

A2J: Anchor-to-Joint Regression Network for 3D Articulated Pose Estimation from a Single Depth Image

Fu Xiong^{1*}, Boshen Zhang^{1*}, Yang Xiao^{1†}, Zhiguo Cao¹, Taidong Yu¹, Joey Tianyi Zhou², and Junsong Yuan³

¹ National Key Laboratory of Science and Technology on Multi-Spectral Information Processing, School of Artificial Intelligence and Automation, Huazhong University of Science and Technology

² IHPC, A*STAR, Singapore ³CSE Department, State University of New York at Buffalo

xiongf, zhangbs, Yang_Xiao, zgcao, taidongyu@hust.edu.cn,

joey.tianyi.zhou@gmail.com, jsyuan@buffalo.edu

Abstract

For 3D hand and body pose estimation task in depth image, a novel anchor-based approach termed Anchor-to-Joint regression network (A2J) with the end-to-end learning ability is proposed. Within A2J, anchor points able to capture global-local spatial context information are densely set on depth image as local regressors for the joints. They contribute to predict the positions of the joints in ensemble way to enhance generalization ability. The proposed 3D articulated pose estimation paradigm is different from the state-of-the-art encoder-decoder based FCN, 3D CNN and point-set based manners. To discover informative anchor points towards certain joint, anchor proposal procedure is also proposed for A2J. Meanwhile 2D CNN (i.e., ResNet-50) is used as backbone network to drive A2J, without using time-consuming 3D convolutional or deconvolutional layers. The experiments on 3 hand datasets and 2 body datasets verify A2J's superiority. Meanwhile, A2J is of high running speed around 100 FPS on single NVIDIA 1080Ti GPU.

1. Introduction

With the emergence of low-cost depth camera, 3D hand and body pose estimation from a single depth image draws much attention from computer vision community with wide-range application scenarios (e.g., HCI and AR) [32, 33]. Despite recent remarkable progress [20, 42, 26, 18, 19, 7, 33, 50, 41, 3], it is still a challenging task due to the issues of dramatic pose variation, high similarity among the different joints, self-occlusion, etc [20, 42, 37].

Most of state-of-the-art 3D hand and body pose estimation approaches rely on deep learning technology. Never-

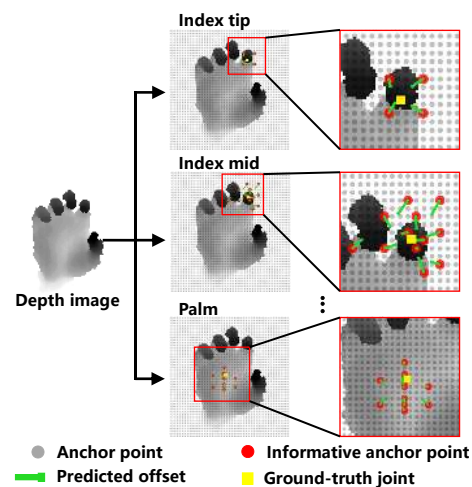


Figure 1. The main idea of anchor-based 3D pose estimation paradigm within A2J. The densely set anchor points predict the positions of joints with weighted aggregation. The different joints possess the different informative anchor points of high weights (i.e., > 0.02), which reveals A2J's adaptive characteristics.

theless, they still suffer from some defects. First, encoder-decoder based FCN manners [2, 43, 27, 4, 42, 41, 26] are generally trained with non-adaptive ground-truth Gaussian heatmap for different joints and with relatively high computational burden. Meanwhile, most of them cannot be fully end-to-end trained towards 3D pose estimation task [35]. Secondly, 3D CNN models [16, 10, 26] are difficult to train with costly voxelizing procedure, due to the large number of convolutional parameters. Additionally, point-set based approaches [14, 17] require some extra time-consuming pre-processing treatments (e.g., point sampling).

Thus, we attempt to address 3D hand and body pose estimation problem using a novel anchor-based approach termed Anchor-to-Joint regression network (A2J). The pro-

*Fu Xiong and Boshen Zhang devote the equal contribution.

†Yang Xiao is the corresponding author (Yang_Xiao@hust.edu.cn).

posed A2J network has end-to-end learning ability. The key idea of A2J is to *predict 3D joint position by aggregating the estimation results of multiple anchor points, in spirit of ensemble learning to enhance generalization ability*. Specifically, the anchor points can be regarded as the local regressors towards the joints from different viewpoints and distances. They are densely set on depth image to capture the global-local spatial context information together. Each of them will contribute to regress the positions of all the joints, but with different weights. The joint is localized by aggregating the outputs of all the anchor points. Since different joints may share the same anchor points, the articulated characteristics among them can be well maintained.

For a specific joint, not all of the anchor points contribute equally. Accordingly, an anchor proposal procedure is proposed to discover the informative anchor points towards the certain joint by weight assignment. During training, both factors of estimation error of anchor points and spatial layout of informative anchor points are concerned. In particular, *the picked up informative anchor points are encouraged to uniformly surround the corresponding joint to alleviate overfitting*. Accordingly, the main idea of the proposed anchor-based 3D pose estimation paradigm within A2J is shown in Fig. 1. We can see that, generally different joints possess different informative anchor points. Furthermore, the visible “index tip” joint holds few informative anchor points. While, the invisible “index mid” joint and the “palm” joint on the relatively flat area possess much more ones, in order to capture richer spatial contexts. This actually reveals A2J’s adaptive property.

Technically, A2J network consists of 3 branches driven by 2D CNN backbone network (i.e., ResNet-50 [21]) without deconvolutional layers. In particular, the 3 branches take charges of predicting in-plane offsets between the anchor points and joints, estimating depth value of the joints, and informative anchor point proposal respectively. The main reasons to build A2J on 2D CNN for 3D pose estimation lie in 3 folders: (1) 3D information is already involved in depth image, using 2D CNN can still reveal 3D characteristics of the original depth image data; (2) compared to 3D CNN and point-set network, 2D CNN can be pre-trained on large-scale datasets (e.g., ImageNet [9]), which may help to enhance its visual pattern capturing capacity for depth image; (3) 2D CNN is of high running efficiency without time-consuming 3D convolution operation and preprocessing procedures (e.g., voxelizing and point sampling).

