# Orientation Keypoints for 6D Human Pose Estimation

Martin Fisch, Ronald Clark

**Abstract**—Most realtime human pose estimation approaches are based on detecting joint positions. Using the detected joint positions, the yaw and pitch of the limbs can be computed. However, the roll along the limb, which is critical for application such as sports analysis and computer animation, cannot be computed as this axis of rotation remains unobserved. In this paper we therefore introduce orientation keypoints, a novel approach for estimating the full position and rotation of skeletal joints, using only single-frame RGB images. Inspired by how motion-capture systems use a set of point markers to estimate full bone rotations, our method uses virtual markers to generate sufficient information to accurately infer rotations with simple post processing. The rotation predictions improve upon the best reported mean error for joint angles by 48% and achieves 93% accuracy across 15 bone rotations. The method also improves the current state-of-the-art results for joint positions by 14% as measured by MPJPE on the principle dataset, and generalizes well to in-the-wild datasets.

**Index Terms**—Computer Vision, Pose Estimation, Pose Tracking, 6D Estimation.

✦

## 1 INTRODUCTION

Human motion capture (MoCap) has been a major enabling technology across both the arts and sciences. Motion capture has played a key role in kinematic analysis for sports and medicine, has created engaging user experiences with devices like the Kinect, and has been an essential part of the visual effects industry for years. In the past, MoCap required sophisticated purpose-built studios with multi-camera capture systems. However, recent advances in computer vision have led to new ways of doing MoCap that are far less restrictive than traditional methods. These approaches can capture 3D human poses from single RGB cameras and have spurred interest in next-generation applications such as personal digital sports coaches, and the possibility of capturing high-quality animation directly on consumer smartphones. However, despite their great promise, existing single-camera human pose estimation approaches have failed to achieve a level of fidelity that matches that of traditional MoCap systems.

Much of the current research in 3D pose estimation focuses on localizing joint keypoints with convolutional neural networks (CNNs). However, as shown in Figure 1 (b), most methods only detect keypoints at the joint locations. Using these detected points, the yaw, $\Psi$, and pitch, $\theta$, can be computed however, one degree of freedom is left unobserved, i.e., the roll, $\Phi$, around the axis. Therefore, most keypoint-based human pose estimation approaches can only observe five degrees-of-freedom for each joint, although there are six degrees of freedom, i.e., $(x, y, z, \Phi, \Psi, \theta)$.

In order to address this problem and estimate the full six the degrees of freedom, we propose a method that takes inspiration from traditional MoCap systems. MoCap systems use a large set of markers attached to the body, as shown in Figure 1 (a). Groups of these markers are used

- *Department of Computing, Imperial College London.*
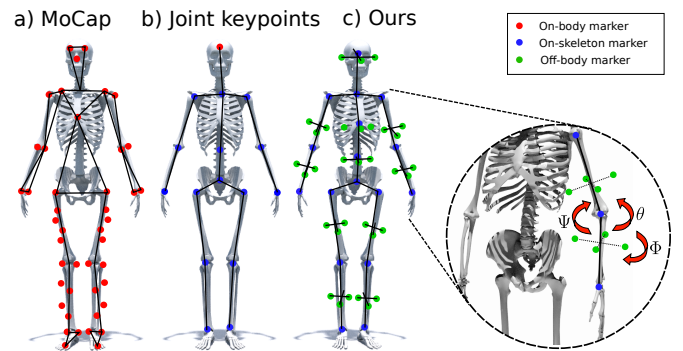
*Manuscript received July 10, 2020.*



Fig. 1. Our orientation keypoints approach uses a set of virtual point markers, similar to that used in motion capture systems, to accurately infer joint positions and orientations. a) Motion capture systems place sufficiently many (physical) markers to fully observe rigid joint rotations, including roll $\Phi$, yaw $\Psi$ and pitch $\theta$ of each bone. b) Existing joint keypoint-based human pose estimation methods place keypoints at joints, which only allows $\Psi$ and $\theta$ to be estimated, leaving roll unobserved. c) Our approach estimates virtual markers or "orientation keypoints" capturing the full joint rotation.

to compute the orientation of each bone. For example, the upper leg has four markers attached, which are used to solve the femur's position and orientation.

Inspired by this, we introduce a novel keypoint-based approach that solves for complete kinematic transforms of a human skeleton, with six degrees of freedom at each bone as illustrated in Figure 1 (c). Key to our approach is an additional set of *orientation keypoints* that provide sufficient information for inferring full joint orientation.

Specifically, our contributions in this paper are threefold:

1) We introduce orientation keypoints as a novel approach to solving for full 3-axis joint rotations.
2) We propose a neural network model that accurately localizes these points in 3D from monocular images,

achieving state-of-the-art accuracy for joint rotation and position estimation.

3) We demonstrate that our approach generalizes well to other datasets even without retraining, and with fine-tuning we set state-of-the-art benchmarks on these additional sets as well.

## 2 RELATED WORK

Deep learning related pose estimation initially focused on estimating 2D poses, with many of these techniques subsequently extended to 3D. For relevance, we focus on monocular single-image pose estimation works here.

**Pose Estimation in 2D**. Many approaches [58], [63], [67] predict discrete heatmaps estimating the probability of a joint occurring at each pixel instead of a continuous regression. Converting between the heatmaps and point coordinates can be done in several ways. The most common is to either take a "hard" argmax of the maps, or via integral regression (i.e., a "soft" argmax) [58]. We use soft argmax for this paper as it delivered slightly more accurate quantitative results. The architecture of these approaches vary, [67] use a ResNet backbone and three convolution transpose layers to upsample into 64x64 pixel heatmaps while [45] introduce the "stacked hourglass" architecture and [17] use Mask-RCNN and pixel-by-pixel masks to predict keypoints. One of the most successful approaches has been the cascaded pyramid network (CPN) which aims to address the problem of hard keypoints, integrate feature representations and use online hard keypoint mining loss (OHKM) [9]. Augmenting the keypoint heatmaps with part affinity fields has also shown to be very beneficial especially when predicting poses for a variable number of people [5]. In contrast to these works, our approach operates in 3D but can also be used to predict 2D keypoints simply by dropping the z-dimension in the predictions.

**Pose Estimation in 3D** Multiple 3D interpretations typically exist for a single 2D skeleton [32]. Therefore many approaches use a preconstructed model to map 2D detections to 3D [1], [7], [24], [62], [69]. The first approaches along these lines created a pose dictionary from 3D MoCap data to generate paired 2D projections from different angles and generate depth values with a lookup [1], while others have used a nearest neighbor search [7], [24]. Geometric information, such as bone length priors and projection consistency, can also be utilized for converting from 2D to 3D [4]. Other approaches take this a step further by trying to establish the correspondence between the 2D image and a 3D human model. This has been done, for example, by using distance matrix regression [44] or by directly predicting dense correspondences between pixels in the image to UV coordinates on a body mesh [55] or landmark locations on the body [31]. This differs from our research because we do not use body landmarks or a body model but instead chose detached points to maximize the angular perspective.

Another popular paradigm is to "lift" 2D detections to 3D using a learned network [41], [44], [55], [62]. In fact, [41] showed that lifting ground truth 2D locations to 3D can be solved with a low error rate with a relatively simple network. Other approaches perform a direct 3D prediction of keypoints from the images [16], [48], [51], [60]. Many of

these approaches are voxel-based which can be memory intensive and requires discretizing the space at a suitable resolution. To overcome this, [51] propose a fine discretization of the 3D space around the subject and train a network to predict per voxel likelihoods for each joint. Combining 2D information with the 3D predictions can also help improve accuracy, and therefore [48], [60] fuse direct image 3D features with 2D estimation while [16] embed 3D pose cues in a learned latent space. In this paper we consider two models, one which predicts 3D keypoints through regression and one which predicts through per dimension heatmaps.

**Weak Supervision and Generative Approaches.** As obtaining ground-truth labels for keypoints can be challenging, several works have focussed on using other signals for training. Geometric constraints can be used to train on in-the-wild datasets in a self-supervised manner [12], [22], [29], [72]. Other, weaker supervision signals can also be used such the ordinal depths of human joints, acquired from supplementary human annotations [50]. Adversarial training has also been quite popular as it enables using unlabelled data for training or training 3D predictions with only 2D annotations [66]. This is usually accomplished by generating 3D pose predictions for images with only 2D annotations and using a discriminator which distinguishes implausible poses [8], [25], [66], [68]. In our approach we do not use weak supervision, as some works have reported convergerce issues with GAN-type losses, however, this could easily be included in our framework in the future.

**Estimation from video.** While most works have focussed on the single-frame setting, utilizing the temporal regularity of video can help to improve pose accuracy. The temporal regularity can be integreted in various ways. Some approaches [74] use explicit temporal smoothness constraints, while others have used recurrent LSTM / GRU units [14], [27], [57], and dilated temporal convolutions [53]. While our approach only relies on single frames, the technique can easily be extended to most multi-frame settings.

**Joint rotation prediction.** Existing research which directly estimates joint angles, can capture the full six degrees of freedom when used in conjunction with kinematic constraints of a skeleton model. Here, 3D joint positions are typically computed by using for the forward kinematics. However, these approaches significantly underperform location-based methods, as convolutional neural networks have not proven adept at modeling the non-linearities complexities of angular representations.

Various parameterizations can be used for the joint angles, such as quaternions [54], Euler angles [43] or by regressing 3x3 rotation matrices [70]. The estimated joint angles are then usually mapped onto a skeleton using forward kinematics [43]. The angular constraints within the skeleton's kinematic chain can be formulated in a differentiable manner and embedded directly into the network itself [73]. There is also a body of related work which predict 6D position and rotation of objects by estimating virtual 3D bounding box vertices [56], [61]. This is conceptually the closest to our approach, but we calculate 15 rotations of a highly complex kinematic chain using a mix of joint and virtual markers detached from the shape of the limb.

