# Multiview-Consistent Semi-Supervised Learning for 3D Human Pose Estimation

Rahul Mitra[1*], Nitesh B. Gundavarapu[2*], Abhishek Sharma[3], Arjun Jain[3,4]

[1]IIT Bombay [2]University of California, San Diego [3]Axogyan AI, Bangalore [4] IISc Bangalore

## Abstract

*The best performing methods for 3D human pose estimation from monocular images require large amounts of in-the-wild 2D and controlled 3D pose annotated datasets which are costly and require sophisticated systems to acquire. To reduce this annotation dependency, we propose Multiview-Consistent Semi Supervised Learning (MCSS) framework that utilizes similarity in pose information from unannotated, uncalibrated but synchronized multi-view videos of human motions as additional weak supervision signal to guide 3D human pose regression. Our framework applies hard-negative mining based on temporal relations in multi-view videos to arrive at a multi-view consistent pose embedding. When jointly trained with limited 3D pose annotations, our approach improves the baseline by 25% and state-of-the-art by 8.7%, whilst using substantially smaller networks. Lastly, but importantly, we demonstrate the advantages of the learned embedding and establish view-invariant pose retrieval benchmarks on two popular, publicly available multi-view human pose datasets, Human 3.6M and MPI-INF-3DHP, to facilitate future research.*

## 1. Introduction

Over the years, the performance of monocular 3D human pose estimation has improved significantly due to increasingly sophisticated CNN models [54, 33, 52, 44, 43, 29, 47]. For training, these methods depend on the availability of large-scale 3D-pose annotated data, which is costly and challenging to obtain, especially under in-the-wild setting for articulated poses. The two most popular 3D-pose annotated datasets, Human3.6M [14] (3.6M samples) and MPI-INF-3DHP [28] (1.3M samples), are biased towards indoor-like environment with uniform background and illumination. Therefore, 3D-pose models trained on these datasets don't generalize well for real-world scenarios [8, 54].

Limited training data, or costly annotation, poses serious challenges to not only deep-learning based methods, but other machine-learning methods as well. Semi-

supervised approaches [10, 21, 41, 24] have been extensively used in the past to leverage large-scale unlabelled datasets along with small labelled dataset to improve performance. Semi-supervised methods try to exploit the structure/invariances in the data to generate additional learning signals for training. Unlike classical machine-learning models that use fixed feature representation, deep-learning models can learn a suitable feature representation from data as part of training process too. This unique ability calls for semi-supervised approaches to encourage better features representation learning from large-scale unlabelled data for generalization. Intuitively its more appealing to leverage semi-supervised training signals that are more relevant to the final application. Therefore, given the vast diversity of computer-vision tasks, it remains an exciting area of research to innovate novel semi-supervision signals.

To this end, we leverage projective multiview consistency to create a novel metric-learning based semi-supervised framework for 3D human-pose estimation. Multiview consistency has served as a fundamental paradigm in computer vision for more than 40 years and gave rise to some of the most used algorithms such as stereo [39], structure from motion [19], motion capture [31], simultaneous localization and mapping [4], etc. From human-pose estimation perspective, the intrinsic 3D-pose of the human-body remains the same across multiple different views. Therefore, a deep-CNN should *ideally* be able to map 2D-images corresponding to a common 3D-pose, captured from different viewpoints, to nearby points in an embedding space. Intuitively, such a deep-CNN is learning feature representations that are invariant to different views of the human-pose. Therefore, we posit that perhaps it can learn to project 2D images, from different viewpoints, into a *canonical 3D-pose* space in $\mathcal{R}^N$. In Fig. 1b, we show a few embedding distances between different images from the Human3.6M [14] and provide empirical evidence to the aforementioned hypothesis via a novel cross-view pose-retrieval experiment. Unfortunately, embedding-vectors, $x$, from such a space do not translate directly to the 3D coordinates of human-pose. Therefore, we learn another transformation function from embedding to pose space and regress with small 3D-pose supervision while training. Since, the em-
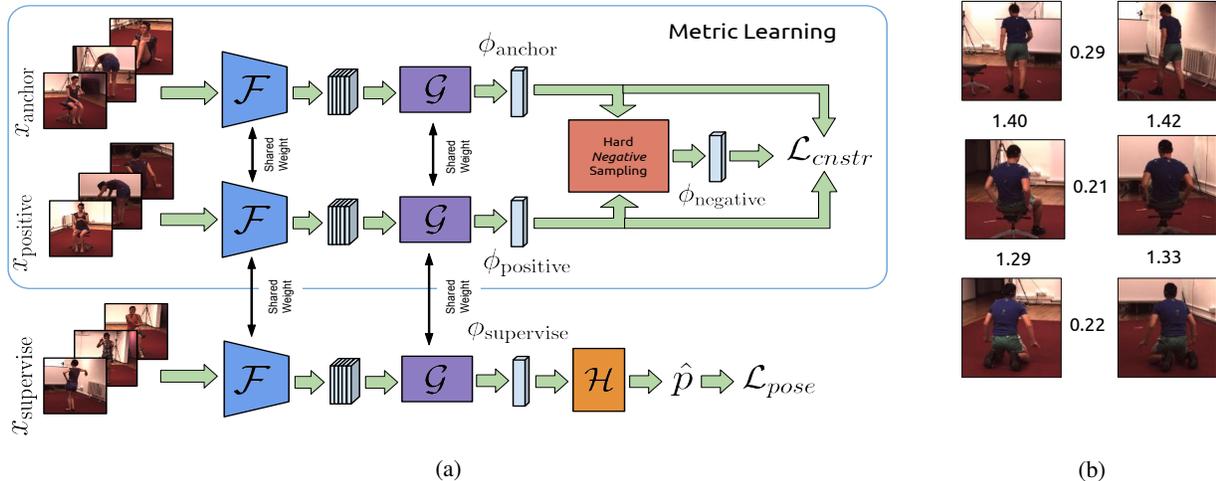
---

[*]- equal contribution

Figure 1: **(a)** Training framework for learning our pose embedding and subsequent *canonical* pose estimation. $x_{\text{anchor}}$ and $x_{\text{positive}}$ are a batch of *anchor* and *positive* image pairs taken from different camera views. $x_{\text{supervise}}$ is the batch of images with 3D-pose supervision. $\mathcal{F}$ is the ResNet based feature extractor. $\mathcal{G}$ maps features extracted from $\mathcal{F}$ to our embedding $\phi$. The *Hard Negative Sampling* module performs in-batch hard mining as given in Eq. 1. Module $\mathcal{H}$ regresses pose $\hat{p}$ from our embedding $\phi$. **(b)** Distances between a few images in our learned embedding space. Each column represents images in the same pose from different view. Images across rows have different poses. The numbers between a pair of images represent its embedding distance. The distance is low for pairs with the same pose irrespective of viewpoint and high for those having different poses.

bedding is shared between the pose-supervision and semi-supervised metric-learning, it leads to better generalizeable features for 3D-pose estimation. We name our proposed framework as *Multiview Consistent Semi-Supervised* learning, or **MCSS** for short.