A2J is experimented on 3 hand datasets (i.e., HANDS 2017 [48], NYU [37], and ICVL [36]) and 2 body pose datasets (i.e., ITOP [20] and K2HPD [42]) to verify its superiority. The experiments reveal that, both for 3D hand and body pose estimation tasks A2J generally outperforms the state-of-the-art methods on effectiveness and efficiency simultaneously. Meanwhile, A2J can online run with the high

speed around 100 FPS on a single NVIDIA 1080Ti GPU.

The main contributions of this paper include:

- A2J: an anchor-based regression network for 3D hand and body estimation from a single depth image. It is of end-to-end learning capacity;
- An informative anchor proposal approach is proposed, concerning the joint position prediction error and anchor spatial layout simultaneously;
- 2D CNN without deconvolutional layers is used to drive A2J to ensure high running efficiency.

A2J’s code is available at <https://github.com/zhangboshen/A2J>.

2. Related Works

The existing 3D hand and body pose estimation approaches can be mainly categorized into non-deep learning and deep learning based groups. The state-of-the-art non-deep learning based ones [33, 22, 13, 46] generally follow the 2-step technical pipeline of first extracting hand-crafted feature, and then executing classification or regression. One main drawback is that, hand-crafted feature is often not representative enough. This tends to lead non-deep learning based method to be inferior to deep learning based manner. Since the proposed A2J falls into deep learning group, next we will introduce and discuss this paradigm from the perspectives of 2D and 3D deep learning respectively.

2D deep learning based approach. Due to end-to-end working manner, deep learning technology holds strong fitting ability for visual pattern characterization. 2D CNN has already achieved great success for 2D pose estimation [38, 4, 27, 43, 44]. Recently it has also been introduced to 3D domain, resorting to global regression [18, 19, 29, 28, 7, 15, 20] or local detection [37, 25, 42, 41, 39] ways. The global regression manner cannot well maintain local spatial context information due to the global feature aggregation operation within fully-connected layers. Local detection based paradigm of promising performance generally chooses to address this problem via encoder-decoder model (e.g., FCN), setting local heatmap for each joint. Nevertheless, heatmap setting is still not adaptive for the different joints. And, the deconvolution operation is time consuming. Furthermore, most of the encoder-decoder based methods cannot be fully end-to-end trained [44].

3D deep learning based approach. To better reveal the 3D property within depth image for performance enhancement, one recent research trend is to resort to 3D deep learning. The paid efforts can be generally categorized into 3D CNN based and point-set based families. 3D CNN based methods [16, 10, 26] voxelizes the depth image into volumetric representation (e.g., occupancy grid models [24]). 3D convolution or deconvolution operation is then executed to capture 3D visual characteristics. However, 3D CNN is relatively hard to tune due to the large number of convo-

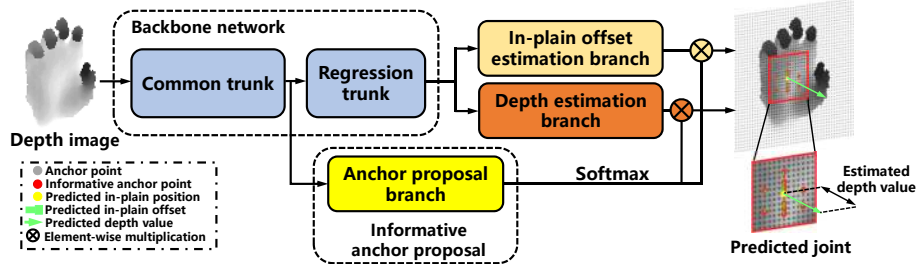


Figure 2. The main technical pipeline of A2J. A2J consists of backbone network and 3 functional branches. The backbone network is built on ResNet-50. And, the 3 branches are in-plain offset prediction branch, depth estimation branch, and anchor proposal branch.

Symbol	Definition
A	Anchor point set.
a	Anchor point $a \in A$.
J	Joint set.
j	Joint $j \in J$.
K	Number of joints.
$S(a)$	In-plain position of anchor point a .
$P_j(a)$	Response of anchor a towards joint j .
$O_j(a)$	Predicted in-plain offset towards joint j from anchor point a .
$D_j(a)$	Predicted depth value of joint j by anchor point a .

Table 1. Symbol definition within A2J.

lutional parameters. Meanwhile, 3D voxelization operation also leads to high computational burden both on memory storage and running time. Another way for 3D deep learning is point-set network [6, 30], transferring depth image into point cloud as input. Nevertheless some time-consuming procedures (e.g., point sampling and KNN search) are required [6, 30], which weakens running efficiency.

Accordingly, A2J belongs to 2D deep learning based group. The dense anchor points capture the global-local spatial context information in ensemble way, without using computationally expensive deconvolutional layers. 2D CNN is used as the backbone network for high running efficiency, also aiming to transfer knowledge from RGB domain.

3. A2J: Anchor-to-Joint Regression Network

The main technical pipeline of A2J is shown in Fig. 2. And, the symbols within A2J are defined in Table 1. A2J consists of 2D backbone network (i.e., ResNet-50), and 3 functional branches: in-plain offset estimation branch, depth estimation branch, and anchor proposal branch. The 3 branches predict $O_j(a)$, $D_j(a)$, and $P_j(a)$ respectively.

