

Monocular 3D Human Pose Estimation In The Wild Using Improved CNN Supervision

Dushyant Mehta¹, Helge Rhodin², Dan Casas³, Pascal Fua²,
Oleksandr Sotnychenko¹, Weipeng Xu¹, and Christian Theobalt¹

¹MPI for Informatics, Germany ²EPFL, Switzerland ³Universidad Rey Juan Carlos, Spain

Abstract

We propose a CNN-based approach for 3D human body pose estimation from single RGB images that addresses the issue of limited generalizability of models trained solely on the starkly limited publicly available 3D pose data. Using only the existing 3D pose data and 2D pose data, we show state-of-the-art performance on established benchmarks through transfer of learned features, while also generalizing to in-the-wild scenes. We further introduce a new training set for human body pose estimation from monocular images of real humans that has the ground truth captured with a multi-camera marker-less motion capture system. It complements existing corpora with greater diversity in pose, human appearance, clothing, occlusion, and viewpoints, and enables an increased scope of augmentation. We also contribute a new benchmark that covers outdoor and indoor scenes, and demonstrate that our 3D pose dataset shows better in-the-wild performance than existing annotated data, which is further improved in conjunction with transfer learning from 2D pose data. All in all, we argue that the use of transfer learning of representations in tandem with algorithmic and data contributions is crucial for general 3D body pose estimation.

1. Introduction

We present an approach to estimate the 3D articulated human body pose from a single image taken in an uncontrolled environment. Unlike marker-less 3D motion capture methods that *track* articulated human poses from *multi-view* video sequences, [75, 61, 62, 72, 6, 20, 63, 13] or use *active* RGB-D cameras [57, 5], our approach is designed to work from a single low-cost RGB camera.

Data-driven approaches using Convolutional Neural Networks (CNNs) have shown impressive results for 3D pose regression from monocular RGB, however, in-the-wild

scenes and motions remain challenging. Aside from the difficulty of the 3D pose estimation problem, it is further stymied by the lack of suitably large and diverse annotated 3D pose corpora. For 2D joint detection it is feasible to obtain ground truth annotations on in-the-wild data on a large scale through crowd sourcing [55, 4, 30], consequently leading to methods that generalize to in-the-wild scenes [16, 74, 70, 71, 45, 11, 7, 42, 37, 22, 24, 12]. Some 3D pose estimation approaches take advantage of this generalizability of 2D pose estimation, and propose to lift the 2D keypoints to 3D [69, 76, 9, 73, 36, 80, 83, 79, 60, 59, 14]. This approach however is susceptible to errors from depth ambiguity, and often requires computationally expensive iterative pose optimization. Recent advances in direct CNN-based 3D regression show promise, utilizing different prediction space formulations [65, 35, 81, 44, 40] and incorporating additional constraints [81, 67, 83, 78]. However, we show on a new in-the-wild benchmark that existing solutions have a low generalization to in-the-wild conditions. They are far from the accuracy seen for 2D pose prediction in terms of correctly located keypoints.

Existing 3D pose datasets use marker-based motion capture, MoCap, for 3D annotation [27, 58], which restricts recording to skin-tight clothing, or markerless systems in a dome of hundreds of cameras [32], which enables diverse clothing but requires an expensive studio setup. Synthetic data can be generated by retargeting MoCap sequences to 3D avatars [15], however the results lack realism, and learning based methods pick up on the peculiarities of the rendering leading to poor generalization to real images.

Our contributions towards accurate in-the-wild pose estimation are twofold. First, in Section 4, we explore the use of transfer learning to leverage the highly relevant mid- and high-level features learned on the readily available in-the-wild 2D pose datasets [4, 31] in conjunction with the existing annotated 3D pose datasets. Our experimentally validated mechanism of feature transfer shows better accuracy and generalizability compared to naïve weight initialization from 2D pose estimation networks and domain adaptation based approaches. With this we show previously unseen levels of accuracy on established benchmarks, as well as

This work was funded by the ERC Starting Grant project CapReal (335545). Dan Casas was supported by a Marie Curie Individual Fellow grant (707326), and Helge Rhodin by the Microsoft Research Swiss JRC. We thank The Foundry for license support.

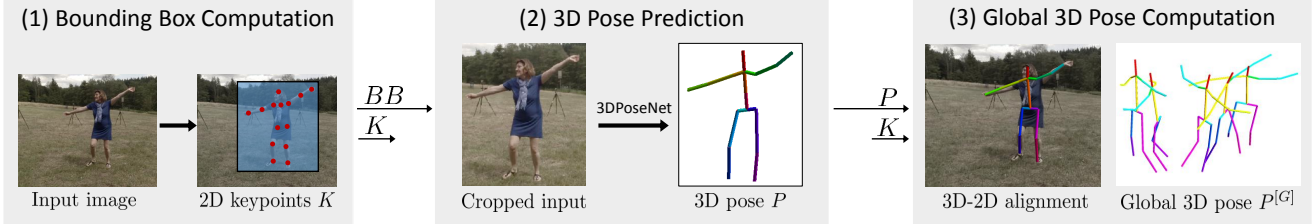


Figure 1. We infer 3D pose from single image in three stages: (1) extraction of the actor bounding box from 2D detections; (2) direct CNN-based 3D pose regression; and (3) global root position computation in original footage by aligning 3D to 2D pose.

generalizability to in-the-wild scenes, with only the existing 3D pose datasets.

Second, in Section 5, we introduce the new MPI-INF-3DHP dataset¹ real humans with ground truth 3D annotations from a state-of-the-art markerless motion capture system. It complements existing datasets with everyday clothing appearance, a large range of motions, interactions with objects, and more varied camera viewpoints. The data capture approach eases appearance augmentation to extend the captured variability, complemented with improvements to existing augmentation methods for enhanced foreground texture variation. This gives a further significant boost to the accuracy and generalizability of the learned models.

The data-side supervision contributions are complemented by CNN architectural supervision contributions in Section 3.2, which are orthogonal to in-the-wild performance improvements.

Furthermore, we introduce a new test set, including sequences outdoors with accurate annotation, on which we demonstrate the generalization capability of the proposed method and validate the value of our new dataset.

The components of our method are thoroughly evaluated on existing test datasets, demonstrating both state-of-the-art results in controlled settings and, more importantly, improvements over existing solutions for in-the-wild sequences thanks to the better generalization of the proposed techniques.

2. Related Work

There has been much work on learning- and model-based approaches for human body pose estimation from monocular images, with much of the recent progress coming through CNN based approaches. We review the most relevant approaches, and discuss their relation with our work.

3D pose from 2D estimates. Deep CNN architectures have dramatically improved 2D pose estimation [28, 42], with even real-time solutions [74]. Graphical models [19, 1] continue to find use in modeling multi-person relations [45]. 3D pose can be inferred from 2D pose through geometric and statistical priors [41, 64]. Optimization

of the projection of a 3D human model to the 2D predictions is computationally expensive and ambiguous, but the ambiguity can be addressed through pose priors and it further allows incorporation of various constraints such as inter-penetration constraints [9], sparsity assumptions [73, 80, 82], joint limits [17, 2], and temporal constraints [50]. Simo-Serra *et al.* [60] sample noisy 2D predictions to ambiguous 3D shapes, which they disambiguate using kinematic constraints, and improve discriminative 2D detection from likely 3D samples [59]. Li *et al.* look up the nearest neighbours in a learned joint embedding of human images and 3D poses [36] to estimate 3D pose from an image. We choose to use the geometric relations between the predicted 2D and 3D skeleton pose to infer the global subject position.

Estimating 3D pose directly. Additional image information, e.g. on the front-back orientation of limbs, can be exploited by regressing 3D pose directly from the input image [65, 35, 81, 26]. Deep CNNs achieve state-of-the-art results [81, 66, 44]. While CNNs dominate, regression forests have also been used to derive 3D *posebit descriptors* efficiently [47]. The input and output representations are important too. To localize the person, the input image is commonly cropped to the bounding box of the subject before 3D pose estimation [26]. Video input provides temporal cues, which translate to increased accuracy [67, 83]. The downside of conditioning on motion is the increased input dimensionality, and requires motion databases with sufficient motion variation, which are even harder to capture than pose data sets. In controlled conditions, fixed camera placement provides additional height cues [78]. Since monocular reconstruction is inherently scale-ambiguous, 3D joint positions relative to the pelvis, with normalized subject height are widely used as the output. To explicitly encode dependencies between joints, Tekin *et al.* [65] regressing to a high-dimensional pose representation, learned by an auto encoder. Li *et al.* [35] report that predicting positions relative to the parent joint of the skeleton improves performance, but we show that a pose-dependent combination of absolute and relative positions leads to further improvements. Zhou *et al.* [81] regress joint angles of a skeleton from single images, using a kinematic model.

Addressing the scarcity and limited appearance variability of datasets. Learning-based methods require large

¹MPI-INF-3DHP dataset available at gvv.mpi-inf.mpg.de/3dhp-dataset

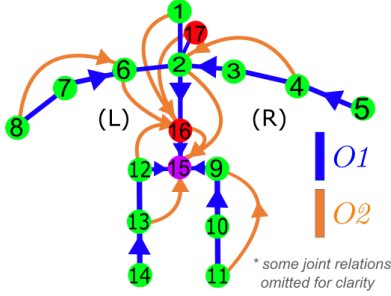


Figure 2. 3D pose, represented as a vector of 3D joint positions, is expressed variously as 1) P : relative to the root (joint #15), 2) $O1$ (blue): relative to first order and, 3) $O2$ (orange): relative to second order parents in the kinematic skeleton hierarchy.

annotated dataset corpora. 3D annotation [26] is harder to obtain than 2D pose annotation. Some approaches treat 3D pose as a hidden variable, and use pose priors and projection to 2D to guide the training [10, 76]. Rogez *et al.* render mosaics of in-the-wild human pose images using projected mocap data [52]. Chen *et al.* [15] render textured rigged human models, but still require domain adaptation to in-the-wild images for generalization. Other approaches use the estimated 2D pose to look up a suitable 3D pose from a dictionary [14], or use the ground truth 2D pose based dictionary lookup to create 3D annotations for in-the-wild 2D pose data [53], but neither address the 2D to 3D ambiguity. Our new dataset complements the existing datasets, through extensive appearance and pose variation, by using marker-less annotation and provides an increased scope for augmentation.