The proposed framework fits really well with the practical requirements of our problem because it's relatively easier to obtain real-world time-synchronized video streams of humans from multiple viewpoints *vs*. setting up capture rigs for 3D-annotated data out in-the-wild. An alternative approach could be to setup a calibrated multi-camera capture rig in-the-wild and use triangulation from 2D-pose annotated images to obtain 3D-pose. But, it still requires hand-annotated 2D-poses or an automated 2D-pose generation system . In [18], a pre-trained 2D-pose network has been used to generate *pseudo* 3D-pose labels for training a 3D-pose network. Yet another approach exploits relative camera extrinsics for cross-view image generation via a latent embedding [36]. We, on the other hand, don't assume such requirements to yield a more practical solution for the limited data challenge.

We use MCSS to improve 3D-pose estimation performance with limited 3D supervision. In Sec. 5, we show the performance variation as 3D supervision is decreased. Sec. 6 demonstrates the richness of view-invariant MCSS embedding for capturing human-pose structure with the help of a carefully designed cross-view pose-retrieval task on Human3.6M and MPI-INF-3DHP to serve as a bench-

mark for future research in this direction. A summary of our contributions is,

- Proposed a novel Multiview-Consistent Semi-Supervised learning framework for 3D-human-pose estimation.
- Achieved state-of-the-art performance on Human 3.6M dataset with limited 3D supervision.
- Formulated a cross-view pose-retrieval benchmark on Human3.6M and MPI-INF-3DHP datasets.

## 2. Related Work

This section first reviews prior approaches for learning human-pose embedding followed by a discussion of previous weakly supervised methods for monocular 3D human pose estimation to bring out their differences with our approach.

### 2.1. Human Pose Embedding

Historically, human-pose embeddings have been employed in tracking persons [48, 23]. Estimation of 3D human pose and viewpoint from input silhouettes via learning a low dimension manifold is shown in [9]. 2D-pose regression and retrieval by pose similarity embedding is shown in [22, 32], but they require 2D-pose labels. In [42], the need for 2D-pose labels is eliminated by using human motion videos and temporal ordering as weak supervision in a

metric learning framework. Unlike the aforementioned approaches, we learn a 3D-pose embedding by leveraging intrinsic 3D-pose consistency from synchronized multi-view videos. In [45], a 3D-pose embedding is learnt using an over-complete auto-encoder for better structure preservation, but it still requires requires 3D-pose labels for the entire dataset.

## 2.2. Weakly Supervised 3D Human Pose Estimation

Majority of supervised 3D-pose Estimation algorithms [29, 27, 47, 35, 44, 53, 43, 54, 40] require 3D-pose labels in conjunction with either 2D-pose labels or a pre-trained 2D pose estimator to learn a deep-CNN mapping from images to 3D-pose or images to 2D-pose followed by 2D-to-3D lifting task. Some methods refine these pose estimates using either temporal cues, anthropometric constraints, geometric constraints or additional supervision [8, 51, 50, 3, 25, 12]. A complete decoupling between 2D and 3D pose estimation is presented in [50] with the use of generative lifting network followed by a back-projection constraint to achieve generalization. Another line of work focuses on augmenting 2D/3D-pose labels using mesh representation [15, 16, 3, 20] or a dense pose representation [1, 11] to improve pose estimation. All the aforementioned approaches require large amount of annotated 2D and/or 3D labels while our method is designed for limited 3D-pose labels only.

**Strong 2D and limited/no 3D supervision** In recent years, weak-supervision from limited 3D-pose labels along with in-the-wild 2D-pose labels has gained popularity, because labelling 2D-pose is easier than labelling 3D-pose [34, 18, 37, 6, 49, 5]. A weak-supervision in the form of re-projection constraints on the predicted 3D pose is proposed [34]. Mostly, such approaches take advantage of multi-view images during training by means of geometric constraints [37, 18, 5], domain adaptation and adversarial constraints [5], or cross-view reprojection constraints [6]. In [6], a latent 3D-pose embedding is learned by reconstructing 2D-pose from the embedding in a different view. A shallow network with limited 3D-pose supervision is learned to regress 3D-pose from the embedding. A network with pre-trained weights for 2D-pose estimation is used for 3D-pose estimation in [37] followed by multi-view geometric consistency loss. Pseudo 3D-pose labels are generated in [18] for training, while adversarial losses between the 2D skeleton and re-projection of predicted 3D-pose on different views is used for learning in [5]. In [49], starting with 2D pose inputs, a lifting network is trained with siamese loss on the embedding from multiple views to achieve a weak supervision for 3D-pose. Unlike us, [37, 5, 18] require strong 2D-pose estimation systems trained on MPII or COCO datasets while [34, 49, 6] directly work on 2D-pose detections. We, on the other hand, don't need any 2D-pose labels or pre-trained 2D-pose estimation systems.

**limited/no 2D and limited 3D supervision** - To alleviate the need for a large amount of 2D-pose labels, [36] learns an unsupervised geometry aware embedding and estimates 3D-pose from embedding with limited 3D supervision. Novel view synthesis using multi-view synchronized videos of human motions is used to learn a geometry-aware embedding. This method however still requires camera extrinsics and background extraction and performs worse than our approach.