**Direct mesh regression**

There are two types of methods in this line of work,

the first type of approach does not regress mesh vertices / correspondences directly but use the "Skinned Multi-Person Linear" (SMPL) model [38], or the newer SMPL-X [49], which represent the mesh as a set of low dimensional parameters. In particular, SMPL represents the mesh as 10 shape parameters and 72 (24*3) pose parameters. The pose is specified as relative rotation vectors which define the angle between successive joints. For pose estimation, SMPL-based methods either regress the shape and joint angles directly from an image [11], or from various types of intermediate representation, including 2D keypoints and silhouette [52] or semantic segmentation of the body parts [47]. Other works do not regress the parameters directly, but predict more easily observable quantities such as joint positions and use an optimization to find the joint angles and shape parameters [3]. In this sense, our method is complementary to SMPL in that the joint angles estimated using our OKPs can be used to drive the SMPL model. However, unlike general SMPL approaches, our method does not require mesh annotations for training.

The second type of method directly regresses correspondences between a dense set of points in the RGB image and uv-locations on the canonical human mesh. This dense set of points lies on the mesh and the network only predicts those which are not occluded. In order to recover body pose (i.e. joint position and angles), these methods typically also make use of SMPL [71]. Specifically, they first regress dense uv-coordinates, followed by an estimation of the SMPL parameters conditioned on this information. As these methods only predict correspondences for visible (non-occluded) points in the image, estimating the pose of the full skeleton is difficult, and requires a complex setup such as that in [71] to obtain good performance. Furthermore, these methods require training data where dense correspondences have been established between the image pixels and mesh vertices. In general, this data is very difficult to obtain.

**Deep rotation estimation.** As rotations play an important role in many tasks, there has been a growing theoretical interest in finding out what representations work best with deep networks. The fundamental issue with representing rotations is that lie in the special orthogonal group, $SO(3)$, which consists of all orthogonal $3 \times 3$ matrices. This constraint means that any parameterization using less than 5 dimension is guaranteed to be discontinuous, which creates problems when training deep networks [75]. To address this, [75] propose a continuous representation for rotations using a 6D over-parameterization and the Gram-Schmidt procedure to recover the rotation matrix. More recently, in concurrent work, [34] showed that using a 9D over-parameterization followed by SVD to recover the rotation outperforms the Gram-Schmidt procedure of [75] in terms of accuracy. In contrast, we introduce a representation for bone poses in $SE(3)$, including both the translation and rotation, using a 12D over-parameterization called orientation keypoints. We are also the first to apply this concept of rotation over-parameterization to human pose estimation.

## 3 APPROACH

In this section we describe our approach for estimating human pose. In Section 3.1 we present orientation keypoints
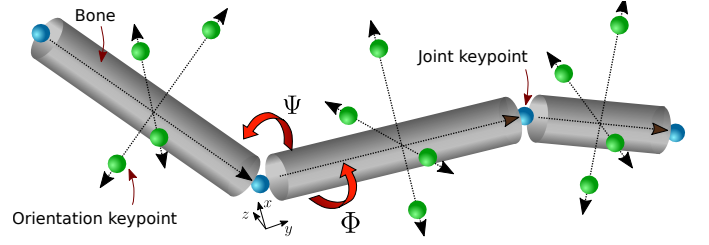


Fig. 2. Illustration of the arrangement of the joint keypoints (blue) and orientation keypoints (green) with respect to their parent bones. Without the orientaiton keypoints, the roll $\Phi$ is not observable.

(OKPS), the main component of our approach which is used as an intermediate representation from which bone rotations and translations can be computed. We then describe the networks we use to predict our OKPS and introduce our novel crosshairs architecture in Section 3.2. Finally, we describe how we post-process the OKPS detections to obtain bone translations and rotations.

### 3.1 Orientation keypoints

Conventional keypoint-based approaches for 3D human pose estimation focus on predicting a set of joint keypoints (JKPS),

$$P_{jkps} = \left( \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \ldots, \mathbf{x}_{N_j} \right) \qquad (1)$$

where each $p_{j_i} = (x, y, z)^\top$ describes the 3D location of a specific joint. The number of JKPS, $N_j$ is typically around 17, corresponding to joints on the human body such as the elbows, ankles, knees, wrists, hip, neck, etc. However, as alluded to in the introduction, this set of keypoints does not capture all the degrees of freedom of the bones of the human body. Therefore, we introduce an extra set of "orientation keypoints" (OKPS),

$$P_{okps} = \left( \mathbf{x}_{N_j}, \mathbf{x}_{N_j+1}, \mathbf{x}_{N_j+2}, \ldots, \mathbf{x}_{N_j+N_o} \right), \qquad (2)$$

and combine the two sets $P = P_{jkps} \bigcup P_{okps}$ for the purpose of pose estimation.

We define the OKPS as clusters of points rigidly attached to a particular joint keypoint to provide information about the two axes ignored by conventional JKPS. They differ from dense pose correspondences and landmarks in that they do not directly correspond to a specific body part or shape but are instead anchored in specific directions from the center of the bone (i.e., forward, or to the side) well offset from the body. These are inspired by "real" MoCap markers which are usually retro-reflective white balls that are attached rigidly to the actor. Orientation keypoints are therefore analogous to MoCap markers, but with the major advantage that no actual marker needs to be attached to the actor – they are simply virtual keypoints detected in relation to natural landmarks on the body. The difference is demonstrated in Figure 1. For example, we assign an orientation keypoint for the lower-left leg set midway between the knee and ankle and well offset from the shin by half the leg bone's length. In this paper, we use four OKPS for each of the 15 free rotations in the 17 joint skeleton. Each OKPS is rigidly attached to the corresponding parent joint

at a distance scaled by the bone length, one forward, back, left, and right, defined in the neutral pose. In general, given a bone $k$ with its first joint keypoint, $x_{j_1}$, the four OKPS for the bone are defined as:

$$\mathbf{x}_{o_1} = \mathbf{x}_{j_1} + T_k \times l_k(0.5, 0.5, 0)^\top$$
$$\mathbf{x}_{o_2} = \mathbf{x}_{j_1} + T_k \times l_k(0.5, -0.5, 0)^\top$$
$$\mathbf{x}_{o_3} = \mathbf{x}_{j_1} + T_k \times l_k(0.5, 0, 0.5)^\top$$
$$\mathbf{x}_{o_4} = \mathbf{x}_{j_1} + T_k \times l_k(0.5, 0, -0.5)^\top \quad (3)$$

where $l_k$ is the length of bone $k$, and $T_k$ is the transformation (i.e. pose) that converts coordinates in the bone's keyframe to world space.

## 3.2 Network design

Our framework's main component is a convolutional neural network detector that localizes both joint and orientation keypoints. Since the latter is offset and virtual, this requires learning depth and perspective even for 2D predictions, and so we choose to predict the full 3D keypoint locations directly in our detector. This also enables direct calculation of the full kinematic rotations from the model without further lifting. We also explore using a two-stage process similar to [41] and [53], where the second stage is a lifter model which transforms 2D detections into 3D predictions (without a detector depth branch), or a refiner model which further hones initial 3D predictions.

### 3.2.1 Detector models

For the detector, we experiment with two models; a simple Resnet based 3D regression baseline and a more sophisticated novel architecture, which we call Crosshairs, capable of providing accurate 3D estimates while limiting the memory and calculation overhead. We describe them both below.

**Simple regression baseline** This simple baseline detector uses a Resnet50 as the backbone and adds a head connected to the final convolution layer (removing the final pooling and fully connected layer in the base Resnet). Our simple head is composed of four layers using grouped convolutions. The grouping focuses the model at each keypoint and considerably limits the calculations beyond the backbone. The overall architecture consists of,

- *A convolutional layer* with a 1x1 kernel, batchnorm and ReLU. We use 12 x Number of Keypoints (77) = 924 channels. This prepares features for the grouped-by-keypoint convolutions
- *A second grouped convolutional layer* with a 5x5 kernel (2x2 padding), batchnorm and ReLU. We use the same number of channels but with 77 groups. This means for each keypoint there are 12 convolutional filters each only using 12 channels from the previous layer.
- *A grouped convolutional layer* with a 1x1 kernel, batchnorm and ReLU. The 924 channels and 77 groups are the same.
- A final *fully connected grouped convolutional layer*. As the kernel is the full width and height of the layer (12x9) this acts as a fully connected layer, except it is again grouped by keypoint: each keypoint's

xyz prediction is made only from the 12 associated channels. The number of channels is 3 (for XYZ) x Number of Keypoints, each outputting a single value.

This detector is simple, very fast, and quick to train. In terms of accuracy, it is a strong baseline model and we show good results for Orientation Keypoints, even using this baseline detector. As we show in Table 6, it is sufficient to achieve state-of-the-art results in 3D on Human 3.6m when used with orientation keypoints.

**The Crosshairs detector** Conceptually we follow [67] and use a Resnet backbone with convolution transpose layers to recover a higher resolution. Our key innovation is that at the output, each strand uses 1D heatmaps per dimension in place of square and volumetric heatmaps for 2D and 3D estimation. This keeps the computation cost of the head linear with the resolution rather than square or cubic and is a substantial saving, particularly considering the larger number of keypoints we employ.

For example, one *x* strand takes the C5 2048 channel 8x8 tensor and first samples into a 256 channel 8x8 tensor using 1x1 kernel convolutions. The y-dimension is then flattened into a tensor with only the x-dimension using an 8x1 convolution. Flattening is followed by a bottleneck block with padding and a kernel size of 1x9 to immediately provide a global view - in a single dimension and at a low resolution, which is computationally efficient. We then use transpose convolution layers to upsample, in one-dimension, back to the original resolution width. The upsampling approach is similar to [67] but much cheaper, as we are operating in a single dimension. Each convolution and convolution transpose layer is followed by a batchnorm layer [20] and ReLU [15] activation function. We use 256 channels throughout. A final 1x1 convolution layer collapses the channels into a 1D heatmap for each keypoint along the single dimension. Each heatmap represents the network's estimate of the keypoint position along the single axis. We then apply the same technique for the Y dimension by instead flattening the other dimension. As depth is not a native dimension, we use the same principle but a modified flattening procedure. Again, we sample the backbone C5 layer with 1x1 convolutions but with more channels to reshape into a depth dimension, i.e. into a 256 channel 9x9x12 tensor (we arbitrarily decide the tensor depth resolution is equal to the narrowest XY dimension). We then use a convolution layer and a 1x9x12 kernel, which only slides in the depth dimension, to collapse x and y into a 9x1x1 block, which is again a flattened to 1D as before.