Transfer Learning [43] is commonly used in computer vision to leverage features and representations learned on one task to offset data scarcity for a related task. Low and/or mid-level CNN features can be shared also among unrelated tasks [56, 77]. Pretraining on ImageNet [54] is commonly used for weight initialization [25, 66] in CNNs. We explore different ways of using the low and mid-level features learned on in-the-wild 2D pose datasets for further improving the generalization of 3D pose prediction models.

3. CNN-based 3D Pose Estimation

We start by introducing the network architecture, utilized input and output domains, and notation. While the particularities of our architecture are explained in Section 3.2, our main contributions towards in-the-wild conditions are covered in sections 4 and 5.

Given an RGB image, we estimate the global 3D human pose $P^{[G]}$ in the camera coordinate system. We estimate the global positions of the joints of the skeleton depicted in Figure 2, accounting for the camera viewpoint, which goes beyond only estimating in a root-centered (pelvis) coordinate system, as is common in many previous works. Our algorithm consists of three steps, as illustrated in Figure 1.

(1) the subject is localized in the frame with a 2D bounding box BB , computed from 2D joint heatmaps H , obtained with a CNN we call *2DPoseNet*; (2) the root-centered 3D pose P is regressed from the BB -cropped input with a second CNN termed *3DPoseNet*; and (3) global 3D pose coordinates $P^{[G]}$ and perspective correction are computed in closed form using 3D pose P , 2D joint locations K and known camera calibration.

3.1. Bounding Box and 2D Pose Computation

We use our *2DPoseNet* to produce 2D joint location heatmaps H . The heat map maxima provide the most likely 2D joint locations K which can also act as a stand-in person bounding-box BB detector. See Figure 1. The 2D joint locations K are further used for global pose estimation in Section 3.3. In case of an alternative BB detector, K comes from *3DPoseNet*. See *2D Auxiliary Task* in Figure 3.

Our *2DPoseNet* is fully convolutional and is trained on MPII [4] and LSP [31, 30] datasets. We use a CNN structure based on Resnet-101 [23], up to the filter banks at level 4. Striding is removed at level 5, and features in the *res5a* block are halved and identity skip connections removed from *res5b* and *res5c*. For specifics of the network architecture and the training scheme, refer to the supplementary document.

3.2. 3D Pose Regression

The 3D pose CNN, termed *3DPoseNet*, is used to regress root-centered 3D pose P from a cropped RGB image, and makes use of new CNN supervision techniques. Figure 3 depicts the main components of the method, detailed in the following sections.

Network The base network derives from Resnet-101 as well, and is identical to *2DPoseNet* up to *res5a*. We remove the remaining layers from level 5. A 3D prediction stub S comprised of a convolution layer ($k_{5 \times 5}, s_2$) with 128 features and a final fully-connected layer that outputs the 3D joint locations is added on top. Additionally we predict 2D heatmaps H as an auxiliary task after *res5a* and, use intermediate supervision with pose P at *res3b3* and *res4b22*. Refer to the supplementary for specifics of the loss weights for the intermediate and auxiliary tasks.

3.2.1 Multi-level Corrective Skip Connections

We additionally use a skip connection scheme as a training-time regularization architecture. We add skip connections from *res3b3* and *res4b20* to the main prediction P_{deep} , leading to P_{sum} . In contrast to vanilla skip-connections [38], we compare both P_{sum} and P_{deep} to the ground truth, and remove the skip connections after training. We show the improvements due to this approach Section 6.

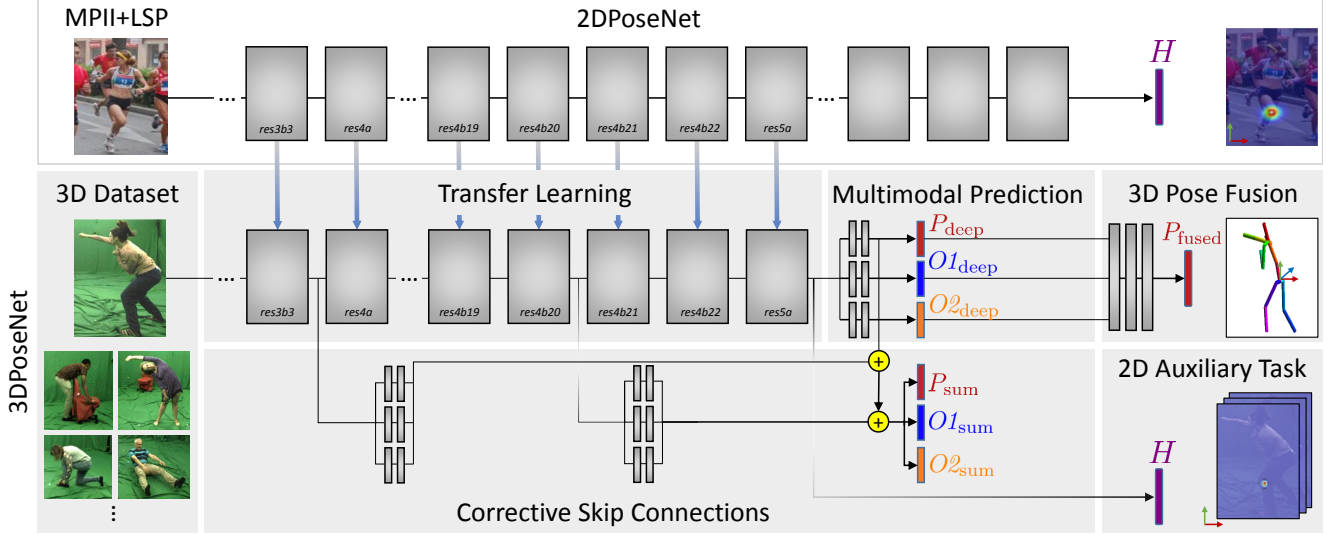


Figure 3. 3D pose Training overview. The main components are 1) regularization through corrective skip connections, and 2D pose prediction as auxiliary task, 2) Multi-modal 3D pose prediction and fusion, 3) a new marker-less 3D pose database with appearance augmentation, and 4) Transfer learning from features learned for 2D pose estimation.

3.2.2 Multi-modal Pose Fusion

Formulating joint location prediction relative to a single local or global location is not always optimal. Existing literature [35] has observed that predicting joint locations relative to their direct kinematic parents (Order 1 parents) improves performance. Our experiments reveal that to not universally hold true. We find that depending on the pose and the visibility of the joints in the input image, the optimal relative joint for each joint’s location prediction differs. Hence, we use joint locations P relative to the root, $O1$ relative to Order 1 parents and $O2$ relative to Order 2 parents along the kinematic tree as the *three modes* of prediction, see Figure 2, and fuse them with fully-connected layers.

For the joint set we consider, the kinematic relationships chosen suffice, as it puts at least one reference joint for each joint in the relatively low entropy torso [33]. We use three identical 3D prediction stubs attached to `res5a` for predicting the pose as P , $O1$ and $O2$, and for each we use corrective skip connections. These predictions are fed into a smaller network with three fully connected layers, to implicitly determine and fuse the better constraints per joint into the final prediction P_{fused} . The network has the flexibility to emphasize different combinations of constraints depending on the pose. This can be viewed as intermediate supervision with auxiliary tasks, yet the separate streams for predicting each mode individually are key to its efficacy.

3.3. Global Pose Computation

The bounding box cropping normalizes subject size and position, which frees 3D pose regression from having to localize the person in scale and image space, but loses global pose information. We propose a lightweight and efficient

way to reconstruct the global 3D pose $P^{[G]} = (R|T) P_{\text{fused}}$ from pelvis-centered pose P_{fused} , camera intrinsics, and K .

Perspective correction. The bounding box cropping can be interpreted as using a virtual camera, rotated towards the crop center and its field of view covering the crop area. Since the *3DPoseNet* only ‘sees’ the cropped input, its predictions live in this rotated view, leading to a consistent orientation error in P_{fused} . To compensate, we compute rotation R that rotates the virtual camera to the original view.

3D localization. We seek the global translation T that aligns P_{fused} and K under perspective projection. We assume weak perspective projection, Π , and solve the linear least squares equation $\sum_i \|K^i - \Pi(T + P_{\text{fused}}^i)\|^2$, where i indexes the joints. This assumption yields global position

$$T = \frac{\sqrt{\sum_i \|P_{[xy]}^i - \bar{P}_{[xy]}\|^2}}{\sqrt{\sum_i \|K^i - \bar{K}\|^2}} \begin{pmatrix} \bar{K}_{[x]} \\ \bar{K}_{[y]} \\ f \end{pmatrix} - \begin{pmatrix} \bar{P}_{[x]} \\ \bar{P}_{[y]} \\ 0 \end{pmatrix}, \quad (1)$$

in terms of distances to the 3D mean \bar{P} and 2D mean \bar{K} over all joints. $P_{[xy]}$ is the x, y part of P_{fused} and single subscripts indicate the respective elements. Please see the supplemental document for the derivation and evaluation.

Our solution can be considered a generalization of *procrustes analysis* for projective alignment. Note that this is different to *perspective-n-point* 6DOF rigid pose estimation [34], structure-from-motion, and from the convex approach of Zhou *et al.* [80], which require iterative optimization.

4. Transfer Learning

We use the features learned with Resnet-101 from ImageNet [54] to initialize both *2DPoseNet* and *3DPoseNet*,

Table 1. Evaluation of the mechanisms of transfer learning from *2DPoseNet* to *3DPoseNet* that were explored in the context of the *Base* network. The table compares the effect of various learning rate multiplier combinations for different parts of the network. For network details, refer to Section 3.2. Human3.6m, Subjects 1,5,6,7,8 used for training, and every 64th frame of 9,11 used for testing. * = weights randomly initialized

Learning Rate Multiplier			Total MPJPE (mm)
up to res4b22	res5a	3D Stub S	
1	1	1*	118.7
1/10	1/10	1*	84.6
1/1000	1/1000	1*	89.2
1/10	1	1*	90.7
1/1000	1	1*	80.7

as common for many vision tasks. While this affords a faster convergence while training, there remains room for improved generalization beyond the gains from potential supervision and dataset contributions. Due to the similarity of the tasks, features learned for 2D pose estimation on in-the-wild MPII and LSP training sets can be transferred to 3D pose estimation. We explore different variants of the, thus far, un-utilized method of improving generalization by transferring weights from *2DPoseNet* to *3DPoseNet*.