Our approach falls in the same category as we don't use any 2D-pose labels. We utilize synchronized videos from multiple views to learn a pose embedding with limited 3D-pose labels such that similar pose samples are mapped close to each other in the embedding space. Unlike [36], we don't require camera extrinsics and background extraction. Moreover, we exploit multiview-consistency to directly obtain a canonical pose instead of performing image-reconstruction, which affords smaller networks, Resnet-18 [13] vs. Resnet-50.

# 3. Proposed Approach

Our proposed MCSS approach consists of two modules- i) Multiview-consistent metric-learning from time synchronised videos (Sec. 3.1) and ii) 3D-pose regression with limited 3D supervision (Sec. 3.2). Both the modules are jointly trained as shown in Fig. 1a. Metric-learning acts as semi-supervision signal to reduces the dependency on large-scale 3D-pose labels while pose-regression encourages to learn pose-specific features.

## 3.1. Multiview-Consistent Metric Learning

We utilize *Hardnet* framework [30] to learn pose embedding. The datasets used for training is divided into images belonging to one of $\mathcal{S} = \{S_1, S_2, \ldots S_n\}$ set of subjects. $\mathcal{P} \subset \mathbb{R}^{16 \times 3}$ is the set of all possible poses and each pose is viewed from $\mathcal{V} = \{v_1, v_2, \ldots v_q\}$ viewpoints. For training *hardnet*, each batch consists of paired *anchor* ($\mathcal{X}_p^{v_a}(S_i) \in \mathcal{X}$) and *positive* ($\mathcal{X}_p^{v_b}(S_i) \in \mathcal{X}$) images, from subject $S_i$, with the same pose, $p \in \mathcal{P}$, taken from two different viewpoints $v_a$ and $v_b$, here $\mathcal{X} \subset \mathbb{R}^{3 \times 256 \times 256}$ is the set of images.

We pass both the *anchor* and *positive* images through feature extractor ($\mathcal{F}_{\theta_\mathcal{F}} : \mathcal{X} \to \Psi; \Psi \subset \mathbb{R}^{512 \times 4 \times 4}$) to generate features $\{\psi_p^{v_a}, \psi_p^{v_b}\} \in \Psi$. The feature extractor network is parameterised by $\theta_\mathcal{F}$. The features are then finally passed through an embedding generating network ($\mathcal{G}_{\theta_\mathcal{G}} : \Psi \to \Phi; \Phi \subset \mathbb{R}^{dim_\phi}$; where $dim_\phi$ is dimension of our embedding). Let's assume we feed *anchor* and *positive* images to $\mathcal{F}$ in batches of $m$. Once corresponding features $\{\phi_{p_1}^{v_{a_1}}, \ldots, \phi_{p_m}^{v_{a_m}}\}$ and $\{\phi_{p_1}^{v_{b_1}}, \ldots, \phi_{p_m}^{v_{b_m}}\}$ are computed, we create a distance matrix $D$ of size of $m \times m$ with $D(i, j) = \|\phi_{p_i}^{v_{a_i}} - \phi_{p_j}^{v_{b_j}}\|_2$. *Negative*s $\phi_{p_{j_{min}}}^{v_{j_{min}}}$ and $\phi_{p_{k_{min}}}^{v_{k_{min}}}$

for each of $\phi_{p_i}^{v_{a_i}}$ and $\phi_{p_i}^{v_{b_i}}$ are then sampled from the current batch which lie closest in the embedding space from $\phi_{p_i}^{v_{a_i}}$ and the $\phi_{p_i}^{v_{b_i}}$ respectively. Mathematically, the sampling is formulated in Eq. 1. Here, $\beta$ denotes the minimum distance between a hard-mined *negative* and *anchor/positive* in embedding space. The threshold $\beta$ is necessary for stable training and to avoid similar poses as negatives.

$$
\begin{aligned}
j_{min} &= \arg\min_{j \neq i} \delta(D(i,j)) * D(i,j); \\
k_{min} &= \arg\min_{k \neq i} \delta(D(k,i)) * D(k,i) \\
\delta(x) &= 1 \text{ if } x > \beta, \; 0 \text{ otherwise} \\
D_{min}^i &= \min(D(i, j_{min}), D(k_{min}, i))
\end{aligned}
\tag{1}
$$

The average contrastive loss is given in Eq. 2, with $\alpha$ being the margin.

$$
\mathcal{L}_{cnstr} = \frac{1}{m} \sum_{i=1}^{m} D(i,i) + \max(0, \alpha - D_{min}^i) \tag{2}
$$

Note that the above learning framework has the following two objectives, namely, a) to bring the *anchors* and *positives* closer and b) to separate out the *negatives* from *anchors* and *positives*. Intuitively, the goal is to learn embedding that captures 3D-pose information while ignoring irrelevant information, such as subject appearance or background. To this end, we propose the following mini-batch selection mechanism to promote the aforementioned goal:

### 3.1.1 Mini-batch Selection

We compose each mini-batch using *anchor* and *positive* pairs from the same subject, and in many cases with overlapping backgrounds, and the *negatives* are also from the same subject since *Hardnet* chooses the hardest negatives from the same mini-batch. The presented mini-batch selection scheme encourages the resulting embedding to capture pose information while discarding subject-appearance and background features when separating the hardest negatives from *anchors* and *positives*. It's due to the inclusion of same personal-appearance and background in both the negatives and anchor/positives, which cannot be used to separate negatives. We take care to not include temporally close images in a mini-batch by sub-sampling and appropriately choosing $\beta$. Specific hyper-parameter choices are detailed in supplementary material. In Sec. 6, we show pose retrieval ability of the learned embedding to show that it has indeed successfully captured 3D-pose information.