Within A2J, anchor points are densely set up on the input depth images with stride $S_t = 4$ pixels to capture the global-local spatial context information as in Fig. 3. Essentially, each of them serves as the local regressor to predict the 3D position of all the joints via in-plain offset prediction branch and depth estimation branch. For certain joint, it is finally localized by aggregating the outputs of all the anchor points. Concerning that maybe not all the anchor points contribute equally to certain joint, the anchor points

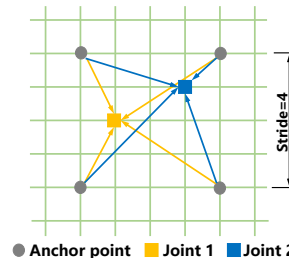


Figure 3. The densely set anchor points on depth image. They will serve for predicting the positions of all joints in ensemble way.

will be assigned weights via anchor proposal branch to discover the informative ones. As consequence, the in-plain position and depth value of joint j can be achieved as the weighted average of the outputs of all anchor points as:

$$\begin{cases} \hat{S}_j = \sum_{a \in A} \tilde{P}_j(a) (S(a) + O_j(a)) \\ \hat{D}_j = \sum_{a \in A} \tilde{P}_j(a) D_j(a) \end{cases}, \quad (1)$$

where \hat{S}_j and \hat{D}_j indicate the estimated in-plain position and depth value of joint j ; $\tilde{P}_j(a)$ can be regarded as the normalized weight of anchor point a towards joint j across all anchor points, and is acquired using softmax by:

$$\tilde{P}_j(a) = \frac{e^{P_j(a)}}{\sum_{a \in A} e^{P_j(a)}}. \quad (2)$$

It is worthy noting that, the anchor point a with $\tilde{P}_j(a) > 0.02$ will be regarded as the informative anchor points for joint j . The selected informative anchor points can reveal A2J's adaptive characteristics as in Fig. 1. Joint position estimation loss and anchor point surrounding loss are used to supervise A2J's end-to-end training. Under their joint supervision, informative anchor points with the spatial layout that surrounds the joint will be picked up to enhance generalization ability. Next, we will illustrate the proposed A2J regression network and its learning procedure in details.

3.1. A2J regression network

Here, the 3 functional branches and backbone network within A2J will be illustrated in details respectively.

3.1.1 In-plain offset and depth estimation branches

Essentially, these 2 branches play the role of predicting the 3D positions of joints. Since in-plain position estimation and depth estimation are of different properties, we choose to execute them separately. Specifically, one is to estimate $O_j(a)$ between anchor points and joints. And, the other is to estimate $D_j(a)$ towards joints. As in Fig. 4, they are built upon the output feature map of regression trunk within backbone network to involve semantic feature. Four 3×3 intermediate convolutional layers (with BN and ReLU) are consequently set to aggregate richer local context information without reducing in-plain size. Since the feature map is a $16 \times$ downsampling of the input depth image on in-plain size (illustrated in Sec. 3.1.3) and anchor point setting stride $S_t = 4$ as in Fig. 3, one feature map point corresponds to $4 \times 4 = 16$ anchor points on depth image. An output convolutional layer with the feature map in-plain size is then set towards all the 16 corresponding anchor points in column-wise manner. Suppose K joints exist, in-plain offset estimation branch is of $16 \times K \times 2$ output channels. And, depth estimation branch is of $16 \times K \times 1$ output channels.

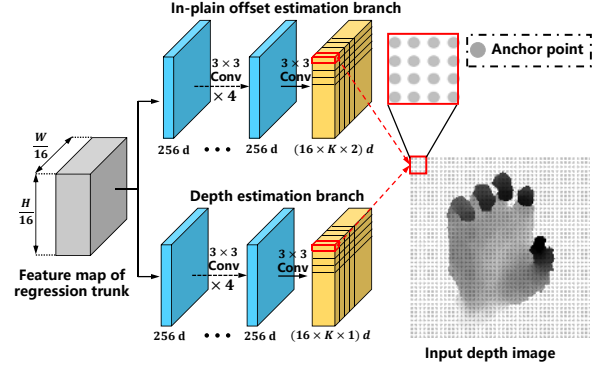


Figure 4. In-plain offset and depth estimation branches. They both contain 4 intermediate convolutional layers with 256 channels, and 1 output convolutional layer (with $16 \times K \times 2$ or $16 \times K \times 1$ channels). W and H indicate width and height of the input depth image. d means dimensionality.

3.1.2 Anchor proposal branch

This branch discovers informative anchor points for the certain joint by weight assignment as Eqn. 2. As in Fig. 5, anchor proposal branch is built upon the output feature map of common trunk within backbone network to involve relatively fine feature. As the 2 branches introduced in Sec. 3.1.1, 4 intermediate convolutional layers and 1 output convolutional layer are consequently set for predicting $P_j(a)$ for the anchor points without losing in-plain size. Accordingly, the output layer of this branch is of $16 \times K \times 1$ channels.

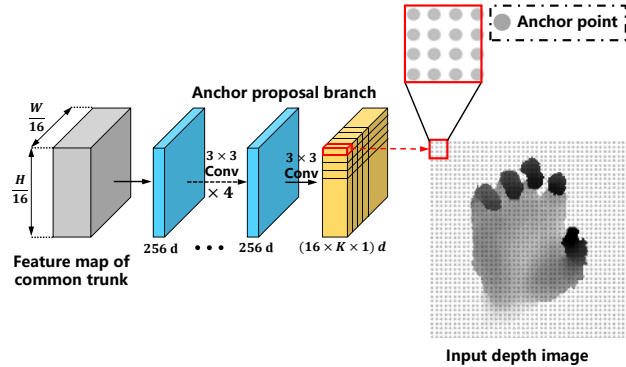


Figure 5. Anchor proposal branch with 4 intermediate convolutional layers with 256 channels, and 1 output convolutional layer with $16 \times K \times 1$ channels. W and H indicate width and height of the input depth image. d means dimensionality.

3.1.3 Backbone network architecture

ResNet-50 [21] pre-trained on ImageNet is used as the backbone network. In particular, layers 0-3 correspond to the common trunk in Fig. 2. And, layer 4 corresponds to regression trunk. Some modifications are executed to make ResNet-50 more suitable for pose estimation. First, the convolutional stride in layer 4 is set to 1. Consequently, the output feature map of layer 4 is a $16 \times$ downsampling of the input depth image on in-plain size. Compared with the raw ResNet-50 with $32 \times$ downsampling, more fine spatial information can be maintained in this way. Meanwhile, the convolution operation within layer 4 is revised as the dilated convolution with a dilation of 2 to enlarge receptive field.