This leads to 3 vectors for each point representing the same data as a volume would but much more efficiently,

$$\mathbf{v}^x = (v_k^x)_{k=1}^{N_x}, \mathbf{v}^y = (v_k^y)_{k=1}^{N_y}, \mathbf{v}^z = (v_k^z)_{k=1}^{N_z} \quad (4)$$

The keypoints are then recovered by using a coordinate-weighted softmax [40]:

$$\mathbf{x}_k = \left( \frac{\sum_i e^{v_i^x} w_i^x}{\sum_i e^{v_i^x}}, \frac{\sum_i e^{v_i^y} w_i^y}{\sum_i e^{v_i^y}}, \frac{\sum_i e^{v_i^z} w_i^z}{\sum_i e^{v_i^z}} \right) \quad (5)$$

where $w_i^x = \frac{i - 0.5 N_x}{0.5 N_x}$ so that $w_i \in [-1, 1]$.

The introduction of orientation keypoints adds points that may lie outside the subject's silhouette and a tighter
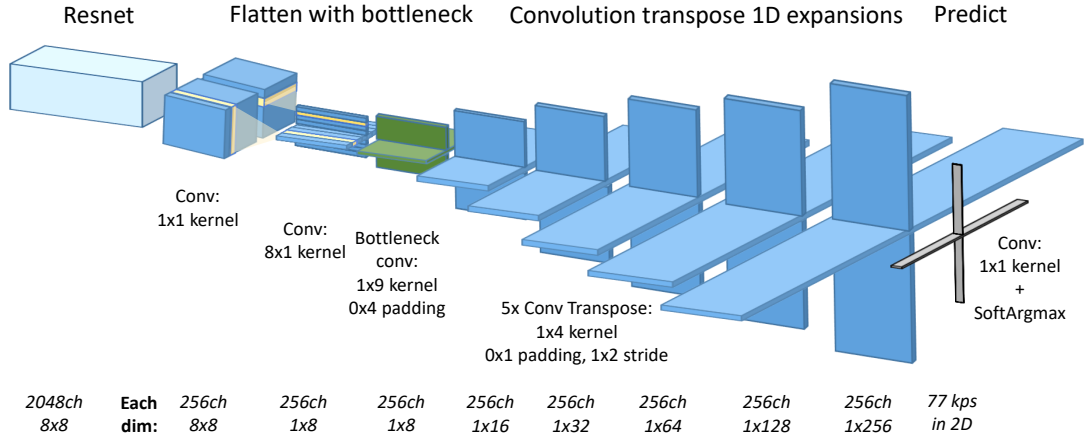
Fig. 3. Overview of the crosshairs detector head. We visualize only a 2D detector with single strands for X and Y attached to the C5 Resnet layer. Each strand flattens the backbone layer along one dimension, filters through a bottleneck with a full width kernel, and then expands the single dimension back to the original resolution with convolution transpose layers. Batch normalization and ReLU activation layers follow convolutions. A 1x1 convolution layer and soft argmax make predictions.

fitting bounding box. Rather than use a larger bounding box and lose effective resolution, we instead map the soft argmax layer output to a 25% wider pixel range than the underlying image - thus each heatmap covers a wider area than the image itself.

For the benefits of intermediate supervision and higher resolution access, we propose using multiple crosshairs, one attached to each layer group of Resnet. We aggregate crosshair strands with a 1x1 kernel convolution layer combining the concatenated high-resolution 1D feature maps, followed by batchnorm, ReLU, and a bottleneck block with a 1x5 kernel. This produces the final predictions.

### 3.2.2 Lifter/refiner regression model

The lifter takes as input the set of normalized keypoints predicted by the detector (either a 77x3 or 77x2 matrix depending on whether a 2D or 3D detector is used) and outputs the keypoints in metric space. For the architecture of the lifter/refiner block we follow [41] and use a similar architecture. This entails an inner block with a linear layer, followed by batch normalization [20], dropout [19] and rectified linear units [15]. The outer blocks contain two inner blocks and a residual connection. A first linear layer converts from the number of keypoint inputs, flattened to a single dimension, into the network's linear width. A final layer converts from the width to the number of predictions. We use two outer blocks in our lifter/refiner, but widen the network compared to [41], increasing the size of each linear layer by 50% from 1024 to 1536, which approximately doubles the total parameters. This helps accommodate the 5-6x as many keypoint inputs and outputs needed for orientation keypoints.

### 3.3 Post processing

For inference, we take the average of the predictions and the horizontally flipped predictions and use these to compute the rotations and positions of each joint.

**Positions.** As the detector predicts keypoints in normalized pixel or voxel units, we consider two approaches to predicting real-world 3d positions from orientation keypoints

: (a) we map the rotations onto a full kinematic skeleton, based on the average bone lengths from the training set, and (b) train a second stage refiner network described in Section 3.2.2. This can lift purely from 2D or refine 3D voxels. The first approach more elegantly unites rotations with positions for 6D, easily allows remapping onto different sized individuals in new environments and can deal well with novel poses. Differences in skeleton size however, contribute to the error. The second approach effectively bakes the skeleton size and camera perspectives into an additional neural network during training and is more reliant on the set of training poses, as highlighted in [41], but is more accurate in the Human 3.6m setting.

**Rotations.** We calculate each joint's rotation from the two JKPS and four OKPS associated with each bone with reference to the neutral T-pose positions of these points, normalized in bone length units (i.e., independent of the actual skeleton). From a set of 2D joint and orientation keypoint estimates, the rotations could be determined with a Perspective-n-Point algorithm, such as [33]. As our network also learns depth, we use the predictions to reproject XY detections into voxel 3D-space and use these estimates to solve for the transform, which minimizes the least square error from neutral pose. We found the 3D approach more accurate than PnP and faster. Therefore, we use the method attributed to [64], which works as follows. Here we define $(\mathbf{x}_{k_i})_{i=1}^6$ as the set of 2 JKPS and 4 OKPS associated with bone $k$. We first compute the centroids of the predictions and the T-pose points ($\mathbf{y}$),

$$\bar{\mathbf{x}}_k = \frac{1}{6}\sum_{i=1}^{6}\mathbf{x}_{k_i}, \quad \bar{\mathbf{y}}_k = \frac{1}{6}\sum_{i=1}^{6}\mathbf{y}_{k_i} \tag{6}$$

as well as the variance for re-scaling,

$$\sigma^2 = \frac{1}{6}\sum_{i=1}^{6}\|\mathbf{x}_{k_i} - \bar{\mathbf{x}}_k\|^2 \tag{7}$$

and the covariance,

$$M_k = \frac{1}{6} \sum_{i=1}^{6} \left(\mathbf{x}_{k_i} - \bar{\mathbf{x}}_k\right)^T \left(\mathbf{y}_{k_i} - \bar{\mathbf{y}}_k\right). \qquad (8)$$

We then take the Singular Value Decomposition (SVD) of the covariance,

$$M_k = \mathrm{USV}^\mathrm{T} \qquad (9)$$

and compute the rotation of bone $k$ as,

$$R_k = UV^T \qquad (10)$$

As JKPS detections are more accurate than OKPS detections, we double the JKPS correspondences' weight in the solution. Using the predicted rotations, we can then trivially infer joint positions to scale from a given skeleton (i.e., bone lengths). Specifically, we use the average bone lengths from the five-subject H3.6m training set.

### 3.4 Summary of configurations

In the previous sections we defined various detectors options (i.e. 2D/3D and crosshairs/regression) and post-processing methods (i.e. mapping to skeleton or using refiner module). In this paper, we consider four main configurations of these components. Figure 4 shows the overall pipeline for each configuration.
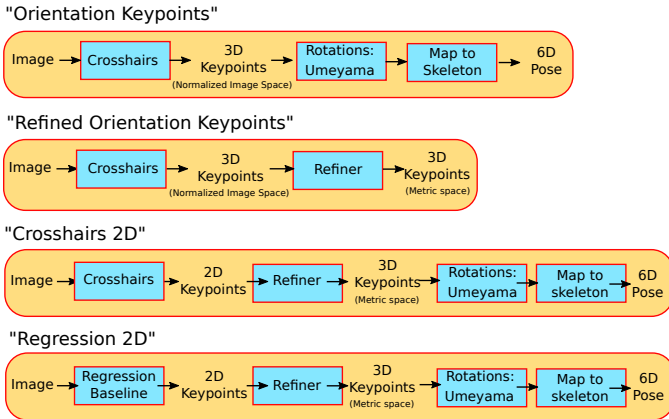


Fig. 4. Summary of the pipeline configurations used in this paper. Configuration names correspond to those used in the experiments.

All configurations output the 6D pose for each bone of the skeleton. The "Refined orientation keypoints" configuration is special in that it is only used for evaluating the keypoint accuracy when using a refiner for converting from normalized image-space 3D keypoints to metric 3D keypoints.

## 4 DATA AND TRAINING

For 3D pose estimation research, the Human3.6m dataset is the most commonly used and includes 3.6 million accurate 3D Human poses acquired by recording the performance of five female and six male subjects, under four different viewpoints, introduced in [6] and [21]. Other 3D pose estimation approaches, focused on joint keypoints only, usually train on more diverse datasets and only finetune keypoint locations with Human3.6m, which helps preserve

generalization. We found that just training on Human3.6m data was problematic as the model would quickly saturate and struggle to generalize well in validation. While the set may have many images, there are only five training subjects, and the high parameter models may memorize the specific subjects instead of learning more general rules. We also use the MPII Human Pose dataset [2] during training to help with generalization.

### 4.1 Losses

We train our networks to minimize the prediction error over a dataset of poses, where the error is the Mean Per Joint Position Error (MPJPE) which is defined as:

$$\mathcal{L}_{mpjpe} = \frac{1}{N} \sum_{j=1}^{N} ||\mathbf{x}_j - \mathbf{x}_{j,gt}||_2 \qquad (11)$$

where $N = N_o + N_j$. For the crosshairs, we also experimented with a regularizer. This regularizer forces the orientation keypoints at the bone ends towards their centroid.

$$\mathcal{L}_{cnt} = \sum_k \sum_i ||(\mathbf{x}_{k_i} - \bar{\mathbf{x}}_k) - (\mathbf{x}_{k_i,gt} - \bar{\mathbf{x}}_{k,gt})||_2 \qquad (12)$$

We use this to encourage the model to generate a better prediction structure in line with the orientation algorithm used in post-processing. We expect this to behave similarly to cross-joint loss functions used in the literature, which compare each joint's relative pose to every other joint. However, we use a cross-comparison limited to immediate neighbors. In our experiments, the accuracy improvement was *de minimis* compared to an identical model trained without the additional loss function. We did not try this loss function with the simpler regression detector to preserve simplicity.