### 3.2. Pose Regression

Most 3D-pose estimation approaches focus on regressing for pose in the local camera coordinate system [54, 33, 47, 27, 44, 50, 37]. In our framework, however, 2D-images

captured from different views are all mapped to nearby embedding locations, if their intrinsic 3D-poses are the same. Therefore, 3D-pose regression using our embedding is ambiguous because the local camera coordinate system is lost. Moreover, the relation from our embedding to the view-specific 3D-pose is *one-to-many*. In order to address this issue, we make use of the MoCap system's global coordinate to represent the 3D-poses instead of view-specific 3D-poses. Hence, synchronous frames captured from different views are labelled as one global-coordinate 3D-pose. However, asynchronous frames can contain poses which are rigid transformations of one another with same 2D projections. In such cases, the mapping from our embedding to 3D-pose is again an ill-posed *one-to-many* mapping. In Fig. 2, an example of such ambiguity is illustrated.
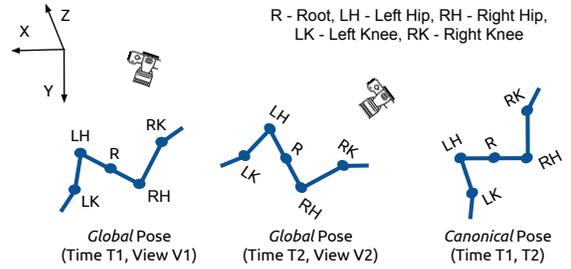


Figure 2: Shows top view of bottom half of a human skeleton taken at two different time instants and view points. The left and middle images show two poses having different joint co-ordinates when presented in the *global* pose while having same projections in their respective cameras. On the contrary *canonical* pose provides provides a uniform representation.

### 3.2.1 Canonical Pose Representation

In order to resolve the aforementioned ambiguities, we formulate a multiview-consistent and rigid rotation invariant 3D-pose representation and refer to it as *canonical* pose. Canonical pose is obtained by constraining the bone connecting the pelvis to the left hip joint to be always parallel to XZ plane. In Human3.6M dataset, the upward direction is +Z axis while XY plane forms the horizontal. Therefore, we rotate the skeleton about the +Z axis until the above mentioned bone is parallel to the XZ plane. We don't require any translation since the joint positions are relative to the pelvis. Mathematically, the transformation from *global* to *canonical* is given in Eq. 3,

$$
\hat{u} = \frac{p_{lh}^{glb} - p_{root}^{glb}}{\|p_{lh}^{glb} - p_{root}^{glb}\|}; \quad \hat{u_{xy}} = \frac{\hat{u_x}\hat{i} + \hat{u_y}\hat{j}}{\|\hat{u_x}\hat{i} + \hat{u_y}\hat{j}\|}
$$
$$
\theta = \cos^{-1}(\hat{u_{xy}} \cdot \hat{i}); \quad p^{can} = R_\theta^z * p^{glb}
\tag{3}
$$

where, $p_{root}^{glb}$ and $p_{lh}^{glb}$ are the root and left-hip joint respectively in the *global* representation. The unit vector along $(p_{lh}^{glb} - p_{root}^{glb})$ is represented as $\hat{u}$, and $\theta$ is the required angle of rotation along the +Z-axis to obtain the canonical pose representation. A positive side-effect of canonical pose representation *vs.* view-specific representation is that our predicted canonical pose doesn't change orientation with variations in camera view. A similar approach to achieve a rotation-invariant pose is suggested in [47]. Note that the canonical pose is constructed directly from MoCap system's coordinates and doesn't require camera extrinsics.

Finally, we regress for canonical pose from the latent embedding $\Phi$ with the help of a shallow network ($\mathcal{H}_{\theta_\mathcal{H}} : \Phi \to \mathcal{P}$), as shown in Fig. 1a. The loss-function is L1-norm between the predicted, $\hat{p}$, and target, $p \in \mathcal{P}$, canonical 3D-pose: $\mathcal{L}_{pose} = \|p - \hat{p}\|_1$.

## 4. Implementation and Training Details

We use the first 4 residual blocks of an ImageNet [38] pre-trained ResNet-18 as our backbone. In addition, we modify the batch-norm (BN) layers by turning off the affine parameters as suggested in [30]. For an input image of size $224 \times 224$ pixels, the output of ResNet is a $512 \times 7 \times 7$ blob, which is further down-sample by 2 using a max-pool operation to get $\Psi$. The embedding network $\mathcal{G}$ is FC layers followed by L2-normalization and it maps $\Psi$ to an embedding of dimension $dim_\phi$ (128 in our case), following usual [30, 46].

For 3D-pose regression, the input data is normalized for each joint. The pose regression network $\mathcal{G}$ consists of FC layers FC(128, 48), with $\Phi \subset \mathbb{R}^{128}$. The margin $\alpha$ for $\mathcal{L}_{conrst}$ is set at 0.6 and $\beta$ at 0.3. Adam [17] optimized is used with default parameters $(0.9, 0.99)$ with initial learning rate $10^{-3}$. The model is trained for 40 epochs with a drop in learning-rate by 0.1 at every 20 epochs. In our joint training frame work, ratio of the batch size for metric learning to pose regression is kept at $3 : 1$ with batch size for regression is 22. A schematic diagram of our network architecture is shown in Fig. 1a.

### 4.1. Datasets

We use the popular Human3.6M [14] and MPI-INF-3DHP [28] datasets for our experiments.

- **Human3.6M [14]** contains 3.6 million frames captured from an indoor MoCap system with 4 cameras ($\mathcal{V}$). It comprises of 11 subjects ($\mathcal{S}$), each performing 16 actions with each action having 2 sub-actions. Following the standard *Protocol 2* [44], we use subjects (S1, S5, S6, S7, S8) for training and (S9, S11) for testing. Like several other methods, we also use cropped subjects' using bounding-boxes provided with the dataset and **temporal sub-sampling** is done

to include every $5^{th}$ and $64^{th}$ frame for training and testing phase, respectively.

- **MPI-INF-3DHP [28]** is generated from a MoCap system with 12 synchronized cameras in both indoor and outdoor settings. It contains 8 subjects($\mathcal{S}$) with diverse clothing. We use the 5 chest height cameras($\mathcal{V}$) for both training and test purposes. Since the test set doesn't contain annotated multi-view data, we use S1-S6 for training and S7-S8 for evaluation.