A naïve initialization of the weights of *3DPoseNet* is inadequate, and there is a tradeoff to be made between the preservation of transferred features and learning new pertinent features. We achieve this through a learning rate discrepancy between the transferred layers and the new layers. We experimentally determine the mechanism for this transfer of features through validation. Table 1 shows the evaluated mechanisms for transfer from *2DPoseNet*. Based on the experiments, we choose to scale down the learning rate of the layers till res4b22 by a factor of 1000. Through similar experiments for the transfer of ImageNet features, we choose to scale down the learning rate of layers till res4b22 by 10.

The same approach can be applied to other network architectures, and our experiments on the learning rate discrepancy serve as a sound starting point for the determination of the transfer learning mechanism. Unlike jointly training with annotated 2D and 3D pose datasets, this approach has the advantage of not requiring the 2D annotations to be consistent between the two datasets, and one can simply use off-the-shelf trained 2D pose networks. In Section 6 we show that our approach outperforms domain adaptation, see Table 5, first row. Additionally, Table 1 validates that the common fine-tuning of the fully-connected layers (third row) and fine-tuning of the complete network (first row) is much less effective than the proposed scheme.

5. MPI-INF-3DHP: Human Pose Dataset

We propose a new dataset captured in a multi-camera studio with ground truth from commercial marker-less mo-



Figure 4. MPI-INF-3DHP dataset. We capture actors using a markerless multi-camera in a green screen studio (left), compute masks for different regions (center left) and augment the captured footage by compositing different textures to the background, chair, upper and lower body areas, independently (center right and right).

tion capture [68]. No special suits and markers are needed, allowing the capture of motions wearing everyday apparel, including loose clothing. In contrast to existing datasets, we record in green screen studio to allow automatic segmentation and augmentation. We recorded 8 actors (4m+4f), performing 8 activity sets each, ranging from walking and sitting to complex exercise poses and dynamic actions, covering more pose classes than Human3.6m. Each activity set spans roughly one minute. Each actor features 2 sets of clothing split across the activity sets. One clothing set is *casual everyday apparel*, and the other is *plain-colored* to allow augmentation.

We cover a wide range of viewpoints, with five cameras mounted at chest height with a roughly 15° elevation variation similar to the camera orientation jitter in other datasets [15]. Another five cameras are mounted higher and angled down 45°, three more have a top down view, and one camera is at knee height angled up. Overall, from all 14 cameras, we capture >1.3M frames, 500k of which are from the five chest high cameras. We make available both true 3D annotations, and a skeleton compatible with the “universal” skeleton of Human3.6m.

Dataset Augmentation. Although our dataset has more clothing variation than other datasets, the appearance variation is still not comparable to in-the-wild images. There have been several approaches proposed to enhance appearance variation. Pishchulin *et al.* warp human size in images with a parametric body model [46]. Images can be used to augment background of recorded footage [49, 15, 27]. Rhodin *et al.* [49] recolor plain-color shirts while keeping the shading details, using intrinsic image decomposition to separate reflectance and shading [39].

We provide chroma-key masks for the background, a chair/sofa in the scene, as well as upper and lower body segmentation for the plain-colored clothing sets. This provides an increased scope for foreground and background augmentation, in contrast to the marker-less recordings of Joo *et al.* [32]. For background augmentation, we use images sampled from the internet. For foreground augmentation, we use a simplified intrinsic decomposition. Since for plain colored clothing the intensity variation is solely due to shading, we use the average pixel intensity as a surrogate for the shading component. We composite cloth like textures with

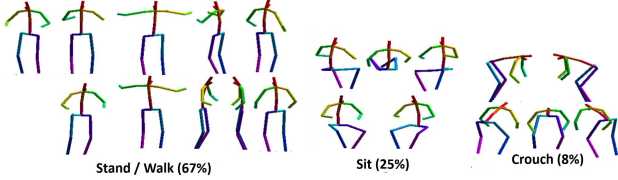


Figure 5. Representative poses (centroids) of the 20 K-means pose clusters of the Human3.6m test set (subjects S9,S11), visually grouped into three broad pose classes, which are used also to perform per-class evaluation. Upright poses are dominant, with complex poses such as sitting and crouching only accounting for 25% and 8% of the poses respectively. Our multimodal fusion scheme significantly improves the latter two, yielding a 3.5mm improvement for Sit and 5.5mm for Crouch class.

the pixel intensity of the upper body, lower body and chair marks independently, for a photo-realistic result. Figure 4 shows example captured and augmented frames.

Test Set. We found the existing test sets for (monocular) 3D pose estimation to be restricted to limited settings due to the difficulty of obtaining ground truth labels in general scenes. HumanEva [58] and Human3.6m [27] are recorded indoors and test on similar looking scenes as the training set, the Human3D+ [15] test set was recorded with sensor suits that influence appearance and lacks global alignment, and the MARCoNI set [17] is markerless through manual annotation, but shows mostly walking motions and multiple actors, which are not supported by most monocular algorithms. We create a new test set with ground truth annotations coming from a multi-view markerless motion capture system. It complements existing test sets with more diverse motions (standing/walking, sitting/reclining, exercise, sports (dynamic poses), on the floor, dancing/miscellaneous), camera view-point variation, larger clothing variation (*e.g.* dress), and outdoor recordings from Robertini *et al.* [51] in unconstrained environments. This makes the test set suitable for testing the generalization of various methods. See Figure 6 for a representative sample. We use the “universal” skeleton for evaluation.

Alternate Metric. In addition to the Mean Per Joint Position Error (MPJPE) widely used in 3D pose estimation, we concur with [27] and suggest a 3D extension of the Percentage of Correct Keypoints (PCK) [71, 70] metric used for 2D Pose evaluation, as well as the Area Under the Curve (AUC) [25] computed for a range of PCK thresholds. These metrics are more expressive and robust than MPJPE, revealing individual joint mispredictions more strongly. We pick a threshold of 150mm, corresponding to roughly half of head size, similar what is used in MPII 2D Pose dataset. We propose evaluating on the common set of joints across 2D and 3D approaches (joints 1-14 in Figure 2), to ensure evaluation compatibility with existing approaches. Joints are grouped by bilateral symmetry (ankles, wrists, shoulders, etc), and can be evaluated by scene setting or activity class.



Figure 6. Representative frames from MPI-INF-3DHP test set. We cover a variety of subjects with a diverse set of clothing and poses in 3 different settings: studio with green screen (right); studio without green screen (left); and outdoors (center).

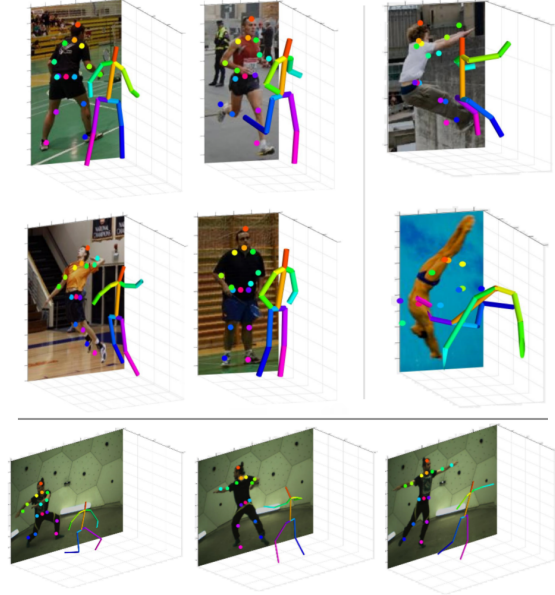


Figure 7. Qualitative evaluation on representative frames of the LSP test set. We succeed in challenging cases (left), with only few failure cases (right). The *Dance1* sequence of the PanopticDataset [32], is also well reconstructed (bottom).

6. Experiments and Evaluation

We evaluate the contributions proposed in the previous sections using the standard datasets Human3.6m and HumanEva, as well as our new MPI-INF-3DHP test set. Additionally, we qualitatively observe the performance on LSP [30] and the CMU Panoptic [32] datasets, demonstrating robustness to general scenes. Refer to Figure 7. Also refer to the supplementary video for global 3D pose results.

We evaluate the impact of training *3DPoseNet* on Human3.6m, and unaugmented and augmented variants of MPI-INF-3DHP, both with and without transfer learning from *2DPoseNet*. We only use Human3.6m compatible camera views from MPI-INF-3DHP for training. Further details are in the supplemental document.

6.1. Impact of Supervision Methods

Multi-level corrective skip connections. In Table 2 we compare a baseline method without any skip connections, a network with vanilla skip connections, and our proposed

Table 2. Activity-wise results (MPJPE in mm) on Human3.6m [27]. Adding our model components one-by-one on top of the *Base* network shows successive improvement of the total accuracy. Significant relative improvements greater than 5mm are underlined. Models are trained on Human3.6m, with network weights initialized from ImageNet, unless specified otherwise. The version marked with MPI-INF-3DHP is trained with Human3.6m and MPI-INF-3DHP. Evaluation with all 17 joints, on every 64th frame, without rescaling to a person specific skeleton.

	Direct	Discuss	Eating	Greet	Phone	Posing	Purch.	Sitting	Sit Down	Smoke	Take Photo	Wait	Walk	Walk Dog	Walk Pair	Total
Base + Regular Skip	113.34	112.26	97.40	110.50	108.63	112.09	105.67	125.97	173.41	109.34	120.87	107.75	97.30	126.05	117.45	115.29
Base	98.98	100.14	86.07	101.83	101.34	96.74	94.89	125.28	158.31	100.21	112.49	99.57	83.39	109.61	95.79	104.32
+ Corr. Skip	<u>92.57</u>	99.08	85.46	<u>95.43</u>	96.93	<u>89.56</u>	95.67	123.54	160.98	97.13	<u>107.56</u>	<u>93.86</u>	<u>76.99</u>	110.93	<u>88.73</u>	101.09
+ Fusion	93.80	99.17	84.73	95.60	94.48	89.40	93.15	119.94	<u>154.61</u>	95.94	106.09	94.13	77.25	108.82	87.38	99.79
+ Transfer <i>2DPoseNet</i>	<u>59.69</u>	<u>69.74</u>	<u>60.55</u>	<u>68.77</u>	76.36	<u>59.05</u>	<u>75.04</u>	96.19	<u>122.92</u>	<u>70.82</u>	<u>85.42</u>	<u>68.45</u>	54.41	<u>82.03</u>	59.79	74.14
+ MPI-INF-3DHP	57.51	68.58	59.56	67.34	78.06	56.86	69.13	99.98	117.53	69.44	82.40	67.96	55.24	76.50	61.40	72.88

Table 3. Evaluation by scene-setting of our design choices on MPI-INF-3DHP test set with weight transfer from ImageNet. Training on our markerless dataset improves accuracy significantly, in particular with the proposed augmentation strategy. Fusion yields an additional gain. *GS* indicates sequences with green screen.