### 4.2 Metrics

To evaluate the performance of the model, we use MPJPE defined above, along with two angular metrics. Specifically, we define mean average accuracy for rotations as:

$$MAA = \frac{1}{N} \sum_{j=1}^{N} 1 - \frac{\theta_{sep}(R_{j,gt}, R_j)}{\pi} \qquad (13)$$

Randomly drawn uniformly distributed rotation predictions average approximately 30% accuracy (with some variation depending on how uniformity is defined for a rotation).

Angular separation maps values to $[0, \pi]$ radians and formally is defined as:

$$\theta_{sep}(R_m, R_n) = ||log(R_m R_n^T)|| \qquad (14)$$

We define MPJAS as the maximum angular separation of points transformed by two rotations:

$$MPJAS = \frac{1}{N} \sum_{j=1}^{N} \theta_{sep}(R_{j,gt}, R_j) \qquad (15)$$

Randomly drawn rotation predictions average approximately 2.2 radians error. MPJAS-15 represents the mean joint comparisons between ground truth and predictions for the 15 free rotations in the typical 17 joint skeleton.

## 4.3 Human 3.6m preparation

Following previous researchers, we extract frames from each video at a downsampled rate of 10HZ (i.e., 1 per 5 frames) and excluded a corrupted sequence for subject 11. For validation and testing, we tried using all frames, and 1/65 frames (i.e., 1 per 13 at 10hz) as different benchmark papers use different frequencies. For each image we use a bounding box from ground truth joint data. We also get comparable but slightly less accurate results when using a fixed pixel bounding box similar to [41], mainly due to effective resolution loss as the figures have smaller sizes. Orientation keypoints are calculated, based on the provided angle data and projected into 2D screen coordinates. For 3x4 aspect ratio resolutions, we preserve scale and crop the width. For depth detections, we convert depth information into depth equivalent relative to the root.

## 4.4 MPII Human Pose

MPI-3DHP [42] uses multiple cameras and markerless technology to estimate ground truth data in more varied scenes than Human 3.6m; this capture method is somewhat less accurate than markers as used in Human 3.6m and does not provide joint rotations, and therefore we can only test position accuracy. The testing is facilitated by a very similar skeleton to Human 3.6m. For the MPII Human Pose dataset [2] we bulk adjust annotations for closer consistency with Human3.6m, namely the feet and head.

During the keypoint detector training, we augment the data with random horizontal flipping, rotating by up to +/-30%, random cropping, positional, and color jitter. For the lifter/refiner, we use predictions from the first stage and subtract the root location from all 3D keypoints, as is conventional in the literature. We augment with horizontal flipping and randomly increasing the detector error by 0-100% vis-à-vis the ground truth.

## 4.5 MPI-3DPW preparation

MPI-3DPW [65] is a recent dataset that tries to capture full 6D skeletal poses in the wild by using IMUs and a single camera markerless algorithm. The 3DPW dataset uses a different skeleton than the 17 joint skeleton commonly used in the literature. We therefore needed to reconcile the skeleton with Human 3.6m. While the joint keypoints can generally be matched up, they are placed differently, most noticeably at the hips and the ankles. For the root hip position, we take the midpoint of LeftUpLeg and RightUpLeg to match Human3.6m. The other differences still meaningfully impact the MPJPE accuracy: we show substantial improvements in accuracy using our original model from just changing the bone lengths (i.e., narrow hip bone) and even more by fine-tuning to learn the new keypoint locations.

We also need to calculate rotation accuracy taking into account the additional joints in the kinematic chain (i.e., 3DPW uses a multi-segmented spine) and different orientation conventions (Human 3.6m follows Vicon convention, 3DPW seems to be zero rotation in a T-pose). To put 3DPW in a common basis, we recalculated the ground truth joint rotations by realigning the rotation matrices (calculated

### TABLE 1
Rotation results on Human3.6M. We average the 15 free rotations in the typical 17 joint skeleton. We convert [70] as they report on a different basis. Low MPJAS (in radians) and high MAA (accuracy) is better.

|  | MPJAS | MAA |
|---|---|---|
| Yoshiyasu et. al. [70] ACCV '18 | 0.424 | 86.5% |
| - pelvis (root) only | 0.226 | 92.8% |
| Orientation Keypoints (ours) |  |  |
| 2D detections + PnP | 0.265 | 91.6% |
| 3D detections + SVD | 0.213 | 93.2% |
| - pelvis (root) only | 0.145 | 95.4% |

from the provided rotation vectors): we align the Y vector to point from parent to child joint (reversed for lower body), keep the parent annotation Z vector as forward and then orthogonalize and normalize the X and Z vectors.

## 4.6 Training regimen

We begin with transfer learning, using an off-the-shelf Resnet-50 [18] backbone used for CPN [10] as a human keypoint detector on COCO [37] and then discard the head. For the simple regression detector, we warmup the new head with 1k iterations, and then train the head and layers C4 and C5 of Resnet for 20k iterations at 0.001/0.0001 learning rate, and then another 80k iterations at 0.00025/0.0001. Each iteration is a 64 sample batch split 75/25 Human3.6m and MPII. For crosshairs we initially train for 40k iterations at 0.001/0.00005 learning rates for the head and backbone, respectively, using L2 loss. We drop the head learning rate to 0.00025 and train for another 40k iterations. We then shift the data mix to 75/25 Human36m/MPII and train for another 80k iterations using L1 loss instead. We use the Adam optimizer [26] and batch normalization [20].

For the second stage refiner, we train for 80 epochs on L2 loss using 0.25 dropout, 0.1 momentum, and the Adam optimizer, starting at a learning rate of 0.001 and 0.98 learning rate gamma.

## 5 RESULTS

### 5.1 Predicted skeletal rotations

The literature is mostly devoid of published metrics on rotational prediction accuracy, as even papers purporting to solve for joint rotations instead choose only to show mean positional error. One exception is [70], which includes results for Human3.6m. Their best-stated result equates to 0.424 radians MPJAS-15[1]. Our best result of 0.213 radians MPJAS-15 is a 48% improvement on their results and is, as far as we are aware, the state-of-the-art in predicting skeletal rotations by a considerable margin.

**Table 1** shows summary prediction results for joint rotations on Human3.6m. We show results with predicted keypoint detection using different strategies to generate rotations based on MPJAS-15, the rotational error across 15 rotations. First, we use 2D detections of the joint and orientation keypoints to solve with PnP, which achieves 0.265 radians error MPJAS-15 across all actions. This equates

---

1. Based on converting their reported metric. As their code is not publicly available, we rely on their brief description in the paper

TABLE 2
Mean Per Joint Position Error (MPJPE) in mm between the ground-truth 3D joints on Human 3.6M for single frame RGB images without depth information. Best results are highlighted in bold. The "Orientation keypoints" results correspond to positions predicted by fitting rotations from the detector to the average training skeleton. "Refined orientation kps" are taken from after the second stage refiner network. Note that these reference methods use a variety of different training approaches, some include additional data and additional weak supervision.

| Protocol #1 | Dir | Disc | Eat | Greet | Phone | Photo | Pose | Punch | Sit | SitD | Smoke | Wait | Walk | WalkD | WalkT | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Zhou et al. [73] 6 | 91.8 | 102.4 | 97.0 | 98.8 | 113.3 | 125.2 | 90.0 | 93.8 | 132.2 | 159.0 | 106.9 | 94.4 | 79.0 | 126.0 | 99.0 | 107.3 |
| Moreno-Noguer [44] | 66.1 | 61.7 | 84.5 | 73.7 | 65.2 | 67.2 | 60.9 | 67.3 | 103.5 | 74.6 | 92.6 | 69.6 | 71.5 | 78.0 | 73.2 | 74.0 |
| Pavlakos et al. [51] | 67.4 | 71.9 | 66.7 | 69.1 | 72.0 | 77.0 | 65.0 | 68.3 | 83.7 | 96.5 | 71.7 | 65.8 | 59.1 | 74.9 | 63.2 | 71.9 |
| Yoshiyasu et al. [70] | 63.3 | 71.6 | 61.4 | 70.4 | 69.9 | 83.2 | 63.1 | 68.8 | 76.8 | 98.9 | 68.2 | 67.5 | 57.7 | 73.7 | 57.1 | 70.0 |
| XNect [43] | 50.2 | 61.9 | 58.3 | 58.2 | 68.8 | 54.1 | 61.5 | 76.8 | 91.7 | 63.4 | 74.6 | 58.5 | 48.3 | 65.3 | 53.2 | 63.0 |
| Martinez et al. [41] | 51.8 | 56.2 | 58.1 | 59.0 | 69.5 | 78.4 | 55.2 | 58.1 | 74.0 | 94.6 | 62.3 | 59.1 | 49.5 | 65.1 | 52.4 | 62.9 |
| Yang et al. [68] | 51.5 | 58.9 | 50.4 | 57.0 | 62.1 | 65.4 | 49.8 | 52.7 | 69.2 | 85.2 | 57.4 | 58.4 | 60.1 | 43.6 | 47.7 | 58.6 |
| Chen et al. [8] | 45.9 | 53.5 | 50.1 | 53.2 | 61.5 | 72.8 | 50.7 | 49.4 | 68.4 | 82.1 | 58.6 | 53.9 | 41.1 | 57.6 | 46.0 | 56.9 |
| Pavlakos et al. [50] | 48.5 | 54.4 | 54.4 | 52.0 | 59.4 | 65.3 | 49.9 | 52.9 | 65.8 | 71.1 | 56.6 | 52.9 | 44.7 | 60.9 | 47.8 | 56.2 |
| Luvizon et al. [40] | 49.2 | 51.6 | 47.6 | 50.5 | 51.8 | 60.3 | 48.5 | 51.7 | 61.5 | 70.9 | 53.7 | 48.9 | 44.4 | 57.9 | 48.9 | 53.2 |
| Pavllo et al. [53] | 47.1 | 50.6 | 49.0 | 51.8 | 53.6 | 61.4 | 49.4 | 47.4 | 59.3 | 67.4 | 52.4 | 49.5 | 39.5 | 55.3 | 42.7 | 51.8 |
| Refined jkps only (ours) | 45.1 | 50.4 | 47.4 | 49.1 | 54.7 | 61.9 | 46.3 | 45.4 | 61.5 | 71.4 | 51.2 | 48.0 | 38.6 | 52.3 | 42.5 | 51.1 |
| Orientation keypoints (ours) | 44.4 | 48.9 | 42.6 | 45.5 | 49.8 | 50.9 | 43.0 | 44.4 | 56.6 | 62.3 | 48.3 | 44.1 | 38.8 | 49.5 | 42.1 | 47.4 |
| Refined orientation kps (ours) | **40.7** | **45.5** | **39.5** | **42.3** | **48.1** | **49.2** | **40.3** | **39.6** | 56.7 | 61.3 | **45.8** | **41.2** | 35.3 | 46.8 | **36.8** | **44.6** |