## 5. Quantitative Evaluation for Pose Estimation

We perform the same quantitative experiment as presented in [36] to assess the benefits of the learned embedding in 3D-pose estimation on Human 3.6M dataset. We evaluate using three well adopted metrics, MPJPE, PA-MPJPE and Normalized MPJPE (N-MPJPE) (introduced in [37]) which incorporates a scale normalization to make the evaluation independent of person's height. We compare our proposed approach and its variants against a baseline which only uses $\mathcal{L}_{pose}$. In addition, we compare our method against the approach proposed by Rhodin et al. [36] and [37], although it estimates human poses in the camera coordinate system. We also report the performance of Rhodin et al. [36] using ResNet-18 as the feature extractor instead of ResNet-50. It is to be noted that [36] uses additional information at training time in the form of relative camera rotation and background extraction which requires sophisticated, well calibrated setup. We acknowledge the existence of more accurate methods like [5, 18, 7] than [36, 37] on Human3.6M when abundant 2D and limited 3D labels are available. For comparison with these approaches, however, we report results from [6] that requires limited 3D supervision but complete 2D supervision from both Human3.6M and MPII [2] dataset. Since, our focus is advancing the research in monocular 3D-pose estimation without using 2D labels under limited 3D-pose labels, we restrict our comparison to cases with limited supervision from both 2D and 3D labels. We don't include the results of [34] as it requires multiple temporally adjacent frames at inference-stage and uses pre-trained 2D-pose estimation models learned from large-scale 2D-pose annotated datasets.

In order to show performance variation as a function of 3D-pose supervision, we report N-MPJPE values for models trained using different amount of 3D-pose labels, in Fig. 3. In this experiment, 3D-pose supervision is reduced gradually using all 5 subjects, to S1+S5, only S1, 50% S1, 10%S1 and finally 5% S1. MVSS clearly outperforms the baseline by a margin of **37.34 N-MPJPE** when only S1 is used for supervision. Moreover, MVSS degrades gracefully as 3D-pose supervision is reduced, which validates the importance of $\mathcal{L}_{conrst}$ in providing weak supervision to capture 3D-pose. Qualitative comparison of our method against the baseline is shown in Fig. 3.
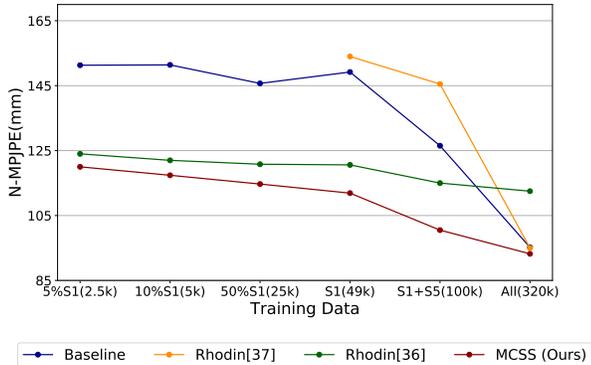
Figure 3: N-MPJPE *vs*. 3D-pose supervision on test split of Human3.6M. Our proposed model outperforms the baseline and the current state-of-the-art Rhodin et al. [36].

In Tab. 1, we compare MPJPE, N-MPJPE and PA-MPJPE values of our approach against baseline and [36]. Clearly, our method outperforms [36] by **22.4** N-MPJPE when fully supervised on 3D data and by **10.7** N-MPJPE with 3D-pose supervision limited to S1. For MPJPE however, the margin is **10.75**. Interestingly as mentioned in [36], the performance of [37] drastically falls when pre-trained model from strong 2D-pose supervision is not used (reported in Tab. 1 as Rhodin [37]* and Rhodin [37]).

As part of ablation studies, we also compare the performance of our learning framework when target pose is represented in MoCap's(*global* pose) against our *canonical* representation in Tab. 2. We observe dramatic decrease in performance, 45 MPJPE, which validates the importance of canonical representation. We also show the results for a deeper ResNet-34 [13] back-end network. We observe a slight drop in performance, 3 MPJPE points, perhaps due to over-fitting.

An additional benefit of our proposed framework is in the use of a much smaller ResNet-18 feature extractor as compared to ResNet-50 used in Rhodin et al [36]. This affords an inference time of 24.8ms *vs*. 75.3ms by [36] on a NVIDIA 1080Ti GPU. Note that Rhodin et al. [36] shows degradation in performance when using the smaller ResNet-18 backbone. We attribute it to direct latent embedding similarity learning instead of generative modelling that requires more representation capacity.

## 6. Analysis of Learned Embedding

In this section, we demonstrate the quality of our learned embedding in capturing 3D human-pose by showing i) pose based cluster formation in our embedding space through retrieval tasks, ii) the correlation between embedding and pose distances. In Fig. 5, we show qualitative image retrieval results based on embedding distance. We can clearly see that the closest images from other subjects and other

| Super-vision | Method | N-MPJPE | MPJPE | PA-MPJPE |
|---|---|---|---|---|
| All | Rhodin [37]* | 63.30 | 66.80 | 51.60 |
| | Chen [6]* | NA | 80.20 | 58.20 |
| | Baseline | 95.07 | 97.90 | 77.18 |
| | Rhodin [37] | 95.40 | NA | NA |
| | Rhodin [36] | 115.00 | NA | NA |
| | MCSS(Ours) | 92.60 | 94.25 | 72.48 |
| S1 | Rhodin [37]* | 78.20 | NA | NA |
| | Chen [6]* | NA | 91.90 | 68.00 |
| | Baseline | 149.28 | 154.78 | 113.69 |
| | Rhodin [37] | NA | 153.30 | 128.60 |
| | Rhodin [36] | 122.60 | 131.70 | 98.20 |
| | Rhodin [36]-Res18 | 136.00 | NA | NA |
| | MCSS(Ours) | **111.94** | **120.95** | **90.76** |

Table 1: Comparing N-MPJPE and MPJPE values between different approaches on Human 3.6M dataset when supervised on all 5 subjects and on only S1. **Note:** Pre-trained ImageNet weights are used to initialize the networks by all the methods. Methods or its variants marked with '*' are supervised with large amount of in-the-wild 2D annotations from MPII [2] dataset either during training or by means of a pre-trained 2D pose estimator. All other methods use much weaker supervision by assuming no 2D annotations and *MCSS* outperforms the state-of-the-art [36] in such settings. NA is assigned against a method if the corresponding result is not reported by the authors.

viewpoints to the query image in embedding space share similar poses. We additionally provide T-SNE [26] plots of our learned embedding space and experiments on generalization to novel view-point in the supplementary material.