3D dataset	Network architecture	Studio GS	Studio no GS	Outdoor	All	
		3DPCK	3DPCK	3DPCK	3DPCK	AUC
Human3.6m	Base + Corr. Skip	22.2	33.9	18.5	25.1	8.7
	Base + Corr. Skip + Fusion	22.3	34.2	20.0	26.0	9.5
Ours Unaug.	Base + Corr. Skip	66.9	38.2	27.9	46.8	20.9
	Base + Corr. Skip + Fusion	67.6	39.6	28.5	47.8	21.8
Ours Aug.	Base + Corr. Skip	71.1	51.7	36.1	55.4	26.0
	Base + Corr. Skip + Fusion	73.5	53.1	37.9	57.3	28.0

corrective skip regularization on Human3.6m test set. We observe that networks using vanilla skip connections perform markedly worse than the baseline, while corrective skip connections yield more than 5mm improvement for 7 classes of activities (marked as underlined). We verified that the effect is not due to a higher effective learning rate seen by the core network due to the additional loss term.

Multimodal prediction and fusion. The multi-modal fusion scheme yields noticeable improvement across all datasets tested in tables 2 and 3. Since upright poses dominate in pose datasets, and the activity classes are often diluted significantly by upright poses, the true extent of improvement by the multi-modal fusion scheme is masked. To show that the fusion scheme indeed improves challenging pose classes, we cluster the Human3.6m test set by pose as shown in Figure 5, which visualizes the centroid of each cluster. Then we group the clusters visually into three pose classes, namely Stand/Walk, Sit and Crouch, going by the cluster representatives. For the Stand/Walk class, adding fusion has minimal effect, going from 88.4mm to 88.8mm. However, for Sit class fusion leads to a 3.5mm improvement, from 118.9mm to 115.4mm. Similarly, Crouch class has the highest improvement of 5.5mm, going from 156mm to 150.5mm. The improvement is not simply due to additional training, and is less pronounced if predicting *P*, *O1* and *O2* with a common stub, even with more features in the

Table 4. Comparison of results on Human3.6m [27] with the state of the art. Human3.6m, Subjects 1,5,6,7,8 used for training, and 9,11 used for testing. ^S = Scaled to test subject specific skeleton, computed from T-pose. ^T = Uses Temporal Information, ^{J14/J17} = Joint set evaluated, ^A = Uses Best Alignment To GT per frame, ^{Act} = Activitywise Training, ^{1/10/64} = Test Set Frame Sampling

Method	Total MPJPE (mm)
Deep Kinematic Pose[81] ^{J17,B}	107.26
Sparse. Deep. [83] ^{T,J17,B,10,Act}	113.01
Motion Comp. Seq. [67] ^{T,J17,B}	124.97
LinKDE [27] ^{J17,B,Act}	162.14
Du et al. [78] ^{T,J17,B}	126.47
Rogez et al. [52] ^{(J13),B,64}	121.20
SMPLify [9] ^{J14,B,A,(First cam.)}	82.3
3D=2D+Matching [14] ^{J17,B}	114.18
Distance Matrix [40] ^{J17,B}	87.30
Volumetric Coarse-Fine[44] ^{J17,B,S*}	71.90
LCR-Net [53] ^{J17,B}	87.7
Full model (w/o MPI-INF-3DHP) ^{J17,B}	74.11
Full model (w/o MPI-INF-3DHP) ^{J17,B,S}	68.61
Full model (w/o MPI-INF-3DHP) ^{J14,B,A}	54.59

fully-connected layer. Details in the supplementary.

6.2. Transfer Learning

Our approach of transferring representations from *2DPoseNet* to *3DPoseNet* yields 64.7% 3DPCK on MPI-INF-3DHP test-set when trained with only Human3.6m data, compared to 63.7% 3DPCK of the model trained on our augmented training set without transfer learning. It also shows state of the art performance on Human3.6m test set with an error of ≈ 74 mm, demonstrating the dual advantage of the approach in improving both the accuracy of pose estimation and generalizability to in-the-wild scenes. Combining our dataset and transfer learning leads to the best results at $\approx 72.5\%$ 3DPCK. See Table 5.

In contrast to existing approaches countering data scarcity, transfer learning does not require complex dataset synthesis, yet exceeds the performance of Chen *et al.* [15] (with synthetic data and domain adaptation, 28.8% 3DPCK, after procrustes alignment) and our base model trained with

Table 5. Evaluation on MPI-INF-3DHP test set with weight transfer from *2DPoseNet*, by scene setting. Training our full model on our dataset paired with Human3.6m yields best accuracy over all. GS indicates sequences with green screen background.

3D dataset	Method	Studio GS	Studio no GS	Outdoor	All	
		3DPCK	3DPCK	3DPCK	3DPCK	AUC
Human3.6m	Domain adapt.	44.1	42.6	35.2	41.4	17.7
	Ours (full model)	70.8	62.3	58.5	64.7	31.7
Ours Aug.	Ours (full model)	82.6	66.7	62.0	71.7	36.4
Ours Unaug.	Ours (full model)	84.1	68.9	59.6	72.5	36.9
Ours Aug.	Ours, w/o persp. corr.	81.9	68.6	67.4	73.5	37.6
+ Human3.6m	Ours, w/o GT BB	80.4	71.2	69.8	74.4	39.6
	Ours (full model)	84.6	72.4	69.7	76.5	40.8

the synthetic data of Rogez *et al.* [52] (21.7% 3DPCK). Our approach also performs better than domain adaptation [21] to in-the-wild data (Table 5). Details in the supplementary.

6.3. Benefit of MPI-INF-3DHP

Evaluating on MPI-INF-3DHP test-set, without any transfer learning from *2DPoseNet*, we see in Table 3 that our dataset, even without augmentation, leads to a $\approx 9\%$ 3DPCK improvement on outdoor scenes over Human3.6m. However, our augmentation strategy is crucial for improved generalization, as seen from the gains in 3DPCK across scene settings in Table 3, giving 57.3% 3DPCK overall.

Even when combined with transfer learning, we see in Table 5 that our dataset (both augmented and unaugmented) consistently performs better than Human3.6m. The best performance of 76.5% 3DPCK on MPI-INF-3DHP test set and of 72.88mm on Human3.6m is obtained when the two datasets are combined with transfer learning.

6.4. Other Components

Bounding box computation. On MPI-INF-3DHP test set, we additionally evaluate our best performing network using bounding boxes computed from *2DPoseNet*. As shown in Table 5, the performance drops to 74.4% 3DPCK from 76.5% 3DPCK due to the additional difficulty.

Perspective correction. Table 5 shows that perspective correction also has a significant impact, without which, the performance drops to 73% 3DPCK from 76.5%.

6.5. Quantitative Comparison

Human3.6m. Table 4 shows comparison of our method with existing methods, all trained on Human3.6m. Altogether, with our supervision contributions and transfer learning, we are the state of the art (74.11mm, without scaling), while also generalizing to in-the-wild scenes. Note that the Volumetric coarse to fine approach [44] requires estimates of the bone lengths to convert their predictions from pixels to 3D space. Complementing Human3.6m with our

augmented MPI-INF-3DHP dataset further reduces the error to 72mm.

HumanEva. The improvements on Human3.6m are confirmed with a 30.8 and 33.5 MPJPE score on the S1 Box and Walk sequences of HumanEva, after alignment. See supplemental document.

MPI-INF-3DHP. We also evaluated some of the existing methods on our test set. Deep Kinematic Pose [81], attains 13.8% 3DPCK overall. Our full model attains significantly higher accuracy: without transfer learning and trained on Human3.6m obtains 26% 3DPCK, and 64.7% 3DPCK with transfer learning. The large discrepancy in performance between Human3.6m and our new in-the-wild test set highlights the importance of a new benchmark to test generalization to natural images and motions.

7. Discussion

Despite the demonstrated competitive results, our method and others have limitations. Most training sets, also [15], have a strong bias towards chest height cameras. Thus, estimating 3D pose from starkly different camera views is still a challenge. Our new dataset provides diverse viewpoints, which can support development towards viewpoint invariance in future methods. Similar to related approaches, our per-frame estimation exhibits temporal jitter on video sequences. In future, we will investigate integration with model-based temporal tracking to further increase accuracy and temporal smoothness. At less than 250 ms per frame, our approach is much faster than model based methods which work offline in the order of minutes. There still remains scope for improvement towards real time, through smaller input resolution and shallower networks.

We also show that joining forces with transfer learning, in conjunction with algorithmic and data contributions, will aid progress in 3D pose estimation in many different directions, such as overall accuracy and generalizability.

8. Conclusion

We have presented a fully feedforward CNN-based approach for monocular 3D human pose estimation that attains state-of-the-art on established benchmarks [27, 58] and quantitatively outperforms existing methods on the introduced in-the-wild benchmark. State of the art is attained with enhanced CNN supervision techniques and improved parent relationships in the kinematic chain. Transfer learning from in-the-wild 2D pose data in tandem with a new dataset that includes a larger variety of real and augmented human appearances, activities and camera views, leads to the significantly improved generalization to in-the-wild images. Our method is also the first to efficiently extract global 3D position in non-cropped images, without time consuming iterative optimization.

Supplemental Document: Monocular 3D Human Pose Estimation In The Wild Using Improved CNN Supervision

This document accompanies the main paper, and the supplemental video.

1. Further Discussion of Design Choices Regarding Multi-modal Fusion

To demonstrate that the improvement seen due to the fusion scheme is not simply a result of fine tuning, we compare the result of fusion with components successively removed. Using P , $O1$ and $O2$, we get an MPJPE of 74.49mm on Human3.6m. On removing $O2$, the error increases to 74.77mm, and on removing both $O1$ and $O2$, the error increases to 75.27mm. The comparison here is without any multi-level corrective skip training.