| Protocol #2 (Procrustes) | Dir | Disc | Eat | Greet | Phone | Photo | Pose | Punch | Sit | SitD | Smoke | Wait | Walk | WalkD | WalkT | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Martinez et al. [41] | 39.5 | 43.2 | 46.4 | 47.0 | 51.0 | 56.0 | 41.4 | 40.6 | 56.5 | 69.4 | 49.2 | 45.0 | 38.0 | 49.5 | 43.1 | 47.7 |
| Chen et al. [8] | 36.5 | 41.0 | 40.9 | 43.9 | 45.6 | 53.8 | 38.5 | 37.3 | 53.0 | 65.2 | 44.6 | 40.9 | 32.0 | 44.3 | 38.4 | 44.1 |
| Pavlakos et al. [50] | 34.7 | 39.8 | 41.8 | 38.6 | 42.5 | 47.5 | 38.0 | 36.6 | 50.7 | 56.8 | 42.6 | 39.6 | 32.1 | 43.9 | 36.5 | 41.8 |
| Pavllo et al. [53] | 36.0 | 38.7 | 38.0 | 41.7 | 40.1 | 45.9 | 37.1 | 35.4 | 46.8 | 53.4 | 36.9 | 41.4 | 30.3 | 43.1 | 34.8 | 40.0 |
| Yang et al. [68] | **26.9** | **30.9** | 36.3 | 39.9 | 43.9 | 47.4 | **28.8** | **29.4** | **36.9** | 58.4 | 41.5 | **30.5** | | **29.5** | 32.2 | 37.7 |
| DaNet [71] | 35.7 | 40.4 | 39.0 | 40.3 | 40.5 | 47.4 | 35.1 | 34.9 | 45.2 | 51.7 | 39.6 | 37.8 | 43.4 | 34.4 | 39.8 | 40.5 |
| Refined jkps only (ours) | 34.2 | 38.0 | 38.3 | 37.8 | 39.9 | 45.3 | 35.5 | 33.9 | 47.9 | 54.2 | 40.7 | 35.6 | 28.9 | 39.4 | 33.5 | 38.9 |
| Orientation keypoints (ours) | 33.3 | 36.3 | 34.2 | 35.2 | 36.3 | 38.6 | 32.5 | 32.8 | 42.9 | 50.3 | 36.5 | 33.4 | 29.8 | 37.2 | 32.1 | 36.1 |
| Refined orientation kps (ours) | 31.7 | 34.5 | **32.7** | **33.9** | **35.4** | **38.1** | 31.5 | 30.9 | 42.9 | **49.2** | 35.8 | 32.3 | **27.6** | 35.8 | **30.5** | **34.9** |

TABLE 3
Results on the 3DHP test set. PCK is percentage correct (within 150mm) of 14 joints, a metric commonly used on this dataset, and MPJPE is again the Mean Per Joint Position Error (mm). PMPJE and PPCK correspond to Protocol 2 where the predictions are further aligned with the ground-truth via a rigid transform using Procrustes before computing MPJPE and PCK, respectively. The datasets on which the models have been trained are shown in brackets. The best results are bolded. Our approach outperforms all the existing methods, even without finetuning on 3DHP itself.

| | MPJPE | PMPJPE | PCK | PPCK |
|---|---|---|---|---|
| | *lower is better* | | *higher is better* | |
| *No 3DHP training* | | | | |
| Yang et al. [68] (H3.6m,MPII) | | | 69.0 | |
| Habibe et al. [16] (H3.6m) | 127.0 | 92.0 | 69.9 | 82.9 |
| OriNet [39] (H3.6m) | - | - | 71.3 | - |
| RepNet [35] (H3.6m, LSP) | 97.8 | - | 82.5 | - |
| Orientation kps (H3.6m, MPII) | 97.0 | **67.7** | 81.1 | 93.3 |
| Refined orientation kps (H3.6m, MPII) | **94.0** | 70.7 | 81.7 | **92.2** |
| *Trained with 3DHP* | | | | |
| SPIN [30] (H3.6m,3DHP,LSP,MPII,COCO) | 105.2 | 67.5 | 76.4 | 92.5 |
| Habibe et al. [16] (H3.6m,3DHP) | 90.7 | 65.4 | 81.5 | 91.3 |
| OriNet [39] (MPII, H3.6m, 3DHP) | 89.4 | - | 81.8 | - |
| XNect [43] (Coco, 3DHP) | 92.4 | - | 82.8 | - |
| MargiPose [46] (MPII, H3.6m,3DHP) | 91.3 | - | 85.4 | - |
| MEVA [46] (MPII, H3.6m,3DHP) | 96.4 | 65.4 | - | - |
| HybrIK [35] (H3.6m, 3DHP, COCO) | 91.0 | - | 86.2 | - |
| Orientation kps fine-tuned (H3.6m,MPII,3DHP) | **86.1** | **60.6** | **85.8** | **94.3** |

to 91.6% accuracy and is already a strong result for the first approach. However, using the full 3D predictions from our detector and SVD we achieve 93.2% accuracy (0.213 rad).

## 5.2 Predicted joint positions

While localizing joint positions is only the secondary goal of our method, most research focuses on this metric, and we also show meaningful improvements to the state-of-the-art. We follow most of the literature in training (S1, S5, S6, S7, S8) / test (S9, S11) split and definitions for Protocol 1 as raw prediction relative to root and Protocol 2 as allowing rigid alignment ('Procrustes') of the overall skeleton.

In Section 3.3, we proposed two approaches for estimating the joint positions from the predicted pixel-space values. In Table 2 we can see that, using the first approach, mapping to skeleton, improves the state-of-the-art MPJPE by 4mm under Protocol 1. The second approach reduces the mean error by another 3mm, for a total of 7mm. This is a significant 13.9% improvement on the previous state-of-the-art for single frame estimation and better than any reported MPJPE for video analysis as far as we are aware. Our approach could also be extended to video for potential further improvements. We also establish a new state-of-the-art under protocol 2. Again, a simple mapping of detector

TABLE 4
Results on the 3DPW test set. MPJPE is the Mean Per Joint Position Error (in mm), PMPJPE corresponds to MPJPE with the predictions aligned to the ground-truth using a rigid alignment and MPJAS is the Mean Per Joint Angular Separation which measures the angular separation of points transformed by two rotations (see Equation 15). Our method outperforms the current state of the art, Mesh GraphFormer [36] on MPJPE which is the best metric for realworld performance (i.e. without procrustes alignment).

| Method | MPJPE | PMPJPE | MPJAS |
|---|---|---|---|
| SPIN [30] (H3.6m,M3DHP,LSP,MPII,COCO) | 96.9 | 59.2 | - |
| Mesh Graphormer [36] (H3.6m, M3DHP,UP-3D, COCO, MPII, 3DPW) | 74.7 | **45.6** | - |
| PARE [28] (COCO, MPII, LSPET, MPI-INF-3DHP, H3.6m) | 79.1 | 46.4 | - |
| ROMP [59] (COCO, MPII, LSPET, AICH, MPI-INF-3DHP, H3.6m) | 80.1 | 56.8 | - |
| HybrIK [35] (COCO, MPI-INF-3DHP, H3.6m) | 80.0 | 48.8 | - |
| MEVA [46] (MPII, H3.6m,3DHP) | 86.9 | 54.7 | - |
| MeshTransformer [36] (H3.6m, UP-3D, MuCo-3DHP, COCO, MPII, FreiHAND, 3DPW) | 77.1 | 47.9 | - |
| Orientation kps (H3.6m, MPII) | 115.4 | 76.7 | 0.408 |
| Refined Orientation kps (H3.6m, MPII) | 112.4 | 67.6 | - |
| Orientation kps (H3.6m, MPII, using 3DPW avg bone lengths) | 90.3 | 66.9 | 0.408 |
| SPIN using Orientation kps in the optimization loop (H3.6m, MPII, 3DPW fine-tuned) | 81.8 | 54.3 | - |
| Orientation kps fine-tuned (H3.6m, MPII, 3DPW) | **70.7** | 50.4 | **0.302** |

TABLE 5
Ablation study of the lifter showing the effect of different keypoint detectors, and keypoint types. We report MPJPE, PMPJPE and MPJAS when using groundtruth 2D detections, regressed 2D detections, and predictions from our crosshairs architecture (for both 2D and 3D). We show results using only joint keypoints (JKPS) and when augmenting these with orientation keypoints (J+OKPS). The lifter is retrained for each scenario. Our 3D Crosshairs detector performs the best in all cases apart from using ground-truth 2D detections (as expected).

| | Groundtruth 2D | | Regression 2D (288x384) | | | Crosshairs 2D (288x384) | | | Crosshairs 3D |
|---|---|---|---|---|---|---|---|---|---|
| Input detections | JKPS | J+OKPS | JKPS | JKPS | J+OKPS | JKPS | JKPS | J+OKPS | J+OKPS |
| Output predictions | JKPS | J+OKPS | JKPS | J+OKPS | J+OKPS | JKPS | J+OKPS | J+OKPS | J+OKPS |
| Detector 2d err (% res) | 0% | 0% | 1.75% | 1.75% | 1.75% | 1.60% | 1.60% | 1.60% | 1.60% |
| MPJPE (mm) | 48.3 | 33.3 | 57.5 | 57.8 | 53.0 | 50.9 | 50.6 | 47.1 | 44.6 |
| PMPJPE (mm) | 34.9 | 24.9 | 41.0 | 41.0 | 37.6 | 39.5 | 39.4 | 36.0 | 34.9 |
| MPJAS (radians) | NA | 0.150 | NA | 0.270 | 0.239 | NA | 0.250 | 0.227 | 0.213 |

