMENTATION ABLATION ON H36M.

| Geometry | Color | Occlusion | MPJPE |
|---|---|---|---|
| ✓ | – | – | 58.0 |
| ✓ | ✓ | – | 52.8 |
| ✓ | ✓ | ✓ | **49.3** |

## VI. CONCLUSION

We proposed metric-scale truncation-robust (MeTRo) volumetric heatmaps in the context of 3D human pose estimation. These heatmaps directly represent the metric space around the person instead of being tied to the image space and can be predicted with any standard fully-convolutional network. With a modified weak supervision scheme for 2D labels, careful stride alignment considerations and strong data augmentation, we achieved state-of-the-art results on two important benchmarks: Human3.6M and MPI-INF-3DHP. In carefully controlled experiments, we showed that our approach can implicitly discover scale cues from the data and outperforms a previously proposed explicit bone length based heuristic on all test scenarios except the two outdoor sequences of MPI-INF-3DHP. Future research should consider possibilities for learning similar scale cues from large-scale outdoor data as well. Beyond scale recovery, we demonstrated the second benefit of the MeTRo representation, the prediction ("hallucination") of complete skeletons even when only a part of the body is contained in the image. Given its speed and

robustness to detection noise, we expect our approach to be useful in designing top-down multi-person pose estimation systems in the future.

## VII. ACKNOWLEDGMENTS

## REFERENCES

[1] T. Alldieck, M. Magnor, W. Xu, C. Theobalt, and G. Pons-Moll. Video based reconstruction of 3D people models. In *CVPR*, 2018.

[2] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele. 2D human pose estimation: New benchmark and state of the art analysis. In *CVPR*, 2014.

[3] V. Belagiannis, X. Wang, H. B. B. Shitrit, K. Hashimoto, R. Stauder, Y. Aoki, M. Kranzfelder, A. Schneider, P. Fua, S. Ilic, H. Feussner, and N. Navab. Parsing human skeletons in an operating room. *Machine Vision and Applications*, 2016.

[4] C. Chen and D. Ramanan. 3D human pose estimation = 2D pose estimation + matching. In *CVPR*, 2017.

[5] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. *PAMI*, 40(4), 2018.

[6] Z. Chen, Y. Guo, Y. Huang, and L. Wang. Learning depth-aware heatmaps for 3D human pose estimation in the wild. In *BMVC*, 2019.

[7] A. Dantcheva, F. Bremond, and P. Bilinski. Show me your face and I will tell you your height, weight and body mass index. In *ICPR*, 2018.

[8] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. http://host.robots.ox.ac.uk/pascal/VOC/voc2012/, 2012.

[9] 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 *AAAI*, 2018.

[10] S. Günel, H. Rhodin, and P. Fua. What face and body shapes can tell about height. In *ICCV Workshops*, 2019.

[11] I. Habibie, W. Xu, D. Mehta, G. Pons-Moll, and C. Theobalt. In the wild human pose estimation using explicit 2D features and intermediate 3D representations. In *CVPR*, 2019.

[12] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *ECCV*, 2016.

[13] C. Ionescu, F. Li, and C. Sminchisescu. Latent structured models for human pose estimation. In *ICCV*, 2011.

[14] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu. Human3.6M: Large scale datasets and predictive methods for 3D human sensing in natural environments. *PAMI*, 2014.

[15] U. Iqbal, P. Molchanov, T. Breuel, J. Gall, and J. Kautz. Hand pose estimation via latent 2.5D heatmap regression. *ECCV*, 2018.

[16] H. Jegou, M. Douze, and C. Schmid. Hamming embedding and weak geometric consistency for large scale image search. In *ECCV*, 2008.

[17] L. Ke, M.-C. Chang, H. Qi, and S. Lyu. Multi-scale structure-aware network for human pose estimation. In *ECCV*, 2018.

[18] S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-end training of deep visuomotor policies. *JMLR*, 17(1):1334–1373, 2016.

[19] D. Liu, Z. Zhao, X. Wang, Y. Hu, L. Zhang, and T. Huang. Improving 3D human pose estimation via 3D part affinity fields. In *WACV*, 2019.

[20] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *ICLR*, 2019.

[21] C. Luo, X. Chu, and A. Yuille. OriNet: A fully convolutional network for 3D human pose estimation. In *BMVC*, 2018.

[22] C. Luo, X. Chu, and A. Yuille. OriNet-demo. https://github.com/chenxuluo/OriNet-demo, 2018. [accessed 16-Nov-2018].

[23] D. C. Luvizon, D. Picard, and H. Tabia. 2D/3D pose estimation and action recognition using multitask deep learning. In *CVPR*, 2018.

[24] D. C. Luvizon, H. Tabia, and D. Picard. Human pose regression by combining indirect part detection and contextual information. *Computers & Graphics*, 85:15–22, 2019.

[25] J. Martinez, R. Hossain, J. Romero, and J. J. Little. A simple yet effective baseline for 3D human pose estimation. In *ICCV*, 2017.

[26] D. Mehta et al. Vnect: Real-time 3D human pose estimation with a single RGB camera. *ACM Trans. Graphics*, 36(4):44, 2017.