For P , $O1$ and $O2$ to have different modes of mispredictions, the underlying feature set that they are computed from has to be as different as possible, because each is related to the other with a linear transform. We achieve some degree of decorrelation between the three by using 3 different prediction stubs, one each for P , $O1$ and $O2$ with a convolutional layer ($k_{5 \times 5}$, s_2) with 128 features followed by a fully-connected layer. If we replace these three stubs with a single stub with the convolutional layer having 256 features followed by a fully-connected layer, the resulting MPJPE is 75.30mm after fusion, in contrast to an MPJPE of 74.49mm from fusing the result of 3 prediction stubs. Both of these are without corrective-skip connections.

2. Further Discussion of Multi-level Corrective Skip

Since our multi-level corrective skip scheme adds an additional loss at the last stage (X_{deep} , where X is $P/O1/O2$) of the network, it increases the effective learning rate seen by the core network. To verify that the improvements seen due to the proposed scheme are not caused by this difference in the effective learning rate, we trained a version of the *Base* network with loss weights as the sum of the loss weights for X_{deep} and X_{sum} specified in Table 4. We find that this network performs worse than the *Base* network (107.14mm vs 104.32mm MPJPE on Human3.6m), and does not approach the accuracy attained with multi-level corrective skip scheme (101.09mm).

3. Global Pose Computation

3.1. 3D localization

In this section we describe a simple, yet very efficient, method to compute the global 3D location T of a noisy 3D point set P with unknown global position. We assume

known scaling and orientation parameters, obtained from its 2D projection estimate K in a camera with known intrinsics parameters (focal length f). We further assume that the point cloud spread in depth direction is negligible compared to its distance z_0 to the camera and approximate perspective projection of an object near position $(x_0, y_0, z_0)^\top$ with weak perspective projection (linearizing the pinhole projection model at z_0):

$$\begin{pmatrix} u \\ v \end{pmatrix} = \Pi \begin{pmatrix} x \\ y \\ z \end{pmatrix}, \text{ with } \Pi = \begin{pmatrix} \frac{f}{z_0} & 0 & 0 \\ 0 & \frac{f}{z_0} & 0 \end{pmatrix}. \quad (1)$$

Estimates K and P are assumed to be noisy due to estimation errors. We find the optimal global position T in the least squares sense, by minimizing $T = \arg \min_{(x,y,z)} E(x,y,z)$, with

$$\begin{aligned} E &= \sum_i \|K^i - \Pi((x,y,z)^\top + P^i)\|^2 \\ &= \sum_i \|K^i - \frac{f}{z}((x,y)^\top + P_{[xy]}^i)\|^2, \end{aligned} \quad (2)$$

where P^i and K^i denote the i th joint position in 3D and 2D, respectively, and $P_{[xy]}^i$ the xy component of P^i . It has partial derivative

$$\frac{\partial E}{\partial x} = \frac{2f}{z} \sum_i K_{[x]}^i + \frac{f}{z} (P_{[x]}^i - x), \quad (3)$$

where $P_{[x]}$ denotes the x part of P , and \bar{P} the mean of P over all joints. Solving $\frac{\partial E}{\partial x} = 0$ gives the unique closed-form solutions $x = \bar{K}_{[x]} \frac{z}{f} - \bar{P}_{[x]}$ and equivalently $y = \bar{K}_{[y]} \frac{z}{f} - \bar{P}_{[y]}$, for $\frac{\partial E}{\partial y} = 0$.

Substitution of x and y in E and differentiating with respect to z yields

$$\begin{aligned} \frac{\partial E}{\partial z} &= \frac{f \sum_i (K^i - \bar{K})^\top (P_{[xy]}^i - \bar{P}_{[xy]})}{z^2} \\ &\quad + \frac{f^2 \sum_i \|P_{[xy]}^i - \bar{P}_{[xy]}\|^2}{z^3}. \end{aligned} \quad (4)$$

Finally, solving $\frac{\partial E}{\partial z} = 0$ gives the depth estimate

$$\begin{aligned} z &= f \frac{\sum_i \|P_{[xy]}^i - \bar{P}_{[xy]}\|^2}{\sum_i (K^i - \bar{K})^\top (P_{[xy]}^i - \bar{P}_{[xy]})} \\ &\approx f \frac{\sqrt{\sum_i \|P_{[xy]}^i - \bar{P}_{[xy]}\|^2}}{\sqrt{\sum_i \|K^i - \bar{K}\|^2}}, \end{aligned} \quad (5)$$

where $(K^i - \bar{K})(P^i - \bar{P}) = \|K^i - \bar{K}\| \|P^i - \bar{P}\| \cos(\theta)$ is approximated for $\theta \approx 0$. This is a valid assumption in our case, since the rotation of 3D and 2D pose is assumed to be matching.

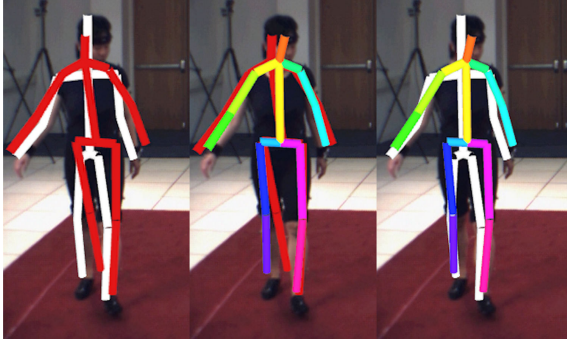


Figure 1. The predicted pose (red) is inaccurate for positions away from the camera center (left), compared against the ground truth (white). Perspective correction (colored) corrects the orientation (center) and is closer to the ground truth (right). Here tested on the walking sequence of HumanEva S1.

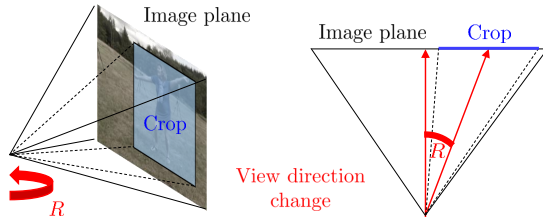


Figure 2. Sketch of the input image cropping and resulting change of field of view. The corresponding rotation R of the view direction is sketched in 2D on the right.

Evaluation on HumanEva: In addition to evaluating centered pose P , we evaluate the global 3D pose prediction $P^{[G]}$ on the widely used HumanEva motion capture dataset — *Box* and *Walk* sequences of *Subject 1* from the validation set. Note that we do not use any data from HumanEva for training. We significantly improve the state of the art for the *Box* sequence (82.1mm [9] vs 58.6mm). Results on the *Walk* sequence are of higher accuracy than *Bogo et al.* [9], but lower than the accuracy of *Bo et al.* [8] and *Yasin et al.* [76], who, however train on HumanEva [8] or use an example database dominated by walking motions [76]. Our skeletal structure does not match that of HumanEva, e.g. the head prediction has a consistent frontal offset and the hip is too wide. To compensate, we compute a linear map of dimension 14×14 (number of joints) that maps our joint positions as a linear combination to the HumanEva structure. The same mapping is applied at every frame, but is computed only once, jointly on the *Box* and *Walk* sequence, to limit the correction to global inconsistencies of the skeleton structure. This fine-tuned result is marked by \sim in Table 1.

3.2. Perspective correction

Our *3DPoseNet* predicts pose P in the coordinate system of the bounding box crop, which leads to inaccuracies

as shown in Figure 1. The cropped image appears if as it was taken from a virtual camera with the same origin as the original camera, but with view direction to the crop center, see Figure 2. To map the reconstruction from the virtual camera coordinates to the original camera, we rotate P by the rotation R between the virtual and original camera. Since the existing training sets provide chest-height camera placements with the same viewpoint, the bias in vertical direction is already learned by the network. We apply perspective correction only in horizontal direction, where a change in cropping and yaw rotation of the person cannot be distinguished by the network. R is then the rotation around the camera up direction by the angle between the original and the virtual view direction, see Figure 2. On our MPI-INF-3DHP test set perspective correction improves the PCK by 3 percent points. On HumanEva the improvement is up to 3 mm MPJPE, see Table 1. The correction is most pronounced for cameras with a large field of view, e.g. Go-Pro and similar outdoor cameras, and when the subject is located at the border of the view. Using the vector from the camera origin to the centroid of 2D keypoints K as the virtual view direction was most accurate in our experiments. However, the crop center can be used instead. Opposed to the Perspective-n-Point algorithm applied by Zhou *et al.* [83], any regression method that works on cropped images could immediately profit from this perspective correction, without computing 2D keypoint detections.

4. CNN Architecture and Training Specifics

4.1. 2DPoseNet

Architecture: The architecture derives from Resnet-101, using the same structure as is until level 4. Since we are interested in predicting heatmaps, we remove striding at level 5. Additionally, the number of features in the *res5a* module are halved, identity skip connections are removed from *res5b* and *res5c*, and the number of features gradually tapered to 15 (heatmaps for 14 joints + root). As shown in Table 2, for *2DPoseNet*, our results on MPII and LSP test sets approach that of the state of the art.

Intermediate Supervision: Additionally, we employ intermediate supervision at *res4b20* and *res5a*, treating the first 15 feature maps of the layers as the intermediate joint-location heatmaps. Further, we use a Multi-level Corrective Skip scheme, with skip connections coming from *res3b3* and *res4b22* through prediction stubs comprised of a 1×1 convolution with 20 feature maps followed by a 3×3 convolution with 15 outputs.

Training: For training, we use the Caffe [29] framework, with the AdaDelta solver with a momentum of 0.9 and weight decay rate of 0.005. We employ a batch size of 7, and use Euclidean Loss everywhere. For the Learning Rate and Loss Weight taper schema, refer to Table 3.

Table 1. Quantitative evaluation on HumanEva-I [58], with different alignment strategies used in the literature. For reference, we also show multi-view existing results. Our models use no data from HumanEva for training, while the other methods listed train/finetune on HumanEva-I. * = Does not use GT Bounding Box information. † = Translation alignment only. \sim = trained or fine-tuned on HumanEva-I.