3.2. Learning procedure of A2J

To generate input of A2J, we follow [26] and use center points to crop the hand region from depth image. For body pose, we follow [11] and use bounding box to crop

the body region. For joint j , in-plain target T_j^i denotes the 2D ground-truth in pixel coordinate transformed according to the cropped region. To make T_j^i and depth target T_j^d be in comparable magnitude, we transform the ground-truth depth Z_j of joint j as:

$$T_j^d = \mu(Z_j - \theta), \quad (3)$$

where μ and θ are the transformation parameters. For hand pose μ is set to 1, and θ is set to the depth of center points. For body pose μ is set to 50 and θ is set as 0, since we do not have depth center. During test, the prediction result will be warped back to world coordinate. A2J is then trained under the joint supervision of 2 loss functions: joint position estimation loss and informative anchor point surrounding loss. Next, we will illustrate these 2 loss functions in details.

3.2.1 Joint position estimation loss

Within A2J, the anchor points serve as the local regressors to predict the 3D position of joints in ensemble way. This

objective loss can be formulated as:

$$\begin{aligned} loss_1 = & \alpha \sum_{j \in J} L_{\tau_1} \left(\sum_{a \in A} \tilde{P}_j(a) (S(a) + O_j(a)) - T_j^i \right) \\ & + \sum_{j \in J} L_{\tau_2} \left(\sum_{a \in A} \tilde{P}_j(a) D_j(a) - T_j^d \right), \end{aligned} \quad (4)$$

where $\alpha = 0.5$ is the factor to balance in-plane offset and depth estimation task; T_j^i and T_j^d are the in-plane and depth targets position of joint j ; and $L_{\tau}(\cdot)$ is the smooth L_1 like loss function [31] given by:

$$L_{\tau}(x) = \begin{cases} \frac{1}{2\tau} x^2, & \text{for } |x| < \tau, \\ |x| - \frac{\tau}{2}, & \text{otherwise.} \end{cases} \quad (5)$$

In Eqn. 4, τ_1 is set to 1 and τ_2 is set to 3 since the depth value is relatively noisy.

3.2.2 Informative anchor point surrounding loss

To enhance the generalization ability of A2J, we intend to let the picked up informative anchor points locate around the joints, in spirit of observing the joints from multiple viewpoints simultaneously. Hence, the informative anchor point surrounding loss is defined by us as:

$$loss_2 = \sum_{j \in J} L_{\tau_1} \left(\sum_{a \in A} \tilde{P}_j(a) S(a) - T_j^i \right). \quad (6)$$

To reveal its effectiveness, we show the informative anchor point spatial layouts with and without using it both for hand and body pose cases in Fig. 6. It can be seen that, informative anchor point surrounding loss can essentially help to alleviate viewpoint bias. Its quantitative effectiveness will also be verified in Sec. 4.3.1.

3.2.3 End-to-end training

The 2 loss functions above jointly supervise the end-to-end learning procedure of A2J, which is formulated as:

$$loss = \lambda loss_1 + loss_2, \quad (7)$$

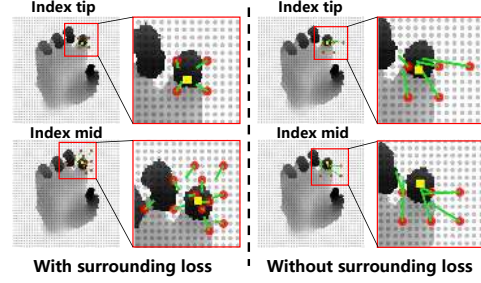
where $loss$ is the loss in all; and $\lambda = 3$ is the weight factor to balance $loss_1$ and $loss_2$.

4. Experiments

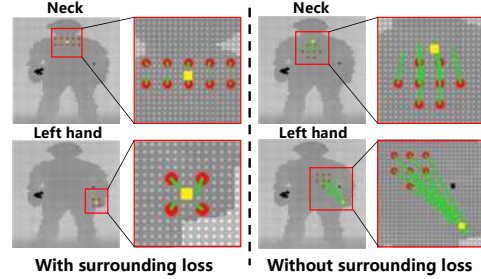
4.1. Experimental setting

4.1.1 Datasets

HANDS 2017 dataset [48]. It contains 957K training and 295K test depth images sampled from BigHand 2.2M [49] and First-Person Hand Action [48] datasets. The ground-truth is the 3D coordinates of 21 hand joints.



(a) Hand cases from NYU dataset



(b) Body cases from ITOP front-view dataset

Figure 6. Effectiveness of anchor point surrounding loss. Grey dot denotes anchor point. Red dot indicate informative anchor point. Green arrow represent in-plane offset. Yellow square corresponds to ground-truth joint.

NYU Hand Pose Dataset [37]. It contains 72K training and 8.2K test depth images with 3D annotation on 36 hand joints. Following [7, 19, 18, 26], we pick 14 of the 36 joints from frontal view for evaluation.

ICVL Hand Pose Dataset [36]. It contains 22K training and 1.5k test depth images. It is augmented to 330K samples by in-plane rotations. 16 hand joints are annotated.

ITOP Body Pose Dataset [20]. It contains 40K training and 10K test depth images both for the front-view and top-view tracks. Each depth image is labelled with 15 3D joint locations of human body.

K2HPD Body Pose Dataset [42]. It contains about 100K depth images. 19 human body joints are annotated with the in-plane manner.

4.1.2 Evaluation metric

For hand, the **average 3D distance error** and **percentage of success frame metrics** [26] are used as the evaluation criteria. For body, **Percent of Detected Joints (PDJ)** [42, 20] and **mean average precision (mAP)** with 10-cm rule [42, 20] are used for evaluation.

4.1.3 Implementation details

A2J network is implemented using PyTorch. The input depth image is cropped and resized to a fixed resolution (i.e., 176×176 for hand, and 288×288 for body). Random

Methods	AVG	SEEN	UNSEEN	FPS
Vanora [47]	11.91	9.55	13.89	-
THU VCLab [8]	11.70	9.15	13.83	-
Oasis [14]	11.30	8.86	13.33	48
RCN-3D [47]	9.97	7.55	12.00	-
V2V* [26]	9.95	6.97	12.43	3.5
A2J (Ours)	8.57	6.92	9.95	105.06

Table 2. Performance comparison on HANDS 2017 dataset [48]. “SEEN” and “UNSEEN” denote the cases whether the test subjects are involved in training set. “AVG” indicates the result over all subjects. And, “*” means the ensemble of 10 models.