[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 *3DV*, 2017.

[28] D. Mehta, O. Sotnychenko, F. Mueller, W. Xu, S. Sridhar, G. Pons-Moll, and C. Theobalt. Single-shot multi-person 3D pose estimation from monocular RGB. In *3DV*, 2018.

[29] N. Neverova, R. A. Güler, and I. Kokkinos. Dense pose transfer. In *ECCV*, 2018.

[30] A. Newell, K. Yang, and J. Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, 2016.

[31] A. Nibali, Z. He, S. Morgan, and L. Prendergast. Numerical coordinate regression with convolutional neural networks. *arXiv:1801.07372*, 2018.

[32] A. Nibali, Z. He, S. Morgan, and L. Prendergast. 3D human pose estimation with 2D marginal heatmaps. In *WACV*, 2019.

[33] B. X. Nie, P. Wei, and S. Zhu. Monocular 3D human pose estimation by predicting depth on joints. In *ICCV*, 2017.

[34] G. Pavlakos, X. Zhou, and K. Daniilidis. Ordinal depth supervision for 3D human pose estimation. In *CVPR*, 2018.

[35] G. Pavlakos, X. Zhou, K. G. Derpanis, and K. Daniilidis. Coarse-to-fine volumetric prediction for single-image 3D human pose. In *CVPR*, 2017.

[36] G. Pavlakos, X. Zhou, K. G. Derpanis, and K. Daniilidis. Coarse-to-fine volumetric prediction for single-image 3D human pose: Supplementary material. In *CVPR*, 2017.

[37] J. Redmon and A. Farhadi. YOLOv3: An incremental improvement. *arXiv:1804.02767*, 2018.

[38] H. Rhodin, J. Spörri, I. Katircioglu, V. Constantin, F. Meyer, E. Müller, M. Salzmann, and P. Fua. Learning monocular 3D human pose estimation from multi-view images. In *CVPR*, 2018.

[39] G. Rogez, P. Weinzaepfel, and C. Schmid. LCR-Net: Localization-classification-regression for human pose. In *CVPR*, 2017.

[40] N. Sarafianos, B. Boteanu, B. Ionescu, and I. A. Kakadiaris. 3D human pose estimation: A review of the literature and analysis of covariates. *CVIU*, 152:1–20, 2016.

[41] I. Sárándi, T. Linder, K. O. Arras, and B. Leibe. How robust is 3D human pose estimation to occlusion? In *IROS Workshop - Robotic Co-workers 4.0*, 2018.

[42] I. Sárándi, T. Linder, K. O. Arras, and B. Leibe. Synthetic occlusion augmentation with volumetric heatmaps for the 2018 ECCV PoseTrack challenge on 3D human pose estimation. *arXiv:1809.04987*, 2018.

[43] V. Srivastav, T. Issenhuth, A. Kadkhodamohammadi, M. de Mathelin, A. Gangi, and N. Padoy. MVOR: A multi-view RGB-D operating room dataset for 2D and 3D human pose estimation. In *MICCAI LABELS Workshop*, 2019.

[44] X. Sun, C. Li, and S. Lin. An integral pose regression system for the ECCV2018 PoseTrack Challenge. *arXiv:1809.06079*, 2018.

[45] X. Sun, J. Shang, S. Liang, and Y. Wei. Compositional human pose regression. In *ICCV*, 2017.

[46] X. Sun, B. Xiao, S. Liang, and Y. Wei. Integral human pose regression. In *ECCV*, 2018.

[47] X. Sun, B. Xiao, S. Liang, and Y. Wei. Integral Human Pose Regression (code repository). https://github.com/JimmySuen/integral-human-pose, 2018. [Online; accessed 28-Apr-2019].

[48] D. Tome, C. Russell, and L. Agapito. Lifting from the deep: Convolutional 3D pose estimation from a single image. In *CVPR*, 2017.

[49] A. Toshev and C. Szegedy. DeepPose: Human pose estimation via deep neural networks. In *CVPR*, 2014.

[50] M. Véges and A. Lőrincz. Absolute human pose estimation with depth prediction network. In *IJCNN*, 2019.

[51] S. Vosoughi and M. A. Amer. Deep 3D human pose estimation under partial body presence. In *ICIP*, 2018.

[52] W. Yang, S. Li, W. Ouyang, H. Li, and X. Wang. Learning feature pyramids for human pose estimation. In *CVPR*, 2017.

[53] W. Yang, W. Ouyang, X. Wang, J. Ren, H. Li, and X. Wang. 3D human pose estimation in the wild by adversarial learning. In *CVPR*, 2018.

[54] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang. Random erasing data augmentation. In *AAAI*, 2020.

[55] X. Zhou, Q. Huang, X. Sun, X. Xue, and Y. Wei. Towards 3D human pose estimation in the wild: a weakly-supervised approach. In *ICCV*, 2017.

[56] C. Zimmermann, T. Welschehold, C. Dornhege, W. Burgard, and T. Brox. 3D human pose estimation in RGBD images for robotic task learning. In *ICRA*, 2018.