		S1 Box			S1 Walk		
		$P^{[G]}$ (global)	P (align ^{S,T})	P (align ^{R,S,T})	$P^{[G]}$ (global)	P (align ^{S,T})	P (align ^{R,S,T})
Monocular	Our full model*	117.1	80.5	58.6	121.1	81.0	67.2
	w/o Persp. correct.*	116.1	79.4	58.6	123.9	83.6	67.3
	Our full model* \sim	77.9	38.4	30.8	89.1	48.2	33.5
	Zhou <i>et al.</i> [83] \sim	-	-	-	-	-	34.2
	Bo <i>et al.</i> [8]* \sim	-	-	-	-	54.8†	-
	Yasin <i>et al.</i> [76]* \sim	-	-	-	52.2	-	-
	Bogo <i>et al.</i> [9] \sim	-	-	82.1	-	-	73.3
	Akhter <i>et al.</i> [2]	-	-	165.5	-	-	186.1
Multiview	Ramakris. <i>et al.</i> [48]	-	-	151.0	-	-	161.8
	Amin <i>et al.</i> [3]	47.7	-	-	54.5	-	-
	Rhodin <i>et al.</i> [50]	59.7	-	-	74.9	-	-
	Elhayek <i>et al.</i> [18]	60.0	-	-	66.5	-	-

Table 2. Results of our 2DPoseNet on MPII Single Person Pose [4] dataset and LSP [30] 2D Pose datasets. * = Trained/Finetuned only on the corresponding training set

	MPII		LSP	
	PCK _{h0.5}	AUC	PCK _{0.2}	AUC
Our 2DPoseNet				
w Person Locali.	89.7	61.3	91.2	65.3
w/o Person Locali.	89.6	61.5	91.2	65.5
Stacked Hourgl.[42]	90.9*	62.9*	-	-
Bulat <i>et al.</i> [11]	89.7*	59.6*	90.7	-
Wei <i>et al.</i> [74]	88.5	61.4	90.5	65.4
DeeperCut [25]	88.5	60.8	90.1	66.1
Gkioxary <i>et al.</i> [22]	86.1*	57.3*	-	-
Lifshitz <i>et al.</i> [37]	85.0	56.8	84.2	-
Belagiannis <i>et al.</i> [7]	83.9*	55.5*	85.1	-
DeepCut[45]	82.4	56.5	87.1	63.5
Hu&Ramanan [24]	82.4*	51.1*	-	-
Carreira <i>et al.</i> [12]	81.3*	49.1*	72.5*	-

4.2. 3DPoseNet

Architecture: The core network is identical to 2DPoseNet up to res5a. A 3D Prediction stub is attached on top, comprised of a 5×5 convolution layer with a stride of 2 and 128 features, followed by a fully-connected layer.

Multi-level Corrective Skip: We attach 3D prediction stubs to res3b3 and res4b20, similar to the final prediction stub, but with 96 convolutional features instead of 128. The resulting predictions are added to P_{deep} to get P_{sum} . We add a loss term to P_{deep} in addition to the loss term at P_{sum} .

Table 3. Loss weight and learning rate, LR, taper scheme used for 2DPoseNet. 2DPoseNet also employs Multi-level Corrective Skip connections, and the heatmap H_{sum} is the sum of H_{deep} and the skip connections. Heatmaps H_{4b20} and H_{5a} are used for intermediate supervision.

Base LR	# Iter	Loss Weights ($w \times L(H_{xx})$)			
		H_{sum}	H_{deep}	H_{4b20}	H_{5a}
0.050	60k	1.0	0.5	0.5	0.5
0.010	60k	1.0	0.4	0.1	0.1
0.005	60k	1.0	0.2	0.05	0.05
0.001	60k	1.0	0.2	0.05	0.05
6.6e-4	60k	1.0	0.1	0.005	0.005
0.0001	40k	1.0	0.01	0.001	0.001
2.5e-5	40k	1.0	0.001	0.0001	0.0001
0.0008	60k	1.0	0.0001	0.0001	0.0001
0.0001	40k	1.0	0.0001	0.0001	0.0001
3.3e-5	20k	1.0	0.0001	0.0001	0.0001

Multi-modal Fusion: We add prediction stubs for $O1$ and $O2$, similar to those for P . Note that the predictions for P , $O1$ and $O2$ are done with distinct stubs, and this slight decorrelation of predictions is important. These predictions are at a later finetuning step fed into three fully-connected layers, with 2k, 1k and 51 nodes respectively.

Intermediate Supervision: We use intermediate supervision at 4b5 and res4b20, using prediction stubs comprised of 7×7 convolution with a stride of 3 and 128 features, followed by a fully-connected layer predicting P , $O1$ and $O2$ as a single vector. Additionally, we predict joint location heatmaps and part-label maps using a 1×1

Table 4. Loss weight and LR taper scheme used for *3DPoseNet*. There is a difference in the number of iterations used when training with Human3.6m or MPI-INF-3DHP alone, v.s. when training with the two in conjunction. Part Labels *PL* are used only when training with H3.6m solely. Multi-level skip connections add up with X_{deep} to yield X_{sum} , where X is *P* or *O1 O2*.

Base LR	H3.6m/Our Batch = 5 #Epochs	H3.6m+Our Batch = 6 #Epochs	Loss Weights ($w \times L(A_{bb})$) $X = P/O1/O2$						<i>H</i>	<i>PL*</i>
			X_{4b5}	X_{4b20}	X_{deep}	X_{sum}				
0.05	3 (45k)	2.4 (60k)	50	50	50	100		0.1	0.05	
0.01	1 (15k)	1.2 (30k)	10	10	10	100		0.05	0.025	
0.005	2 (30k)	1.2 (30k)	5	5	5	100		0.01	0.005	
0.001	1 (15k)	0.6 (15k)	1	1	1	100		0.01	0.005	
5e-4	2 (30k)	1.2 (30k)	0.5	0.5	0.5	100		0.005	0.001	
1e-4	1 (15k)	0.6 (15k)	0.1	0.1	0.1	100		0.005	0.001	

Table 5. Loss weight and LR taper scheme used for fine tuning *3DPoseNet* for Multi-modal Fusion scheme.

Base LR	H3.6m/Our Batch = 5 #Epochs	H3.6m+Our Batch = 6 #Epochs	Loss Weights ($w \times L(A_{bb})$) P_{fused}
0.05	(1k)	(2k)	100
0.01	1 (15k)	0.8 (20k)	100
0.005	1 (15k)	0.8 (20k)	100
0.001	1 (15k)	0.8 (20k)	100

convolution layer after *res5a* as an auxiliary task. We don’t use the part-label maps when training with MPI-INF-3DHP dataset.

Training: For training, the solver settings are similar to *2DPoseNet*, and we use Euclidean Loss everywhere. For transfer learning, we scale down the learning rate of the transferred layers by a factor determined by validation. For fine-tuning in the multi-modal fusion case, we similarly downscale the learning rate of the trained network by 10,000 with respect to the three new fully-connected layers. For the learning rate and loss weight taper schema for both the main training and multi-modal fusion fine-tuning stages, refer to Tables 4 and 5. We use different training durations when using Human3.6m or MPI-INF-3DHP in isolation, versus when using both in conjunction. This is reflected in the aforementioned tables.

4.2.1 3D Pose Training Data

In the various experiments on *3DPoseNet*, for the datasets we consider, we select $\approx 37.5k$ frames for each, yielding $\approx 75k$ samples after scale augmentation at 2 scales (0.7 and 1.0).

Human3.6m: We use the H80k [26] subset of Human3.6m, and train with the “universal” skeleton, using subjects S1,5,6,7,8 for training and S9,11 for testing. The predicted skeleton is not scaled to the test subject skeletons at test time.

MPI-INF-3DHP: For our dataset, to maintain compatibil-

ity of view with Human3.6m and other datasets, we only pick the 5 chest high cameras for all 8 subjects, sampling frames such that at least one joint has moved by more than 200mm between selected frames. A random subset of these frames is used for training, to match the number of selected Human3.6m frames.

MPI-INF-3DHP Augmented: The augmented version uses the same frames as the unaugmented MPI-INF-3DHP above, keeping $\approx 25\%$ frames unaugmented, $\approx 40\%$ with only BG and Chair augmentation, and the rest with full augmentation.

4.2.2 Domain Adaptation To In The Wild 2D Pose Data

We use a domain adaptation stub comprised of $conv_{3 \times 3, 256}$, $conv_{3 \times 3, 128}$, fc_{64} and fc_1 layers, and cross entropy domain classification loss. It uses Ganin *et al.*’s [21] gradient inversion approach. The domain adaptation stub is attached after *res4b22* in the network. We found that directly starting out with $\lambda = -1$ performs better than gradually increasing the magnitude of λ with increasing iterations. We train on the Human3.6m training set, with 2D heatmap and part label prediction as auxiliary tasks. Images from MPII [4] and LSP [30, 31] training sets are used without annotations for learning better generalizable features. The generalizability is improved, as evidenced by the 41.4 3DPCK on MPI-INF-3DHP test set, but does not match up with the 64.7 3DPCK attained using transfer learning. Detailed results in main Table 3.

5. MPI-INF-3DHP Dataset

We cover a wide range of poses in our training and test sets, roughly grouped into various activity classes. A detailed description of the dataset is available in Section 4 of the main paper. In addition, Figure 3 samples the various different activity classes, augmentation and subjects represented in our dataset.

Similarly for the test set, we show a sample of the activities and the variety of subjects in Figure 4.

5.1. The Challenge of Learning Invariance to Viewpoint Elevation

In this paper, we only consider the cameras in the training set placed at chest-height, in part to be compatible with the existing datasets, and in part because viewpoint elevation invariance is a significantly more challenging problem. Existing benchmarks do not place emphasis on this. We will release an expanded version of our MPI-INF-3DHP testset with multiple camera viewpoint elevations, to complement the training data.