Methods	Mean error (mm)	FPS
DISCO [1]	20.7	-
Hand3D [10]	17.6	30
DeepModel [51]	17.04	-
JTSC [12]	16.8	-
Global-to-Local [23]	15.60	50
Lie-X [45]	14.51	-
REN-4x6x6 [18]	13.39	-
REN-9x6x6 [18]	12.69	-
DeepPrior++ [28]	12.24	30
Pose-REN [7]	11.81	-
HandPointNet [14]	10.5	48
DenseReg [39]	10.2	27.8
V2V [26]	9.22	35
P2P [17]	9.045	41.8
A2J (Ours)	8.61	105.06

Table 3. Performance comparison on NYU dataset [37]. “Mean error” indicates the average 3D distance error.

in-plane rotation and random scaling for both in-plane and depth dimension are executed for data augment. Random Gaussian noise is also randomly added with the probability of 0.5 for data augment. We use Adam as the optimizer. The learning rate is set to 0.00035 with a weight decay of 0.0001 in all cases. A2J is trained on NYU for 34 epochs with a learning rate decay by 0.1 every 10 epoch, and for 17 epochs on ICVL and HANDS 2017 with a learning rate decay by 0.1 every 7 epoch. For 2 human body datasets, the epoch for training is set as 26 with a learning rate decay by 0.1 every 10 epoch.

4.2. Comparison with state-of-the-art methods

HANDS 2017 dataset: A2J is compared with the state-of-the-art 3D hand pose estimation methods [47, 8, 14, 26], particularly. The performance comparison is listed in Table 2. It can be observed that:

- On this challenging million-scale dataset, A2J consistently outperforms the other approaches both from the perspectives of effectiveness and efficiency. This essentially verifies the superiority of our proposition;
- It is worthy noting that, A2J is significantly superior to the others with the remarkable margin (2.05 at least) towards the “UNSEEN” test case. This phenomenon essentially demonstrates the generalization ability of A2J;
- V2V* is the strongest competitor of A2J, but with 10

Methods	Mean error (mm)	FPS
LRF [36]	12.58	-
DeepModel [51]	11.56	-
Hand3D [10]	10.9	30
CrossingNets [40]	10.2	90.9
Cascade [34]	9.9	-
JTSC [12]	9.16	-
DeepPrior++ [28]	8.1	30
REN-4x6x6 [18]	7.63	-
REN-9x6x6 [18]	7.31	-
DenseReg [39]	7.3	27.8
Pose-REN [7]	6.79	-
HandPointNet [14]	6.935	48
P2P [17]	6.328	41.8
V2V* [26]	6.286	3.5
A2J (Ours)	6.461	105.06

Table 4. Performance comparison on ICVL dataset [36]. “Mean error” indicates the average 3D distance error.

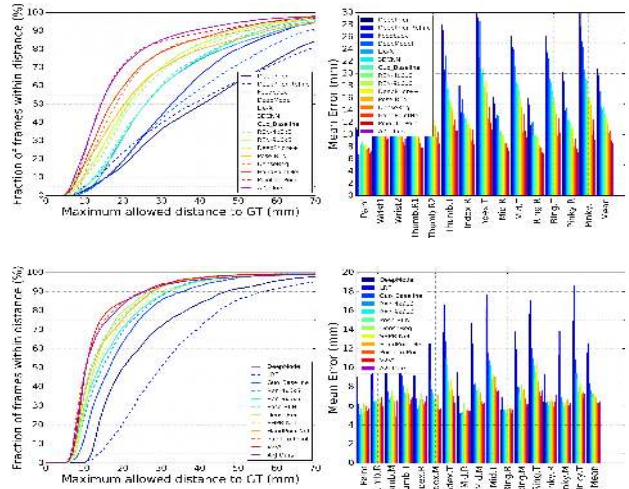


Figure 7. Comparison of A2J with state-of-the-art methods. Left: the percentage of success frames over different error thresholds. Right: 3D distance errors per hand keypoints. Top: NYU dataset. Bottom: ICVL dataset.

models ensemble. As a consequence, it is much slower than A2J with only a single model.

NYU and ICVL datasets: We compare A2J with the state-of-the-art 3D hand pose estimation methods [36, 1, 34, 28, 51, 12, 10, 45, 19, 18, 40, 7, 23, 39, 14, 17, 26] on this 2 datasets specifically. The experimental results are given in Table 3, 4 on the average 3D distance error. Meanwhile, the percentage of success frames over different error thresholds and the error of each joint are also given in Fig. 7. We can summarize that:

- A2J is superior to the other methods in most cases both on accuracy and efficiency. The exceptional case is that, A2J is slightly inferior to V2V* and P2P on ICVL dataset on accuracy but with much higher running efficiency;
- Concerning the good tradeoff between effectiveness and efficiency, A2J essentially takes advantage over the state-of-the-art 3D hand pose estimation approaches.

Method	mAP (front-view)							mAP (top-view)							
	RF [33]	RTW [50]	IEF [5]	VI [20]	CMB [41]	REN- 9x6x6 [18]	V2V* [26]	A2J (Ours)	RF [33]	RTW [50]	IEF [5]	VI [20]	REN- 9x6x6 [18]	V2V* [26]	A2J (Ours)
Head	63.8	97.8	96.2	98.1	97.7	98.7	98.29	98.54	95.4	98.4	83.8	98.1	98.2	98.4	98.38
Neck	86.4	95.8	85.2	97.5	98.5	99.4	99.07	99.20	98.5	82.2	50.0	97.6	98.9	98.91	98.91
Shoulders	83.3	94.1	77.2	96.5	75.9	96.1	97.18	96.23	89.0	91.8	67.3	96.1	96.6	96.87	96.26
Elbows	73.2	77.9	45.4	73.3	62.7	74.7	80.42	78.92	57.4	80.1	40.2	86.2	74.4	79.16	75.88
Hands	51.3	70.5	30.9	68.7	84.4	55.2	67.26	68.35	49.1	76.9	39.0	85.5	50.7	62.44	59.35
Torso	65.0	93.8	84.7	85.6	96.0	98.7	98.73	98.52	80.5	68.2	30.5	72.9	98.1	97.78	97.82
Hips	50.8	90.3	83.5	72.0	87.9	91.8	93.23	90.85	20.0	55.7	38.9	61.2	85.5	86.91	86.88
Knees	65.7	68.8	81.8	69.0	84.4	89.0	91.80	90.75	2.6	53.9	54.0	51.6	70.0	83.28	79.66
Feet	61.3	68.4	80.9	60.8	83.8	81.1	87.60	86.91	0.0	28.7	62.4	51.5	41.6	69.62	58.34
mean	65.8	80.5	71.0	77.4	83.3	84.9	88.74	88.0	47.4	68.2	51.2	75.5	75.5	83.44	80.5

Table 5. Performance comparison on ITOP 3D body pose estimation dataset [20].