### 6.1. Cross-View and Cross-Subject Pose Retrieval

Our learned embedding tries to project similar pose-samples close to each other irrespective of the subject, view-point and background. To validate this claim, we seek mo-

| Supervision | Method | N-MPJPE |
|---|---|---|
| S1 | MCSS | 111.94 |
| | MCSS-global | 157.30 |
| | MCSS-ResNet34 | 115.85 |

Table 2: Comparing N-MPJPE values when pose estimation is done in Mocap's (MCSS-global) and *canonical*(MCSS) representations when only subject S1 is used for supervision. Performance of using ResNet-34 as back-end is reported against MCSS-ResNet34.
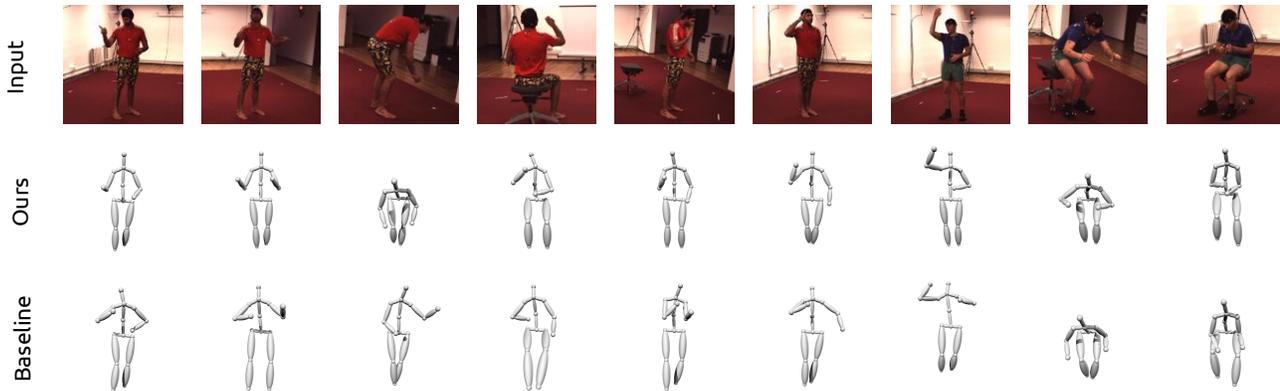
Figure 4: Qualitative results on *canonical* pose estimation by our proposed framework (**MCSS**) against our **Baseline** on Human 3.6M test split (S9, S11). Both the models are trained with supervision from labels of subject S1. Our method produces more accurate estimates for even for challenging poses like 'sitting' and 'bending'.



Figure 5: Qualitative image retrievals on Human 3.6M (S9, S11) and MPI-INF-3DHP (S7, S8) test sets. The first row represents query image and the rows below are the top 3 closest images in embedding space. For the left-most and right-most columns, the retrieval database is composed of images from different subject and viewpoint from that of query's. For the middle two columns, retrieval database is composed of images of same subject but different viewpoint from that of query's. Note how the retrieved poses are very similar to query poses.

tivation from [42], [22] and propose Mean-PA-MPJPE@$K$ to measure the Procrustes Aligned Mean Per Joint Position Error (PA-MPJPE) of $K$ closest neighbours from different views. Since, similar poses in terms of the intrinsic human-body pose can still have different orientations, we use Procrustes Aligned MPJPE to remove this artifact. We compare our model against an *Oracle*, which uses ground truth 3D-pose labels. Given a query image, we ensure that the retrieval database contains images taken from viewpoints other than that of the query image. It is done to clearly bring out the view invariance property of the proposed embedding. First, we report the Mean-PA-MPJPE@$K$ between query pose and its $K$ nearest neighbors in the embedding space. In Fig. 6, we show the comparison of Mean-PA-
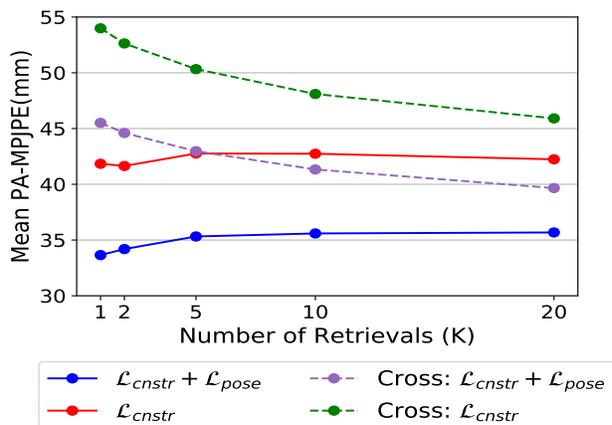


Figure 6: Mean-PA-MPJPE for increasing number of retrievals $K$ on Human3.6M dataset. Prefix 'Cross-' indicates retrieval done on different subjects from that of query. $\mathcal{L}_{pose}$ is from subject S1 in both the cases. All values reported are relative to an Oracle. Low values indicates our retrieved poses are similar to that of the Oracle. PAMPJPE is in mm.

7

MPJPE@$K$ of retrieved poses when retrieval is done from images with:

**Case 1:** all test subjects including that of query's.

**Case 2:** all test subjects except that of query's - *cross*.

We report our results relative to the Oracle. The nearly horizontal plots with low errors suggest that our model picks poses similar to that of the Oracle irrespective of $K$. The error rate is slightly higher for $K = 1, 2$ since our model retrieves images from clusters and does not always pick the one with the lowest error as done by the oracle. The error is lower for **Case 1** than **Case 2** due to the presence of images in the query database that share the exact same pose as that of query, but from different viewpoints. We can also note that upon $\mathcal{L}_{pose}$ from S1, the clustering and mean mpjpe improves in both same subject, **Case 1**, and cross-subject, **Case 2**, settings falling in line with our expectation that small amount of pose supervision improves clustering.

| Method | K=1 | K=5 | K=10 | K=20 |
|--------|-----|-----|------|------|
| $\mathcal{L}_{cnstr}$ | 48.40 | 62.46 | 56.29 | 55.63 |
| Cross-$\mathcal{L}_{cnstr}$ | 82.29 | 83.53 | 80.65 | 76.00 |

Table 3: Mean-PA-MPJPE (mm) for increasing number of retrievals (K) on MPI-INF-3DHP dataset after finetuning with $\mathcal{L}_{cnstr}$. Prefix Cross- indicates retrieval is done on subject other than query's. All values are reported with respect to the Oracle.