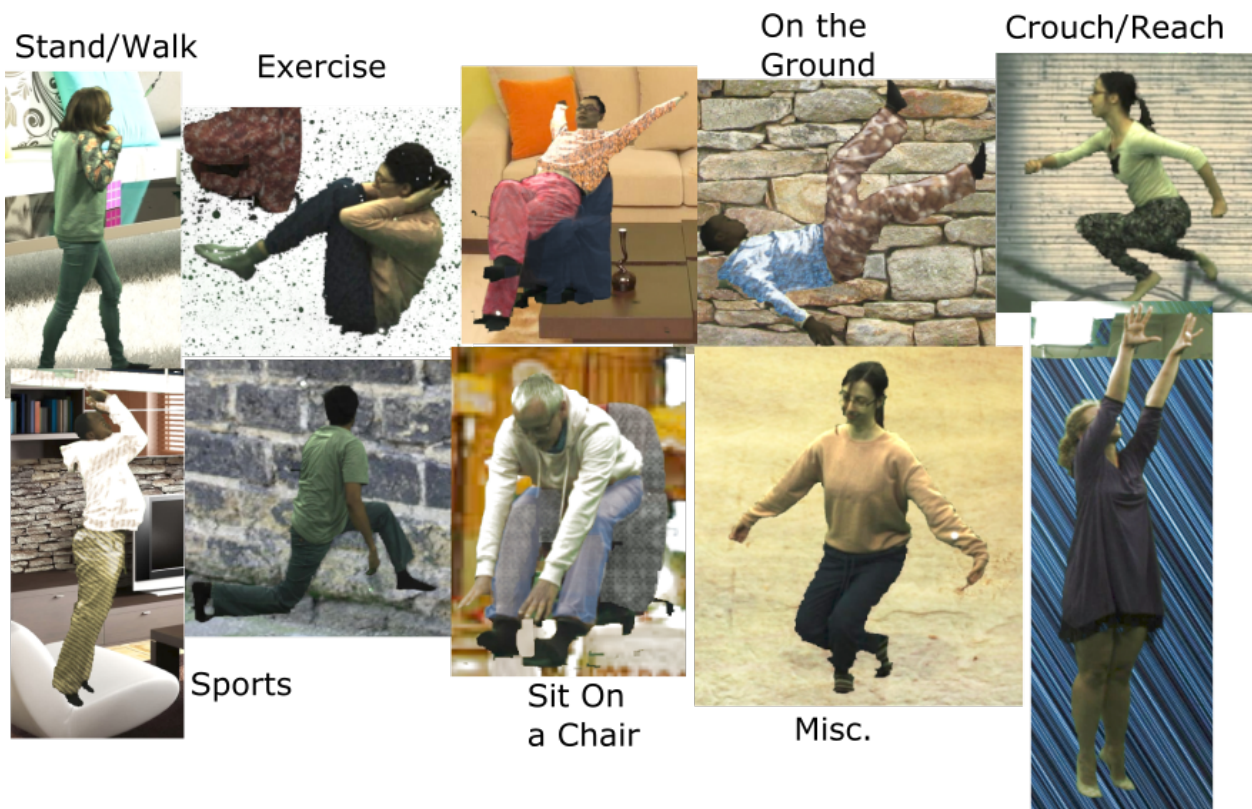


Figure 3. A sample of the activities, clothing, subjects as well as augmentation on MPI-INF-3DHP Training Set.

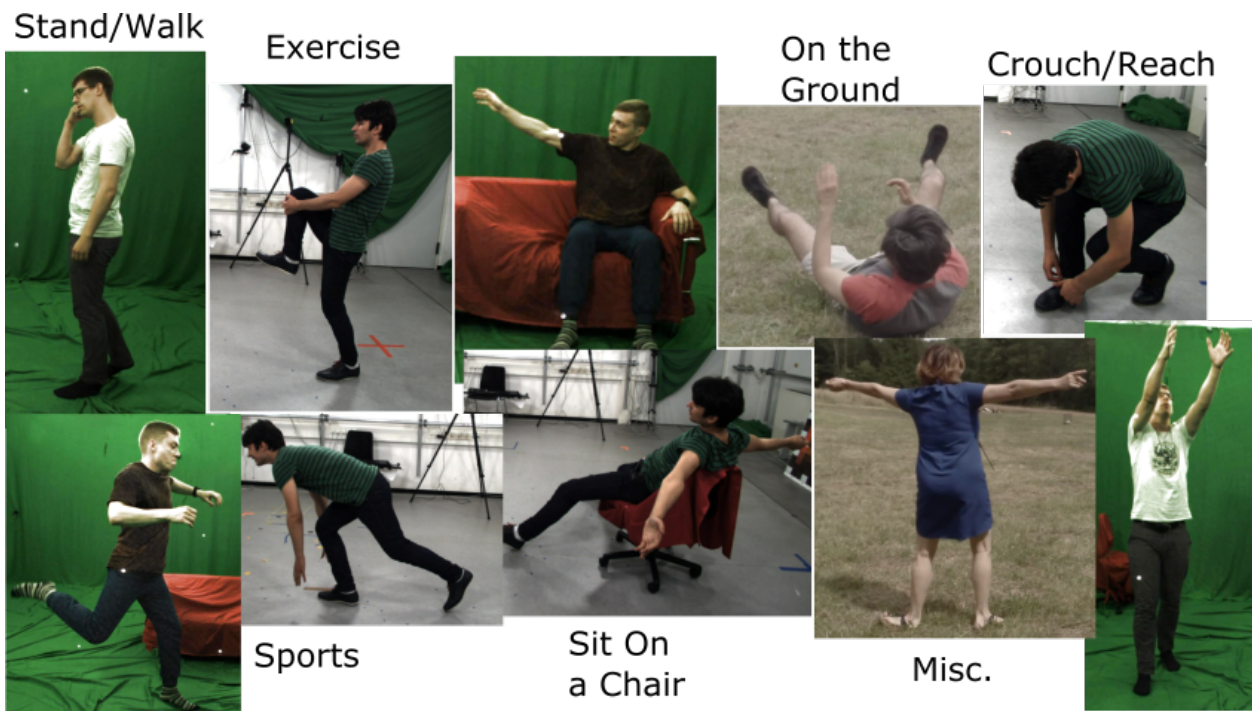


Figure 4. A sample of the activities and subjects in the test set of MPI-INF-3DHP

References

- [1] A. Agarwal and B. Triggs. Recovering 3d human pose from monocular images. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 28(1):44–58, 2006. [2](#)
- [2] I. Akhter and M. J. Black. Pose-conditioned joint angle limits for 3d human pose reconstruction. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1446–1455, 2015. [2](#), [11](#)
- [3] S. Amin, M. Andriluka, M. Rohrbach, and B. Schiele. Multi-view pictorial structures for 3D human pose estimation. In *BMVC*, 2013. [11](#)
- [4] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele. 2D Human Pose Estimation: New Benchmark and State of the Art Analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014. [1](#), [3](#), [11](#), [12](#)
- [5] A. Baak, M. Müller, G. Bharaj, H.-P. Seidel, and C. Theobalt. A data-driven approach for real-time full body pose reconstruction from a depth camera. In *IEEE International Conference on Computer Vision (ICCV)*, 2011. [1](#)
- [6] A. O. Balan, L. Sigal, M. J. Black, J. E. Davis, and H. W. Haussecker. Detailed human shape and pose from images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–8, 2007. [1](#)
- [7] V. Belagiannis and A. Zisserman. Recurrent human pose estimation. *arXiv preprint arXiv:1605.02914*, 2016. [1](#), [11](#)
- [8] L. Bo and C. Sminchisescu. Twin gaussian processes for structured prediction. In *International Journal of Computer Vision*, 2010. [10](#), [11](#)
- [9] F. Bogo, A. Kanazawa, C. Lassner, P. Gehler, J. Romero, and M. J. Black. Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. In *European Conference on Computer Vision (ECCV)*, 2016. [1](#), [2](#), [7](#), [10](#), [11](#)
- [10] E. Brau and H. Jiang. 3D Human Pose Estimation via Deep Learning from 2D Annotations. In *International Conference on 3D Vision (3DV)*, 2016. [3](#)
- [11] A. Bulat and G. Tzimiropoulos. Human pose estimation via convolutional part heatmap regression. In *European Conference on Computer Vision (ECCV)*, 2016. [1](#), [11](#)
- [12] J. Carreira, P. Agrawal, K. Fragkiadaki, and J. Malik. Human pose estimation with iterative error feedback. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [1](#), [11](#)
- [13] J. Chai and J. K. Hodgins. Performance animation from low-dimensional control signals. *ACM Transactions on Graphics (TOG)*, 24(3):686–696, 2005. [1](#)
- [14] C.-H. Chen and D. Ramanan. 3d human pose estimation = 2d pose estimation + matching. In *CVPR 2017-IEEE Conference on Computer Vision & Pattern Recognition*, 2017. [1](#), [3](#), [7](#)
- [15] W. Chen, H. Wang, Y. Li, H. Su, Z. Wang, C. Tu, D. Lischinski, D. Cohen-Or, and B. Chen. Synthesizing training images for boosting human 3d pose estimation. In *International Conference on 3D Vision (3DV)*, 2016. [1](#), [3](#), [5](#), [6](#), [8](#)
- [16] X. Chen and A. L. Yuille. Articulated pose estimation by a graphical model with image dependent pairwise relations. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1736–1744, 2014. [1](#)
- [17] A. Elhayek, E. Aguiar, A. Jain, J. Tompson, L. Pishchulin, M. Andriluka, C. Bregler, B. Schiele, and C. Theobalt. MARCOI - ConvNet-based MARKer-less Motion Capture in Outdoor and Indoor Scenes. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2016. [2](#), [6](#)
- [18] A. Elhayek, E. de Aguiar, A. Jain, J. Tompson, L. Pishchulin, M. Andriluka, C. Bregler, B. Schiele, and C. Theobalt. Efficient ConvNet-based marker-less motion capture in general scenes with a low number of cameras. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3810–3818, 2015. [11](#)
- [19] P. F. Felzenszwalb and D. P. Huttenlocher. Pictorial structures for object recognition. *International Journal of Computer Vision (IJCV)*, 61(1):55–79, 2005. [2](#)
- [20] J. Gall, B. Rosenhahn, T. Brox, and H.-P. Seidel. Optimization and filtering for human motion capture. *International Journal of Computer Vision (IJCV)*, 87(1–2):75–92, 2010. [1](#)
- [21] Y. Ganin and V. Lempitsky. Unsupervised domain adaptation by backpropagation. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 1180–1189, 2015. [8](#), [12](#)
- [22] G. Gkioxari, A. Toshev, and N. Jaitly. Chained predictions using convolutional neural networks. In *European Conference on Computer Vision (ECCV)*, 2016. [1](#), [11](#)
- [23] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [3](#)
- [24] P. Hu, D. Ramanan, J. Jia, S. Wu, X. Wang, L. Cai, and J. Tang. Bottom-up and top-down reasoning with hierarchical rectified gaussians. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [1](#), [11](#)
- [25] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele. Deeppercut: A deeper, stronger, and faster multi-person pose estimation model. In *European Conference on Computer Vision (ECCV)*, 2016. [3](#), [6](#), [11](#)
- [26] C. Ionescu, J. Carreira, and C. Sminchisescu. Iterated second-order label sensitive pooling for 3d human pose estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1661–1668, 2014. [2](#), [12](#)
- [27] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 36(7):1325–1339, 2014. [1](#), [5](#), [6](#), [7](#), [8](#)
- [28] A. Jain, J. Tompson, Y. LeCun, and C. Bregler. Modeep: A deep learning framework using motion features for human pose estimation. In *Asian Conference on Computer Vision (ACCV)*, pages 302–315. Springer, 2014. [2](#)
- [29] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM International Conference on Multimedia*, pages 675–678, 2014. [11](#)
- [30] S. Johnson and M. Everingham. Clustered pose and nonlinear appearance models for human pose estimation.