Method	PHR [2]	CPM [43]	SH [27]	IEML [42]	CMB [41]	A2J (Ours)
PDJ (0.05)	26.8	30.0	41.0	43.2	52.5	76.3
PDJ (0.10)	70.3	58.5	73.7	64.1	84.2	94.4
PDJ (0.15)	84.7	87.8	84.6	88.1	91.7	97.6
PDJ (0.20)	91.3	93.6	89.0	91.0	95.1	98.6
Average	68.3	67.5	72.1	71.6	80.9	91.7

Table 6. Performance comparison on K2HPD dataset [42].

ITOP dataset: We also compare A2J with the state-of-the-art 3D body pose estimation manners [33, 50, 5, 20, 18, 41, 26] on this dataset. The performance comparison is listed in Table 5. We can see that:

- A2J is significantly superior to the other ones both for front-view and top-view tracks, except V2V*. The performance gap is 3.1 at least for front-view case, and 5 at least for top-view case. This reveals that A2J is also applicable to 3D body pose estimation, as well as 3D hand task;

- A2J is inferior to V2V*. However, V2V* actually consists of 10 models ensemble. Thus, compared with A2J with single model it is of much lower running efficiency.

K2HPD dataset: Since this body pose dataset only provides the pixel-level in-plain ground-truth, the depth estimation branch within A2J is removed accordingly. We also compare A2J with the state-of-the-art approaches [2, 43, 27, 42, 41]. The performance comparison is given in Table 6. It can be observed that:

- A2J outperforms the other methods by large margins consistently, corresponding to the difference PDJ thresholds. In average, the performance gap is 10.8 at least. This demonstrates that, A2J is also applicable to 2D case;

- It is worthy noting that, with the decrease of PDJ threshold the advantage of A2J will be enlarged remarkably. This reveals the fact that, A2J is essentially superior to more accurate body pose estimation.

4.3. Ablation study

4.3.1 Component effectiveness analysis

The component effectiveness analysis within A2J is executed on NYU [37] (hand), and ITOP [20] dataset (body). We

Dataset	Component	error / mAP
NYU (hand)	w/o anchor proposal branch	10.08
	w/o informative anchor point surrounding loss	9.00
	Estimate IPO and DV using one branch	8.95
	A2J (Ours)	8.61
ITOP front-view (body pose)	w/o anchor proposal branch	80.1
	w/o informative anchor point surrounding loss	86.4
	Estimate IPO and DV using one branch	87.4
	A2J (Ours)	88.0

Table 7. Component effectiveness analysis within A2J. ‘‘IPO’’ indicates in-plain offset, and ‘‘DV’’ denotes depth value.

will investigate the effectiveness of anchor proposal branch, informative anchor point surrounding loss, and configuration of in-plain offset and depth estimation branches. The results are listed in Table 7. It can be observed that:

- Without using anchor proposal branch, performance will drop remarkably especially for body pose. This verifies our point that, not all the anchor points contribute equally to the certain joints. Actually, anchor point adaptivity is A2J’s essential property to leverage performance;

- Without using informative anchor point surrounding loss, performance will drop especially for body pose. This demonstrates that, informative anchor point spatial layout is an essential issue that should be concerned towards generalization ability;

- When estimating in-plain offset and depth value in one branch, performance will drop to some degree. This may be caused by the fact that, in-plain offset and depth value holds different physical characteristics.

4.3.2 Effectiveness of anchor-based paradigm

To verify the effectiveness of anchor-based 3D pose estimation paradigm, we compare A2J with the global regression based manner [38] and FCN-based approach [44]. Since FCN model is generally used to predict in-plain joint position, this ablation study is executed on K2HPD

Paradigm	Global regression [38]	FCN model [44]	A2J (Ours)
PDJ (0.05)	61.5	70.4	76.3
FPS	145.12	45.48	93.78

Table 8. Performance comparison among the different paradigms on K2HPD dataset [42].

	Pre-train	From scratch	ImageNet pre-training
	NYU (error)	10.08	8.61
	ITOP front-view (mAP)	87.3	88.0

Table 9. Effectiveness of pre-training A2J on ImageNet.

	Backbone	ResNet-18	ResNet-34	ResNet-50
NYU	error	9.32	9.01	8.61
	FPS	192.25	144.63	105.06
ITOP front-view	mAP	87.1	87.8	88.0
	FPS	167.19	122.47	93.78

Table 10. Performance comparison among the backbones.

dataset [42] only with in-plain ground-truth annotation. Global regression manner encodes depth image with 2D CNN, and then regresses in-plain human joint position using fully-connected layers. FCN model is built following [44]. ResNet-50 [21] is employed as the backbone network for them, which is the same as A2J for fair comparison. PDJ (0.05) is used as the evaluation criteria. The performance comparison is listed in Table 8. We can see that:

- Our proposed anchor-based paradigm significantly outperforms the other 2 ones, when using the same ResNet-50 backbone network. We think 2 main reasons lie. First, compared with global regression based manner local spatial context information can be better maintained within A2J. Meanwhile, compared with FCN model A2J possess anchor point adaptivity towards the certain joint;

- A2J runs faster than FCN model, but slower than global regression way. However, its performance advantage over global regression paradigm is significant, actually with better tradeoff between effectiveness and efficiency.