## 6.2. Correlation between Embedding and Pose

In this section, we illustrate the variation exhibited by our learned embedding with change in human pose. To this end, we plot mean embedding distance between a query image and stacks of images with increasing pose difference with that of the query in Fig. 7. Both the query and the image stacks belong to the same subject. One can observe a clear positive co-relation between embedding distance and corresponding pose difference. Further, same view and different view show similar correlations with poses justifying the fact that our learned embedding is multi-view consistent.

## 6.3. Generalization & Limitations:

To test cross-dataset generalization, we applied a model trained on Human 3.6M dataset and performed cross-view pose retrievals on MPI-INF-3DHP dataset. We obtained a mean MPJPE of 119.6mm and 101.9mm for $K = 10$ and $K = 20$ respectively. Further fine-tuning with $\mathcal{L}_{conrst}$ using multi-view images from MPI-INF-3DHP improved the performance to 62.46mm and 56.29mm, see Tab. 3. The dip in performance on cross dataset can be attributed to the fact that our feature extractor and embedding generat-
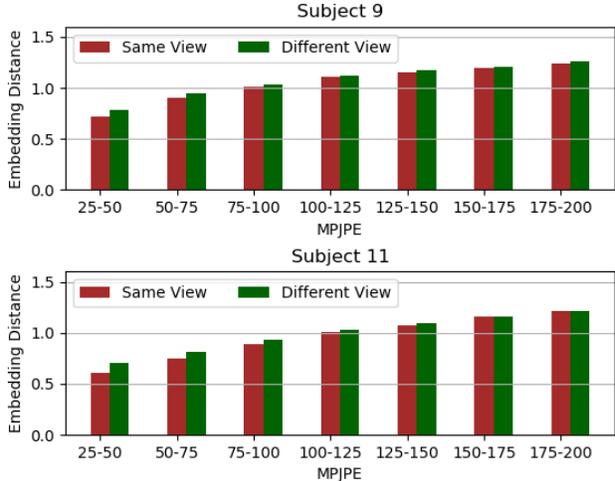


Figure 7: Variation of mean embedding distance with increasing pose variation. We use the show results on (S9, S11) with model being trained with $\mathcal{L}_{cnstr}$ on (S1, S5, S6, S7, S8) and $\mathcal{L}_{pose}$ on (S1). Images are stacked in bins based on the MPJPE difference of their corresponding poses with that of the query. On the Y-axis, the mean embedding distance between the query and the image stacks is plotted. In **Same View**, the query image and image stacks belong to the same viewpoint while in **Different View**, the query stacks belong to different viewpoints. The results are averaged over 200 random queries for each subject.

ing network has learnt a mapping from Human 3.6M images to a pose space and the same mapping is not applicable to the domain of MPI-INF-3DHP images because of huge variation in appearance and more challenging variations of poses. However, upon adding $\mathcal{L}_{conrst}$, as shown in Tab. 3 the weak supervision generalizes to new dataset.

## 7. Conclusion and Future Work

In this paper, we demonstrated a novel Multiview-Consistent Semi-Supervised learning approach to capture 3D human structure for pose estimation and retrieval tasks. With the help of our semi-supervised framework, the need for 3D-pose is reduced. It enables our method to outperform contemporary weakly-supervised approaches even while using a smaller network. Furthermore, we provided strong benchmarks for view-invariant pose retrieval on publicly available datasets.

In future, we plan to use multi-view synchronised videos captured in-the-wild from a larger set of viewpoints to improve generalisation further. We also plan to extend our framework to capture very fine grained pose variations with our embedding by learning distributions of pose variations in temporally consecutive frames using limited 3D annotations.

# References

[1] Rıza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. In *CVPR*, 2018. 3

[2] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *CVPR*, 2014. 5, 6

[3] Anurag Arnab, Carl Doersch, and Andrew Zisserman. Exploiting temporal context for 3d human pose estimation in the wild. In *CVPR*, 2019. 3

[4] Raja Chatila and Jean-Paul Laumond. Position referencing and consistent world modeling for mobile robots. In *ICRA*, 1985. 1

[5] Ching-Hang Chen, Ambrish Tyagi, Amit Agrawal, Dylan Drover, Rohith MV, Stefan Stojanov, and James M Rehg. Unsupervised 3d pose estimation with geometric self-supervision. In *CVPR*, 2019. 3, 5

[6] 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 *CVPR*, 2019. 3, 5, 6

[7] 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 *CVPR*, 2019. 5

[8] Rishabh Dabral, Anurag Mundhada, Uday Kusupati, Safeer Afaque, Abhishek Sharma, and Arjun Jain. Learning 3d human pose from structure and motion. In *ECCV*, 2018. 1, 3

[9] Ahmed Elgammal and Chan-Su Lee. Inferring 3d body pose from silhouettes using activity manifold learning. In *CVPR*, 2004. 2

[10] Rob Fergus, Yair Weiss, and Antonio Torralba. Semi-supervised learning in gigantic image collections. In *NIPS*, 2009. 1

[11] Riza Alp Guler and Iasonas Kokkinos. Holopose: Holistic 3d human reconstruction in-the-wild. In *CVPR*, 2019. 3

[12] 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 *CVPR*, 2019. 3

[13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3, 6

[14] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3. 6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. In *TPAMI*, 2013. 1, 5

[15] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018. 3

[16] Angjoo Kanazawa, Jason Y Zhang, Panna Felsen, and Jitendra Malik. Learning 3d human dynamics from video. In *CVPR*, 2019. 3

[17] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *arXiv*, 2014. 5

[18] Muhammed Kocabas, Salih Karagoz, and Emre Akbas. Self-supervised learning of 3d human pose using multi-view geometry. In *CVPR*, 2019. 2, 3, 5

[19] Jan J Koenderink and Andrea J Van Doorn. Affine structure from motion. In *JOSA A*, 1991. 1

[20] 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 *CVPR*, 2019. 3

[21] Yevhen Kuznietsov, Jorg Stuckler, and Bastian Leibe. Semi-supervised deep learning for monocular depth map prediction. In *CVPR*, 2017. 1