TABLE 6
Detector ablation study. We show the effect of different detector architectures (crosshairs and regression), image resolution, number of crosshair layers for XY dimensions, second-stage refinement and training with feet/hands included. For comparability, evaluation averages when training with feet/hands only includes 17 joints and 15 rotations. MPJPE is reported using both Protocol 1 and Protocol 2 (i.e. Procrustes aligned to ground-truth). We see that higher resolution improves the accuracy, and that our crosshairs significantly ourperforms the baseline regression architecture.

| | MPJAS (15 rot) | MPJPE-17 P1 | P2 | | MPJAS (15 rot) | MPJPE-17 P1 | P2 |
|---|---|---|---|---|---|---|---|
| *256x192 Crosshairs (trained on 17j, 15r)* | | | | *384x288 Regression (trained on 17j, 15r)* | | | |
| Detector to skeleton (1-layer XY) | 0.239 | 52.2 | 39.6 | Detector to skeleton | 0.228 | 52.4 | 40.6 |
| Detector to skeleton (4-layer XY) | 0.235 | 51.9 | 39.1 | | | | |
| Refined detections | NA | 49.2 | 38.0 | Refined detections | NA | 50.0 | 37.6 |
| *384x288 Crosshairs (trained on 17j, 15r)* | | | | *384x288 Crosshairs (trained on 21j, 19r)* | | | |
| Detector to skeleton (1-layer XY) | 0.217 | 48.1 | 37.2 | Detector to skeleton (1-layer XY) | 0.216 | 47.8 | 36.7 |
| Detector to skeleton (4-layer XY) | 0.217 | 48.2 | 37.1 | Detector to skeleton (4-layer XY) | **0.213** | 47.4 | 36.1 |
| Refined detections | NA | 45.5 | 35.8 | Refined detections | NA | **44.6** | **34.9** |

rotations onto the skeleton improves on the previous state-of-the-art, and the refinement stage takes the improvement to 3mm. We also note that the previous state-of-the-art under this protocol, [68], uses a GAN to training their detector, an approach that is complementary to our technique.

Qualitatively, compared to other methods, our approach can also ensure a coherent skeleton - no failure cases with elongated limbs. In **Table 2** we report our results compared to various other methods using similar protocols.

## 5.3 Visualizations

The visualizations provide a more intuitive overview of how our approach performs and demonstrate the rotation information missing from most existing approaches. To express the full range of detections, we rank the full test set from best (low percentile) to worst (high percentile)
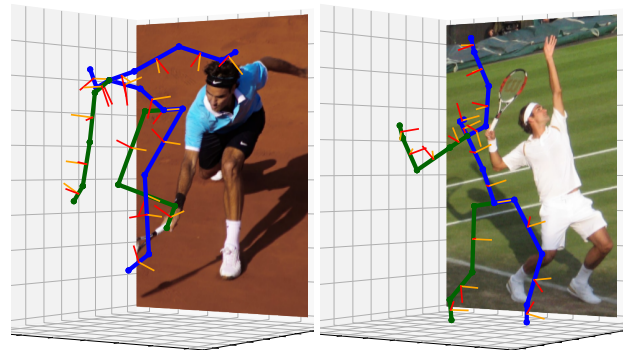


Fig. 5. In the wild example of our 6D human pose prediction (Second photo source: [13])

based on MPJPE-17. In **Figure 6**, we show key percentile predictions, providing the image with the ground truth skeleton and the predictions plotted side-by-side. We also show different angles for a better perspective. To visualize rotations, we plot green and red handles extruding from the bones which show the *forward* and *left* vectors (the bone is always oriented *up*, apart from the legs *down*).

Most of the predictions are quite accurate, both in terms of overall skeletal form and bone orientation. The *25th* percentile and median both included heavily occluded limbs yet the model can make highly accurate predictions. Even the very *worst* result, a pose with the legs heading straight into the camera, is a credible prediction. This example is an outlier, with an error of almost 2x the *99th* percentile.

### 5.4　In-the-wild datasets

We explore how well our approach generalizes to in-the-wild environments by testing on two additional datasets, MPI-3DHP, and MPI-3DPW. We first use the base model from Human 3.6m unchanged and then show a version fine-tuned with a limited number of iterations at 2.5/1.0 e-4 learning rate (head/backbone) on the relevant training set. Testing the unchanged model shows how well orientation keypoints generalizes even when trained on a limited dataset. The fine-tuning allows the model to learn the wider variety of camera angles, distances, and, particularly for 3DPW, learn the significant differences in joint locations.

**Table 3** shows that our method achieves state-of-the-art results on 3DHP. Without training on the 3DHP set, our model is considerably more accurate than the other approaches without 3DHP training and is even competitive with models fully trained on 3DHP. Our refined model, though, does not perform better than our single-stage model predictions. After 2k iterations of fine-tuning our Human 3.6m model with 3DHP mixed into the batches (48/16 samples per batch 3DHP/H3.6m), the accuracy further improves and is state-of-the-art across all metrics by a significant margin. While fine-tuning improves the average metric, the very worst (by MPJPE) failure case worsens further.

In **Table 4** we show (i) results of our H3.6m model without further training, (ii) the same model but using bone lengths from the 3DPW training set, (iii) the results after fine-tuning our model for 8k iterations on the 3DPW training/validation set (48/16 samples 3DPW/MPII) and (iv) using the OKPS from our model as an optimisation signal for the mesh fitting loop of the SPIN model as presented in [30]. For (iv) we simply add an L2 loss to the optimization objective that measures the difference between the SMPL joint positions and angles and those estimated from the OKPs. The model generalizes well despite the very different environment and provides state-of-the-art positional accuracy on this dataset as well. We demonstrate that OKPS can be used to improve other methods, such as SPIN. We also set a benchmark both MPJPE and for accuracy on rotations, 0.302 radians MPJAS, given the noisiness of the provided ground truth, we believe this is a solid result and demonstrates good generalization of our Human 3.6m results. We do note that the method employed to estimate ground-truth for 3DPW has 0.208 radians of error in their test environment, quantitatively comparable to our results

on Human 3.6m, and is likely noisier on the provided in-the-wild data. The test set has significant and visible annotation errors when merely re-projecting the ground truth. Nevertheless, the dataset provides an opportunity to benchmark rotations in a more challenging setting.

We visualize both of these datasets in **Figures 7** and **8**, again showing a range of results based on MPJPE accuracy. The failure cases are illustrative. For 3DHP, the detector gets the entire body reversed for the quantitatively worst prediction. For 3DPW, the worst prediction is for a lunging fencer, and the model rotates the torso circa 120 degrees to reverse the arms. A few frames earlier, before the fencer is at full extension, the model does not make this mistake; incorporating video analysis may resolve these kinds of ambiguities. In other cases, such as 75th and 95th percentile examples, the skeletal scaling is the culprit for MPJPE as the rotations are quite accurate.

We also show two visualizations of in-the-wild images with Table 1. Despite the absence of sports images in the Human3.6m training set, the model has generalized well to a completely different context.

### 5.5　Detector network and ablation

The accuracy of our approach, like that of [41], [53] and others, is impacted by the accuracy of the keypoint detector. The state-of-the-art in detectors is rapidly evolving, and different models offer different accuracy versus computational and training complexity tradeoffs. Higher-resolution input can also improve detector accuracy. We focus this paper on the benefits of using orientation keypoints to solve for 6D human rotations - our method can be applied to different detectors. We therefore explore the most critical aspects of detector impact here in several ways.

**Ground truth analysis**. Following [41] and [53], we use ground truth analysis to provide a baseline as it separates the issue of detection quality, a moving target given rapid advances in the field, and focuses on the information content of orientation keypoints. As shown in **Table 5**, lifting orientation keypoints with ground truth data improves MPJPE by 15mm. This clearly shows the added value of orientation keypoints for human pose estimation irrespective of the detector employed.

**A simpler detector**. We also show, in **Table 5** and **Table 6**, the results of our simple detector, a regression head with Resnet50, on accuracy using both 2D and 3D versions. The detector is less accurate than Crosshairs but still delivers state-of-the-art results as OKPS improves MPJPE accuracy by 4.5mm and 3.4mm.