4.3.3 Effectiveness of the pre-training

One reason for why we build A2J on 2D CNN is that, it can be pre-trained on the large-scale RGB visual datasets (e.g., ImageNet) for knowledge transfer. To verify this point, we compare the performance of A2J with and without pre-training on ImageNet on NYU (hand) and ITOP (body) datasets. The performance comparison is listed in Table 9. It can be observed that, both for hand and body pose cases pre-training A2J on ImageNet can indeed help to leverage the performance.

4.3.4 Backbone network comparison

The comparison among the different backbone networks is further studied. As shown in Table 10, we compare the performance of 3 backbone networks (i.e., ResNet-18, ResNet-34 and ResNet-50). It can be summarized that:

- Deeper network can achieve better results, but with relatively slower running efficiency. However, the performance gap among the different backbones is not huge;

- It is worthy noting that, even using ResNet-18 A2J still can generally achieve the state-of-the-art performance and with extremely fast running speed of 192.25 FPS. This reveals the applicability of A2J towards high real-time running demanding application scenarios.

4.4. Running speed analysis

The average online running speed of A2J for 3D hand pose estimation is 105.06 FPS, including 1.5 ms for reading and warpping image, and 8.0 ms for network forward propagation and post-processing on a single NVIDIA 1080Ti GPU. The running speed for 3D body pose estimation is 93.78 FPS, including 0.4 ms for reading and warpping image, and 10.2 ms for network forward propagation and post-processing. This reveals A2J’s real-time running capacity.

5. Conclusions

In this paper, an anchor-based 3D articulated pose estimation approach for single depth image termed A2J is proposed. Within A2J anchor points are densely set up on depth image to capture the global-local spatial context information, and predict joint’s position in ensemble way. Meanwhile, informative anchor points are extracted to reveal A2J’s adaptive characteristics towards the different joints. A2J is built on 2D CNN without using computational expensive deconvolutional layers. The wide-range experiments demonstrate A2J’s superiority both from the perspectives of effectiveness and efficiency. In future work, we will seek the more effective way to fuse the anchor points.

Acknowledgment

This work is jointly supported by the National Key R&D Program of China (No. 2018YFB1004600), National Natural Science Foundation of China (Grant No. 61876211 and 61602193), the Fundamental Research Funds for the Central Universities (Grant No. 2019kfyXKJC024), the International Science & Technology Cooperation Program of Hubei Province, China (Grant No. 2017AHB051), the start-up funds from University at Buffalo. Joey Tianyi Zhou is supported by Singapore Government’s Research, Innovation and Enterprise 2020 Plan (Advanced Manufacturing and Engineering domain) under Grant A1687b0033 and Grant A18A1b0045. We also thank the anonymous reviewers for their suggestions to enhance the quality of this paper.

References

- [1] Diane Bouchacourt, Pawan K Mudigonda, and Sebastian Nowozin. Disco nets: Dissimilarity coefficients networks. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, pages 352–360, 2016. **6**
- [2] Adrian Bulat and Georgios Tzimiropoulos. Human pose estimation via convolutional part heatmap regression. In *Proc. European Conference on Computer Vision (ECCV)*, pages 717–732, 2016. **1, 7**
- [3] Yujun Cai, Lihao Ge, Jun Liu, Jianfei Cai, Tat-Jen Cham, Junsong Yuan, and Nadia Magnenat Thalmann. Exploiting spatial-temporal relationships for 3d pose estimation via graph convolutional networks. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2019. **1**
- [4] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7291–7299, 2017. **1, 2**
- [5] Joao Carreira, Pulkit Agrawal, Katerina Fragkiadaki, and Jitendra Malik. Human pose estimation with iterative error feedback. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4733–4742, 2016. **7**
- [6] R Qi Charles, Hao Su, Mo Kaichun, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 77–85, 2017. **3**
- [7] Xinghao Chen, Guijin Wang, Hengkai Guo, and Cairong Zhang. Pose guided structured region ensemble network for cascaded hand pose estimation. *arXiv preprint arXiv:1708.03416*, 2017. **1, 2, 5, 6**
- [8] Xinghao Chen, Guijin Wang, Hengkai Guo, and Cairong Zhang. Pose guided structured region ensemble network for cascaded hand pose estimation. *arXiv preprint arXiv:1708.03416*, 2017. **6**
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009. **2**
- [10] Xiaoming Deng, Shuo Yang, Yinda Zhang, Ping Tan, Liang Chang, and Hongan Wang. Hand3d: Hand pose estimation using 3d neural network. *arXiv preprint arXiv:1704.02224*, 2017. **1, 2, 6**
- [11] Hao-Shu Fang, Shuqin Xie, Yu-Wing Tai, and Cewu Lu. Rmpe: Regional multi-person pose estimation. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, pages 2334–2343, 2017. **4**
- [12] Damien Fourure, Rémi Emonet, Elisa Fromont, Damien Muselet, Natalia Neverova, Alain Trémeau, and Christian Wolf. Multi-task, multi-domain learning: application to semantic segmentation and pose regression. *Neurocomputing*, 251:68–80, 2017. **6**
- [13] Varun Ganapathi, Christian Plagemann, Daphne Koller, and Sebastian Thrun. Real-time human pose tracking from range data. In *Proc. European Conference on Computer Vision (ECCV)*, pages 738–751, 2012. **2**
- [14] Lihao Ge, Yujun Cai, Junwu Weng, and Junsong Yuan. Hand pointnet: 3d hand pose estimation using point sets. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8417–8426, 2018. **1, 6**
- [15] Lihao Ge, Hui Liang, Junsong Yuan, and Daniel Thalmann. Robust 3d hand pose estimation in single depth images: from single-view cnn to multi-view cnns. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3593–3601, 2016. **2**
- [16] Lihao Ge, Hui Liang, Junsong Yuan, and Daniel Thalmann. 3d convolutional neural networks for efficient and robust hand pose estimation from single depth images. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 5, 2017. **1, 2**
- [17] Lihao Ge, Zhou Ren, and Junsong Yuan. Point-to-point regression pointnet for 3d hand pose estimation. In *Proc. European Conference on Computer Vision (ECCV)*, pages 475–491, 2018. **1, 6**
- [18] Hengkai Guo, Guijin Wang, Xinghao Chen, and Cairong Zhang. Towards good practices for deep 3d hand pose estimation. *arXiv preprint arXiv:1707.07248*, 2017. **1, 2, 5, 6, 7**
- [19] Hengkai Guo, Guijin Wang, Xinghao Chen, Cairong Zhang, Fei Qiao, and Huazhong Yang. Region ensemble network: Improving convolutional network for hand pose estimation. In *Proc. IEEE International Conference on Image Processing (ICIP)*, pages 4512–4516, 2017. **1, 2, 5, 6**
- [20] Albert Haque, Boya Peng, Zelun Luo, Alexandre Alahi, Serena Yeung, and Li Fei-Fei. Towards viewpoint invariant 3d human pose estimation. In *Proc. European Conference on Computer Vision (ECCV)*, pages 160–177, 2016. **1, 2, 5, 7**
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. **2, 4, 8**
- [22] Li He, Guijin Wang, Qingmin Liao, and Jing-Hao Xue. Depth-images-based pose estimation using regression forests and graphical models. *Neurocomputing*, 164:210–219, 2015. **2**
- [23] Meysam Madadi, Sergio Escalera, Xavier Baró, and Jordi Gonzalez. End-to-end global to local cnn learning for hand pose recovery in depth data. *arXiv preprint arXiv:1705.09606*, 2017. **6**
- [24] Daniel Maturana and Sebastian Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *Proc. IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 922–928. IEEE, 2015. **2**
- [25] Gyeongsik Moon, Ju Yong Chang, Yumin Suh, and Kyoung Mu Lee. Holistic planimetric prediction to local volumetric prediction for 3d human pose estimation. *arXiv preprint arXiv:1706.04758*, 2017. **2**
- [26] Gyeongsik Moon, Ju Yong Chang, and Kyoung Mu Lee. V2v-posenet: Voxel-to-voxel prediction network for accurate 3d hand and human pose estimation from a single depth map. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5079–5088, 2018. **1, 2, 4, 5, 6, 7**