[22] Suha Kwak, Minsu Cho, and Ivan Laptev. Thin-slicing for pose: Learning to understand pose without explicit pose estimation. In *CVPR*, 2016. 2, 7

[23] Chan-Su Lee and Ahmed Elgammal. Modeling view and posture manifolds for tracking. In *CVPR*, 2007. 2

[24] Christian Leistner, Helmut Grabner, and Horst Bischof. Semi-supervised boosting using visual similarity learning. In *CVPR*, 2008. 1

[25] Chen Li and Gim Hee Lee. Generating multiple hypotheses for 3d human pose estimation with mixture density network. In *CVPR*, 2019. 3

[26] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. In *JMLR*, 2008. 6

[27] Julieta Martinez, Rayat Hossain, Javier Romero, and James J Little. A simple yet effective baseline for 3d human pose estimation. In *CVPR*, 2017. 3, 4

[28] 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 *3DV*, 2017. 1, 5

[29] Dushyant Mehta, Srinath Sridhar, Oleksandr Sotnychenko, Helge Rhodin, Mohammad Shafiei, Hans-Peter Seidel, Weipeng Xu, Dan Casas, and Christian Theobalt. Vnect: Real-time 3d human pose estimation with a single rgb camera. In *TOG*, 2017. 1, 3

[30] Anastasiia Mishchuk, Dmytro Mishkin, Filip Radenovic, and Jiri Matas. Working hard to know your neighbor's margins: Local descriptor learning loss. In *NIPS*, 2017. 3, 5

[31] Thomas B Moeslund and Erik Granum. A survey of computer vision-based human motion capture. In *CVIU*, 2001. 1

[32] Greg Mori, Caroline Pantofaru, Nisarg Kothari, Thomas Leung, George Toderici, Alexander Toshev, and Weilong Yang. Pose embeddings: A deep architecture for learning to match human poses. In *arXiv*, 2015. 2

[33] Georgios Pavlakos, Xiaowei Zhou, Konstantinos G Derpanis, and Kostas Daniilidis. Coarse-to-fine volumetric prediction for single-image 3d human pose. In *CVPR*, 2017. 1, 4

[34] Dario Pavllo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3d human pose estimation in video with temporal convolutions and semi-supervised training. In *CVPR*, 2019. 3, 5

[35] Alin-Ionut Popa, Mihai Zanfir, and Cristian Sminchisescu. Deep multitask architecture for integrated 2d and 3d human sensing. In *CVPR*, 2017. 3

[36] Helge Rhodin, Mathieu Salzmann, and Pascal Fua. Unsu-pervised geometry-aware representation for 3d human pose estimation. In *ECCV*, 2018. 2, 3, 5, 6

[37] Helge Rhodin, Jörg Spörri, Isinsu Katircioglu, Victor Con-stantin, Frédéric Meyer, Erich Müller, Mathieu Salzmann, and Pascal Fua. Learning monocular 3d human pose esti-mation from multi-view images. In *CVPR*, 2018. 3, 4, 5, 6

[38] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, San-jeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. In *IJCV*, 2015. 5

[39] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and eval-uation of multi-view stereo reconstruction algorithms. In *CVPR*, 2006. 1

[40] Saurabh Sharma, Pavan Teja Varigonda, Prashast Bindal, Abhishek Sharma, and Arjun Jain. Monocular 3d human pose estimation by generation and ordinal ranking. In *ICCV*, 2019. 3

[41] Rowland R Sillito and Robert B Fisher. Semi-supervised learning for anomalous trajectory detection. In *BMVC*, 2008. 1

[42] Omer Sumer, Tobias Dencker, and Bjorn Ommer. Self-supervised learning of pose embeddings from spatiotemporal relations in videos. In *ICCV*, 2017. 2, 7

[43] Xiao Sun, Jiaxiang Shang, Shuang Liang, and Yichen Wei. Compositional human pose regression. In *ICCV*, 2017. 1, 3

[44] Xiao Sun, Bin Xiao, Fangyin Wei, Shuang Liang, and Yichen Wei. Integral human pose regression. In *ECCV*, 2018. 1, 3, 4, 5

[45] Bugra Tekin, Isinsu Katircioglu, Mathieu Salzmann, Vincent Lepetit, and Pascal Fua. Structured prediction of 3d human pose with deep neural networks. In *BMVC*, 2016. 3

[46] Yurun Tian, Bin Fan, and Fuchao Wu. L2-net: Deep learn-ing of discriminative patch descriptor in euclidean space. In *CVPR*, 2017. 5

[47] Denis Tome, Chris Russell, and Lourdes Agapito. Lifting from the deep: Convolutional 3d pose estimation from a sin-gle image. In *CVPR*, 2017. 1, 3, 4, 5

[48] Raquel Urtasun, David J Fleet, and Pascal Fua. 3d people tracking with gaussian process dynamical models. In *CVPR*, 2006. 2

[49] Márton Véges, Viktor Varga, and András Lőrincz. 3d hu-man pose estimation with siamese equivariant embedding. In *Neurocomputing*, 2019. 3

[50] Bastian Wandt and Bodo Rosenhahn. Repnet: Weakly su-pervised training of an adversarial reprojection network for 3d human pose estimation. In *CVPR*, 2019. 3, 4

[51] Keze Wang, Liang Lin, Chenhan Jiang, Chen Qian, and Pengxu Wei. 3d human pose machines with self-supervised learning. In *TPAMI*, 2019. 3

[52] Min Wang, Xipeng Chen, Wentao Liu, Chen Qian, Liang Lin, and Lizhuang Ma. Drpose3d: Depth ranking in 3d hu-man pose estimation. In *IJCAI*, 2018. 1

[53] Long Zhao, Xi Peng, Yu Tian, Mubbasir Kapadia, and Dim-itris N Metaxas. Semantic graph convolutional networks for 3d human pose regression. In *CVPR*, 2019. 3

[54] Xingyi Zhou, Qixing Huang, Xiao Sun, Xiangyang Xue, and Yichen Wei. Weakly-supervised transfer for 3d human pose estimation in the wild. In *ICCV*, 2017. 1, 3, 4