**Detector accuracy simulation.** We can also use the refiner stage to analyze results by directly scaling the detection errors over the ground truth. We show the impact of detector accuracy on MPJAS and MPJPE in **Figure 9**. This gives a sense of potential improvements from further detector advances as well as potential degradation from a more challenging dataset, and we believe more informative than applying Gaussian noise as it preserves the correlation structure of the detector errors. Using the detections of [41], we plot their stacked hourglass detector as a reference - at their level of accuracy, our method reduces MPJPE by 15%,
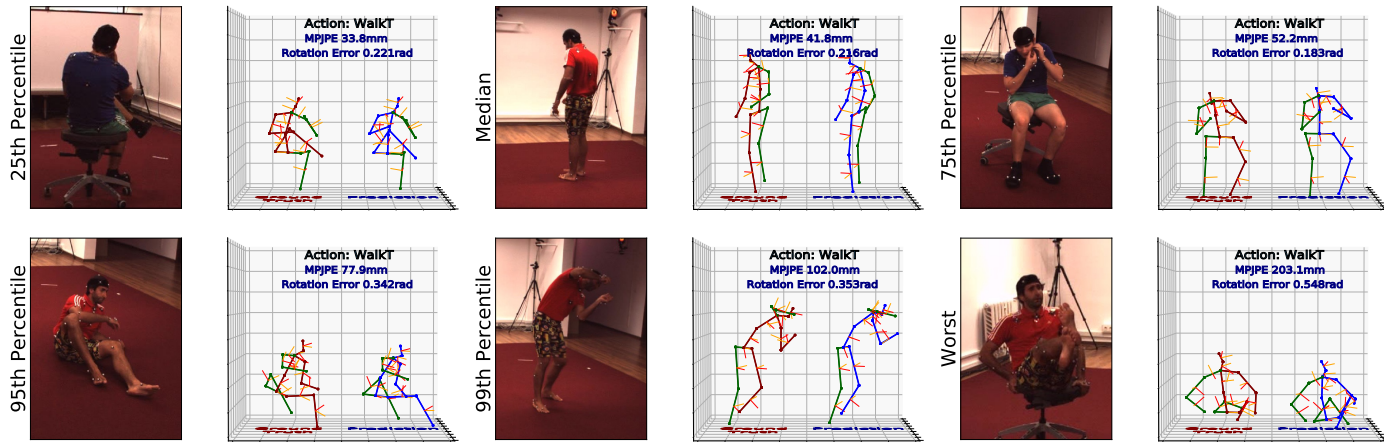
Fig. 6. Model predictions versus ground truth. Kinematic rotations are visualized with protruding bone handles - small red lines are joint forward vector, yellow are left. Samples ranked by Protocol 1 MPJPE on 17 joint skeleton, quoted rotation error is MPJAS-15. We visualize feet and hand predictions but do not include in averages
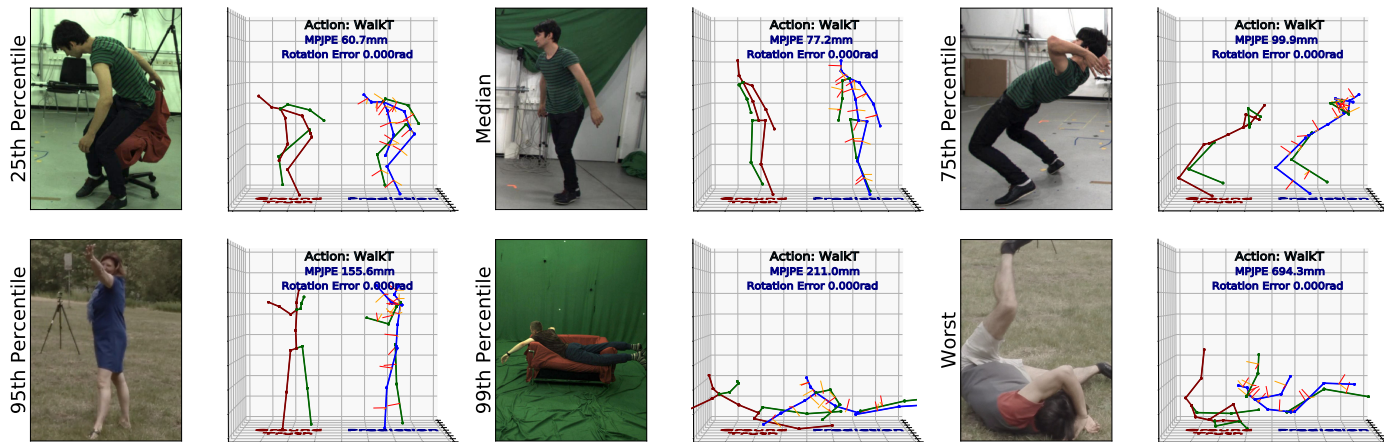


Fig. 7. Model predictions for 6 random scenes on the 3DHP dataset (right skeleton) vs ground truth (left skeleton). We show a range of results including the 25th, 50th, 75th, 95th, 99th percentile and as well the example exhibiting the worst error based on MPJPE accuracy. Our approach works well across this entire range. The failure case in the worst example is due to the fact that the detector has confused the left and right limbs.
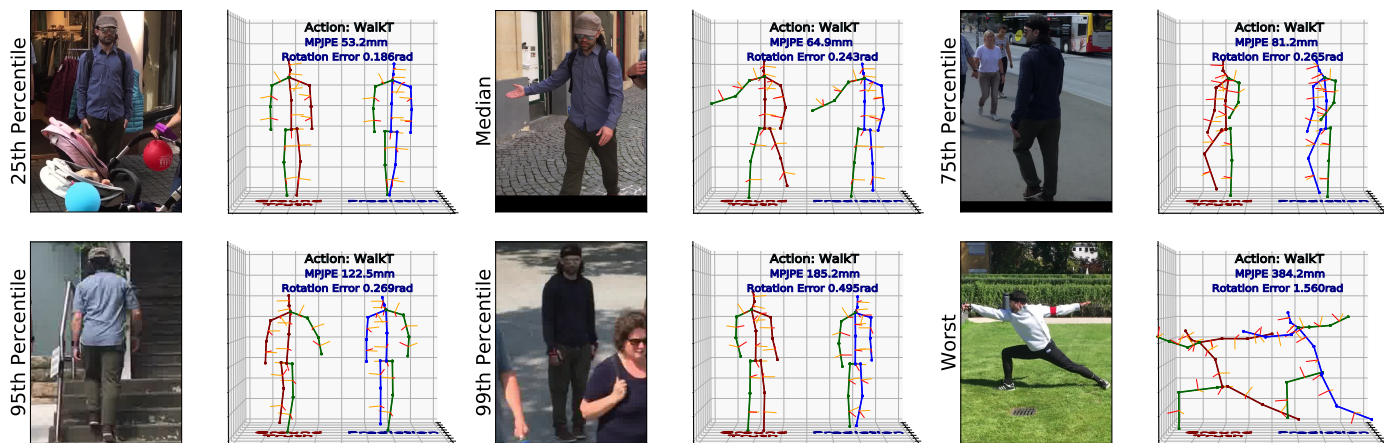


Fig. 8. Model predictions on the 3DPW dataset (right skeleton) vs ground truth (left skeleton). Again we show a range of results including the 25th, 50th, 75th, 95th, 99th percentile and as well the example exhibiting the worst error based on MPJPE accuracy. Overall our approach works well across all the examples, with the failure case again attributed to a swapped limb detection.
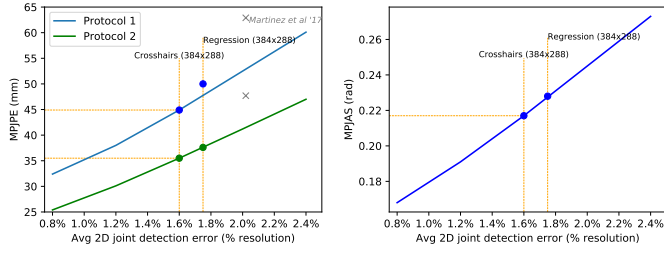
Fig. 9. Orientation keypoints sensitivity to detector accuracy. We plot both detectors we used and the fine-tuned stacked hourglass detections from [41]. The lines are generated from refining detections with errors scaled. Note: detection errors are overstated as 3 actions for S9 are misannotated in source data as highlighted in [23]

while providing full skeletal rotations. Figure 9 also shows that errors in 6D are correlated to errors of the detector.

**Lifting OKPS only from JKPS**. As we show in **Table 5**, we can also use a lifter model to predict both joint and orientation keypoints from only the 2D JKPS: this is sufficient to make an excellent prediction, achieving 0.250 radians MPJAS. This further demonstrates that OKPS is a useful representation to predict rotations. Of course, predicting OKPS from the image benefits from additional clues and is more accurate.

**Further ablation** We decompose the benefits of different resolutions, multi-layer crosshairs, and the extra refiner network in Table 6. Each earlier layer increases the computation cost of flattening as the kernel size increases. While the intermediate supervision has a clear benefit, inference on the validation set of multiple crosshairs only offers a modest detection benefit at the higher resolution. As early layer crosshairs are computationally more expensive, these could be dropped for deployment, leaving a very lightweight head (only 11% incremental multiply-adds over the Resnet-50 backbone). Similarly, the lower resolution version offers competitive accuracy at half the computational cost.

### 5.6 Comparison to mesh-based approaches

Approaches that regress dense mesh correspondences have become an important paradigm for pose estimation, and therefore we provide a numerical comparison to a state-of-the-art method of this type.

In particular, we compare our approach to DaNet [71], which embodies the dense mesh regression paradigm. This method predicts the dense correspondences between 2D pixels and 3D vertices with an HRNet-W48 which has similar capacity to our backbone and then uses this as an intermediate representation to estimate rotations of body joints via a GCN. The results are reported in Table 2. As shown in the table, even though DaNet has a similar capacity, all configurations of our method outperform this approach on H3.6m.

We also compare to HybrIK [35] which is a recent SMPL-based approach. From Table 4 we can see that our approach significantly outperforms HybrIK on 3DPW on all the reported metrics. Additionally, to demonstrate the complementary nature of our research, we used our OKPS predictions as an extra optimisation parameter for the SMPL-based SPIN approach from [30] and improved upon the original results by a substantial 15.1 MPJPE.

## 6 CONCLUSION

We have proposed a novel approach for human pose estimation that makes use of orientation keypoints to parameterize bone orientations and have demonstrated that our method significantly improves upon the state-of-the-art. We can accurately predict skeletal rotations from a single RGB camera image while accurately localizing joints in 3D. Our technique is simple and straightforward to apply, and we believe it can become a vital part of the human pose estimation pipeline.

Our approach offers many opportunities for extension. Further advances in detection methods should improve our results. Exploiting other datasets should improve accuracy, and our method could be used with weak supervision strategies and GANs. Extending our work to video is a natural step as temporal information can help resolve ambiguities. Training to generate a customized skeleton (i.e. bone lengths) could also improve precision. Finally, our post-processing is differentiable and could be directly incorporated for training.

## REFERENCES

[1] Ijaz Akhter and Michael J. Black. Pose-conditioned joint angle limits for 3d human pose reconstruction. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.

[2] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *IEEE Conference on Computer Vision and Pattern Recognition*, June 2014.

[3] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J. Black. Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. In *European Conference on Computer Vision*, 2016.

[4] Ernesto Brau and Hao Jiang. 3d human pose estimation via deep learning from 2d annotations. In *2016 Fourth International Conference on 3D Vision (3DV)*, pages 582–591. IEEE, 2016.

[5] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

[6] Cristian Sminchisescu Catalin Ionescu, Fuxin Li. Latent structured models for human pose estimation. In *International Conference on Computer Vision*, 2011.

[7] Ching-Hang Chen and Deva Ramanan. 3d human pose estimation= 2d pose estimation+ matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7035–7043, 2017.

[8] Xipeng Chen, Kwan-Yee Lin, Wentao Liu, Chen Qian, and Liang Lin. Weakly-supervised discovery of geometry-aware representation for 3d human pose estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[9] Yilun Chen, Zhicheng Wang, Yuxiang Peng, Zhiqiang Zhang, Gang Yu, and Jian Sun. Cascaded pyramid network for multi-person pose estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[10] Yilun Chen, Zhicheng Wang, Yuxiang Peng, Zhiqiang Zhang, Gang Yu, and Jian Sun. Cascaded pyramid network for multi-person pose estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7103–7112, 2018.