- In *British Machine Vision Conference (BMVC)*, 2010. doi:10.5244/C.24.12. 1, 3, 6, 11, 12
- [31] S. Johnson and M. Everingham. Learning effective human pose estimation from inaccurate annotation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2011. 1, 3, 12
- [32] H. Joo, H. Liu, L. Tan, L. Gui, B. Nabbe, I. Matthews, T. Kanade, S. Nobuhara, and Y. Sheikh. Panoptic studio: A massively multiview system for social motion capture. In *ICCV*, pages 3334–3342, 2015. 1, 5, 6
- [33] A. M. Lehrmann, P. V. Gehler, and S. Nowozin. A Non-parametric Bayesian Network Prior of Human Pose. In *IEEE International Conference on Computer Vision (ICCV)*, 2013. 4
- [34] V. Lepetit and P. Fua. *Monocular model-based 3D tracking of rigid objects*. Now Publishers Inc, 2005. 4
- [35] S. Li and A. B. Chan. 3d human pose estimation from monocular images with deep convolutional neural network. In *Asian Conference on Computer Vision (ACCV)*, pages 332–347, 2014. 1, 2, 4
- [36] S. Li, W. Zhang, and A. B. Chan. Maximum-margin structured learning with deep networks for 3d human pose estimation. In *IEEE International Conference on Computer Vision (ICCV)*, pages 2848–2856, 2015. 1, 2
- [37] I. Lifshitz, E. Fetaya, and S. Ullman. Human pose estimation using deep consensus voting. In *European Conference on Computer Vision (ECCV)*, 2016. 1, 11
- [38] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015. 3
- [39] A. Meka, M. Zollhöfer, C. Richardt, and C. Theobalt. Live intrinsic video. *ACM Trans. Graph. (Proc. SIGGRAPH)*, 35(4):109:1–14, 2016. 5
- [40] F. Moreno-Noguer. 3d human pose estimation from a single image via distance matrix regression. In *CVPR 2017-IEEE Conference on Computer Vision & Pattern Recognition*, 2017. 1, 7
- [41] G. Mori and J. Malik. Recovering 3d human body configurations using shape contexts. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 28(7):1052–1062, 2006. 2
- [42] A. Newell, K. Yang, and J. Deng. Stacked hourglass networks for human pose estimation. In *European Conference on Computer Vision (ECCV)*, 2016. 1, 2, 11
- [43] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010. 3
- [44] G. Pavlakos, X. Zhou, K. G. Derpanis, and K. Daniilidis. Coarse-to-fine volumetric prediction for single-image 3D human pose. In *CVPR 2017-IEEE Conference on Computer Vision & Pattern Recognition*, 2017. 1, 2, 7, 8
- [45] L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. Gehler, and B. Schiele. Deepcut: Joint subset partition and labeling for multi person pose estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2, 11
- [46] L. Pishchulin, A. Jain, M. Andriluka, T. Thormählen, and B. Schiele. Articulated people detection and pose estimation: Reshaping the future. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3178–3185. IEEE, 2012. 5
- [47] G. Pons-Moll, D. J. Fleet, and B. Rosenhahn. Posebits for monocular human pose estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2337–2344, 2014. 2
- [48] V. Ramakrishna, T. Kanade, and Y. Sheikh. Reconstructing 3d human pose from 2d image landmarks. In *European Conference on Computer Vision*, pages 573–586. Springer, 2012. 11
- [49] H. Rhodin, C. Richardt, D. Casas, E. Insafutdinov, M. Shafiei, H.-P. Seidel, B. Schiele, and C. Theobalt. Ego-Cap: Egocentric Marker-less Motion Capture with Two Fish-eye Cameras. *ACM Trans. Graph. (Proc. SIGGRAPH Asia)*, 2016. 5
- [50] H. Rhodin, N. Robertini, D. Casas, C. Richardt, H.-P. Seidel, and C. Theobalt. General automatic human shape and motion capture using volumetric contour cues. In *European Conference on Computer Vision (ECCV)*, pages 509–526. Springer, 2016. 2, 11
- [51] N. Robertini, D. Casas, H. Rhodin, H.-P. Seidel, and C. Theobalt. Model-based Outdoor Performance Capture. In *International Conference on Computer Vision (3DV)*, 2016. 6
- [52] G. Rogez and C. Schmid. Mocap-guided data augmentation for 3d pose estimation in the wild. In *Advances in Neural Information Processing Systems*, pages 3108–3116, 2016. 3, 7, 8
- [53] G. Rogez, P. Weinzaepfel, and C. Schmid. Lcr-net: Localization-classification-regression for human pose. In *CVPR 2017-IEEE Conference on Computer Vision & Pattern Recognition*, 2017. 3, 7
- [54] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 3, 4
- [55] B. Sapp and B. Taskar. Modoc: Multimodal decomposable models for human pose estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013. 1
- [56] A. Sharif Razavian, H. Azizpour, J. Sullivan, and S. Carlsson. CNN features off-the-shelf: an astounding baseline for recognition. In *Conference on Computer Vision and Pattern Recognition Workshops*, pages 806–813, 2014. 3
- [57] J. Shotton, T. Sharp, A. Kipman, A. Fitzgibbon, M. Finocchio, A. Blake, M. Cook, and R. Moore. Real-time human pose recognition in parts from single depth images. *Communications of the ACM*, 56(1):116–124, 2013. 1
- [58] L. Sigal, A. O. Balan, and M. J. Black. Humaneva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion. *International Journal of Computer Vision (IJCV)*, 87(1-2):4–27, 2010. 1, 6, 8, 11
- [59] E. Simo-Serra, A. Quattoni, C. Torras, and F. Moreno-Noguer. A joint model for 2d and 3d pose estimation from a

- single image. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3634–3641, 2013. 1, 2
- [60] E. Simo-Serra, A. Ramisa, G. Alenyà, C. Torras, and F. Moreno-Noguer. Single image 3d human pose estimation from noisy observations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2673–2680. IEEE, 2012. 1, 2
- [61] C. Sminchisescu and B. Triggs. Covariance scaled sampling for monocular 3d body tracking. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages I–447. IEEE, 2001. 1
- [62] J. Starck and A. Hilton. Model-based multiple view reconstruction of people. In *IEEE International Conference on Computer Vision (ICCV)*, pages 915–922, 2003. 1
- [63] C. Stoll, N. Hasler, J. Gall, H.-P. Seidel, and C. Theobalt. Fast articulated motion tracking using a sums of Gaussians body model. In *IEEE International Conference on Computer Vision (ICCV)*, pages 951–958, 2011. 1
- [64] C. J. Taylor. Reconstruction of articulated objects from point correspondences in a single uncalibrated image. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 677–684, 2000. 2
- [65] B. Tekin, I. Katircioglu, M. Salzmann, V. Lepetit, and P. Fua. Structured Prediction of 3D Human Pose with Deep Neural Networks. In *British Machine Vision Conference (BMVC)*, 2016. 1, 2
- [66] B. Tekin, P. Márquez-Neila, M. Salzmann, and P. Fua. Fusing 2D Uncertainty and 3D Cues for Monocular Body Pose Estimation. *arXiv preprint arXiv:1611.05708*, 2016. 2, 3
- [67] B. Tekin, A. Rozantsev, V. Lepetit, and P. Fua. Direct Prediction of 3D Body Poses from Motion Compensated Sequences. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2, 7
- [68] The Captury. <http://www.thecaptury.com/>, 2016. 5
- [69] D. Tome, C. Russell, and L. Agapito. Lifting from the deep: Convolutional 3d pose estimation from a single image. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1
- [70] J. J. Tompson, A. Jain, Y. LeCun, and C. Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1799–1807, 2014. 1, 6
- [71] A. Toshev and C. Szegedy. Deeppose: Human pose estimation via deep neural networks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1653–1660, 2014. 1, 6
- [72] R. Urtasun, D. J. Fleet, and P. Fua. Monocular 3d tracking of the golf swing. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 932–938, 2005. 1
- [73] C. Wang, Y. Wang, Z. Lin, A. L. Yuille, and W. Gao. Robust estimation of 3d human poses from a single image. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2361–2368, 2014. 1, 2
- [74] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh. Convolutional Pose Machines. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2, 11
- [75] C. R. Wren, A. Azarbayejani, T. Darrell, and A. P. Pentland. Pfunder: real-time tracking of the human body. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 19(7):780–785, 1997. 1
- [76] H. Yasin, U. Iqbal, B. Krüger, A. Weber, and J. Gall. A Dual-Source Approach for 3D Pose Estimation from a Single Image. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 3, 10, 11
- [77] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *Advances in Neural Information Processing Systems (NIPS)*, pages 3320–3328, 2014. 3
- [78] Y. Yu, F. Yonghao, Z. Yilin, and W. Mohan. Marker-less 3D Human Motion Capture with Monocular Image Sequence and Height-Maps. In *European Conference on Computer Vision (ECCV)*, 2016. 1, 2, 7
- [79] F. Zhou and F. De la Torre. Spatio-temporal matching for human detection in video. In *European Conference on Computer Vision (ECCV)*, pages 62–77, 2014. 1
- [80] X. Zhou, S. Leonardos, X. Hu, and K. Daniilidis. 3D shape estimation from 2D landmarks: A convex relaxation approach. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4447–4455, 2015. 1, 2, 4
- [81] X. Zhou, X. Sun, W. Zhang, S. Liang, and Y. Wei. Deep kinematic pose regression. In *ECCV Workshop on Geometry Meets Deep Learning*, 2016. 1, 2, 7, 8
- [82] X. Zhou, M. Zhu, S. Leonardos, and K. Daniilidis. Sparse representation for 3d shape estimation: A convex relaxation approach. *arXiv preprint arXiv:1509.04309*, 2015. 2
- [83] X. Zhou, M. Zhu, S. Leonardos, K. Derpanis, and K. Daniilidis. Sparseness Meets Deepness: 3D Human Pose Estimation from Monocular Video. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 1, 2, 7, 10, 11