[11] Vasileios Choutas, Georgios Pavlakos, Timo Bolkart, Dimitrios Tzionas, and Michael J Black. Monocular expressive body regression through body-driven attention. In *European Conference on Computer Vision*, pages 20–40. Springer, 2020.

[12] Dylan Drover, Ching-Hang Chen, Amit Agrawal, Ambrish Tyagi, and Cong Phuoc Huynh. Can 3d pose be learned from 2d projections alone? In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018.

[13] elyob. Federer v Safin. https://www.flickr.com/photos/82744013@N00/665277802, 2007.

[14] Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. Recurrent network models for human dynamics. In *International Conference on Computer Vision*, pages 4346–4354, 2015.

[15] Xavier Glorot, Antoine Bordes, and Y. Bengio. Deep sparse rectifier neural networks. *Proceedings of the 14th International Conference on Artificial Intelligence and Statisitics (AISTATS) 2011*, 15:315–323, 01 2011.

[16] Ikhsanul Habibie, Weipeng Xu, Dushyant Mehta, Gerard Pons-Moll, and Christian Theobalt. In the wild human pose estimation using explicit 2d features and intermediate 3d representations. In *IEEE Conference on Computer Vision and Pattern Recognition*, June 2019.

[17] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.

[18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European Conference on Computer Vision*, 2016.

[19] Geoffrey E. Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *CoRR*, abs/1207.0580, 2012.

[20] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 448–456, Lille, France, 07–09 Jul 2015. PMLR.

[21] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7):1325–1339, jul 2014.

[22] Umar Iqbal, Pavlo Molchanov, and Jan Kautz. Weakly-supervised 3d human pose learning via multi-view images in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5243–5252, 2020.

[23] Karim Iskakov, Egor Burkov, Victor Lempitsky, and Yury Malkov. Learnable triangulation of human pose. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.

[24] Hao Jiang. 3d human pose reconstruction using millions of exemplars. *2010 20th International Conference on Pattern Recognition*, pages 1674–1677, 2010.

[25] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *Computer Vision and Pattern Regognition (CVPR)*, 2018.

[26] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 12 2014.

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

[28] Muhammed Kocabas, Chun-Hao P. Huang, Otmar Hilliges, and Michael J. Black. PARE: Part attention regressor for 3D human body estimation. In *Proceedings International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2021.

[29] Muhammed Kocabas, Salih Karagoz, and Emre Akbas. Self-supervised learning of 3d human pose using multi-view geometry. In *IEEE Conference on Computer Vision and Pattern Recognition*, June 2019.

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

[31] Christoph Lassner, Javier Romero, Martin Kiefel, Federica Bogo, Michael J. Black, and Peter V. Gehler. Unite the people: Closing the loop between 3d and 2d human representations. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.

[32] Hsi-Jian Lee and Zen Chen. Determination of 3d human body postures from a single view. *Computer Vision, Graphics, and Image Processing*, 30:148–168, 1985.

[33] V. Lepetit, F. Moreno-Noguer, and P. Fua. Epnp: An accurate o(n) solution to the pnp problem, 2008.

[34] Jake Levinson, Carlos Esteves, Kefan Chen, Noah Snavely, Angjoo Kanazawa, Afshin Rostamizadeh, and Ameesh Makadia. An analysis of svd for deep rotation estimation. *arXiv preprint arXiv:2006.14616*, 2020.

[35] Jiefeng Li, Chao Xu, Zhicun Chen, Siyuan Bian, Lixin Yang, and Cewu Lu. Hybrik: A hybrid analytical-neural inverse kinematics solution for 3d human pose and shape estimation. In *CVPR*, 2021.

[36] Kevin Lin, Lijuan Wang, and Zicheng Liu. Mesh graphormer. *arXiv preprint arXiv:2104.00272*, 2021.

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

[38] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, Oct. 2015.

[39] Chenxu Luo, Xiao Chu, and Alan Yuille. Orinet: A fully convolutional network for 3d human pose estimation. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2018.

[40] Diogo C. Luvizon, David Picard, and Hedi Tabia. 2d/3d pose estimation and action recognition using multitask deep learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[41] Julieta Martinez, Rayat Hossain, Javier Romero, and James J. Little. A simple yet effective baseline for 3d human pose estimation. In *International Conference on Computer Vision*, 2017.

[42] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *International Conference on 3DVision*, 2017.

[43] Dushyant Mehta, Oleksandr Sotnychenko, Franziska Mueller, Weipeng Xu, Mohamed Elgharib, Pascal Fua, Hans-Peter Seidel, Helge Rhodin, Gerard Pons-Moll, and Christian Theobalt. XNect: Real-time multi-person 3D motion capture with a single RGB camera. 2020.

[44] Francesc Moreno-Noguer. 3d human pose estimation from a single image via distance matrix regression. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

[45] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, pages 483–499, Cham, 2016. Springer International Publishing.

[46] Aiden Nibali, Zhen He, Stuart Morgan, and Luke Prendergast. 3d human pose estimation with 2d marginal heatmaps. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2019.

[47] Mohamed Omran, Christoph Lassner, Gerard Pons-Moll, Peter V. Gehler, and Bernt Schiele. Neural body fitting: Unifying deep learning and model-based human pose and shape estimation. In *International Conference on 3DVision*, 2018.

[48] Sungheon Park, Jihye Hwang, and Nojun Kwak. 3d human pose estimation using convolutional neural networks with 2d pose information. In Gang Hua and Hervé Jégou, editors, *Computer Vision – ECCV 2016 Workshops*, pages 156–169, Cham, 2016. Springer International Publishing.

[49] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3D hands, face, and body from a single image. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[50] Georgios Pavlakos, Xiaowei Zhou, and Kostas Daniilidis. Ordinal depth supervision for 3D human pose estimation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

[51] Georgios Pavlakos, Xiaowei Zhou, Konstantinos G Derpanis, and Kostas Daniilidis. Coarse-to-fine volumetric prediction for single-image 3D human pose. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.

[52] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3D human pose and shape from a single color image. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

[53] Dario Pavllo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3d human pose estimation in video with temporal convolutions and semi-supervised training. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.

[54] Dario Pavllo, David Grangier, and Michael Auli. Quaternet: A quaternion-based recurrent model for human motion. In *British Machine Vision Conference (BMVC)*, 2018.

[55] Gerard Pons-Moll, Jonathan Taylor, Jamie Shotton, Aaron Hertzmann, and Andrew Fitzgibbon. Metric regression forests for correspondence estimation. *International Journal of Computer Vision*, 113(3):163–175, 2015. Communicated by Tilo Burghardt, Majid Mirmehdi, Walterio Mayol-Cuevas, and Dima Damen.

[56] Mahdi Rad and Vincent Lepetit. Bb8: A scalable, accurate, robust to partial occlusion method for predicting the 3d poses of challenging objects without using depth. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.

[57] Mir Rayat Imtiaz Hossain and James J Little. Exploiting temporal

information for 3d human pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 68–84, 2018.

[58] Xiao Sun, Bin Xiao, Fangyin Wei, Shuang Liang, and Yichen Wei. Integral human pose regression. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, editors, *Computer Vision – ECCV 2018*, pages 536–553, Cham, 2018. Springer International Publishing.

[59] Yu Sun, Qian Bao, Wu Liu, Yili Fu, Black Michael J., and Tao Mei. Monocular, one-stage, regression of multiple 3d people. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021.

[60] Bugra Tekin, Pablo Márquez-Neila, Mathieu Salzmann, and Pascal Fua. Learning to fuse 2d and 3d image cues for monocular body pose estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3941–3950, 2017.

[61] Bugra Tekin, Sudipta N. Sinha, and Pascal Fua. Real-time seamless single shot 6d object pose prediction. In *IEEE Conference on Computer Vision and Pattern Recognition*, June 2018.

[62] Denis Tome, Chris Russell, and Lourdes Agapito. Lifting from the deep: Convolutional 3d pose estimation from a single image. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

[63] Jonathan Tompson, Ross Goroshin, Arjun Jain, Yann LeCun, and Christoph Bregler. Efficient object localization using convolutional networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.

[64] Shinji Umeyama. Least-squares estimation of transformation parameters between two point patterns. *IEEE Trans. Pattern Anal. Mach. Intell.*, 13:376–380, 1991.

[65] Timo von Marcard, Roberto Henschel, Michael Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *European Conference on Computer Vision*, 2018.

[66] Bastian Wandt and Bodo Rosenhahn. Repnet: Weakly supervised training of an adversarial reprojection network for 3d human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7782–7791, 2019.

[67] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *The European Conference on Computer Vision (ECCV)*, September 2018.

[68] Wei Yang, Wanli Ouyang, Xiaolong Wang, Jimmy Ren, Hongsheng Li, and Xiaogang Wang. 3d human pose estimation in the wild by adversarial learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[69] Hashim Yasin, Umar Iqbal, Bjorn Kruger, Andreas Weber, and Juergen 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.

[70] Yusuke Yoshiyasu, Ryusuke Sagawa, Ko Ayusawa, and Akihiko Murai. Skeleton transformer networks: 3d human pose and skinned mesh from single rgb image. In C.V. Jawahar, Hongdong Li, Greg Mori, and Konrad Schindler, editors, *Computer Vision – ACCV 2018*, pages 485–500, Cham, 2019. Springer International Publishing.

[71] Hongwen Zhang, Jie Cao, Guo Lu, Wanli Ouyang, and Zhenan Sun. Learning 3d human shape and pose from dense body parts. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[72] Xingyi Zhou, Qixing Huang, Xiao Sun, Xiangyang Xue, and Yichen Wei. Towards 3d human pose estimation in the wild: A weakly-supervised approach. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.

[73] Xingyi Zhou, Xiao Sun, Wei Zhang, Shuang Liang, and Yichen Wei. Deep kinematic pose regression. In Gang Hua and Hervé Jégou, editors, *Computer Vision – ECCV 2016 Workshops*, pages 186–201, Cham, 2016. Springer International Publishing.

[74] Xiaowei Zhou, Menglong Zhu, Georgios Pavlakos, Spyridon Leonardos, Konstantinos G. Derpanis, and Kostas Daniilidis. Monocap: Monocular human motion capture using a cnn coupled with a geometric prior. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(4):901–914, Apr. 2019.

[75] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.