- [27] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hour-glass networks for human pose estimation. In *Proc. European Conference on Computer Vision (ECCV)*, pages 483–499, 2016. [1](#), [2](#), [7](#)
- [28] Markus Oberweger and Vincent Lepetit. Deeprior++: Improving fast and accurate 3d hand pose estimation. In *Proc. IEEE International Conference on Computer Vision Workshop (ICCVW)*, volume 840, page 2, 2017. [2](#), [6](#)
- [29] Markus Oberweger, Paul Wohlhart, and Vincent Lepetit. Hands deep in deep learning for hand pose estimation. *arXiv preprint arXiv:1502.06807*, 2015. [2](#)
- [30] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, pages 5099–5108, 2017. [3](#)
- [31] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, pages 91–99, 2015. [5](#)
- [32] Javier Romero, Hedvig Kjellström, and Danica Kragic. Monocular real-time 3d articulated hand pose estimation. In *Proc. IEEE-RAS International Conference on Humanoid Robots (ICHR)*, pages 87–92, 2009. [1](#)
- [33] Jamie Shotton, Andrew Fitzgibbon, Mat Cook, Toby Sharp, Mark Finocchio, Richard Moore, Alex Kipman, and Andrew Blake. Real-time human pose recognition in parts from single depth images. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1297–1304, 2011. [1](#), [2](#), [7](#)
- [34] Xiao Sun, Yichen Wei, Shuang Liang, Xiaoou Tang, and Jian Sun. Cascaded hand pose regression. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 824–832, 2015. [6](#)
- [35] Xiao Sun, Bin Xiao, Fangyin Wei, Shuang Liang, and Yichen Wei. Integral human pose regression. In *Proc. European Conference on Computer Vision (ECCV)*, 2018. [1](#)
- [36] Danhang Tang, Hyung Jin Chang, Alykhan Tejani, and Tae-Kyun Kim. Latent regression forest: Structured estimation of 3d articulated hand posture. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3786–3793, 2014. [2](#), [5](#), [6](#)
- [37] Jonathan Tompson, Murphy Stein, Yann Lecun, and Ken Perlin. Real-time continuous pose recovery of human hands using convolutional networks. *ACM Transactions on Graphics*, 33(5):169, 2014. [1](#), [2](#), [5](#), [6](#), [7](#)
- [38] Alexander Toshev and Christian Szegedy. Deeppose: Human pose estimation via deep neural networks. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1653–1660, 2014. [2](#), [7](#), [8](#)
- [39] Chengde Wan, Thomas Probst, Luc Van Gool, and Angela Yao. Dense 3d regression for hand pose estimation. pages 5147–5156, 2018. [2](#), [6](#)
- [40] Chengde Wan, Thomas Probst, Luc Van Gool, and Angela Yao. Crossing nets: Combining gans and vaes with a shared latent space for hand pose estimation. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. [6](#)
- [41] Keze Wang, Liang Lin, Chuangjie Ren, Wei Zhang, and Wenxiu Sun. Convolutional memory blocks for depth data representation learning. In *Proc. International Joint Conference on Artificial Intelligence (IJCAI)*, pages 2790–2797, 2018. [1](#), [2](#), [7](#)
- [42] Keze Wang, Shengfu Zhai, Hui Cheng, Xiaodan Liang, and Liang Lin. Human pose estimation from depth images via inference embedded multi-task learning. In *Proc. ACM on Multimedia Conference (ACM MM)*, pages 1227–1236, 2016. [1](#), [2](#), [5](#), [7](#), [8](#)
- [43] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4724–4732, 2016. [1](#), [2](#), [7](#)
- [44] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *Proc. European Conference on Computer Vision (ECCV)*, 2018. [2](#), [7](#), [8](#)
- [45] Chi Xu, Lakshmi Narasimhan Govindarajan, Yu Zhang, and Li Cheng. Lie-x: Depth image based articulated object pose estimation, tracking, and action recognition on lie groups. *International Journal of Computer Vision*, 123(3):454–478, 2017. [6](#)
- [46] Mao Ye, Xianwang Wang, Ruigang Yang, Liu Ren, and Marc Pollefeys. Accurate 3d pose estimation from a single depth image. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, pages 731–738, 2011. [2](#)
- [47] Shanxin Yuan, Guillermo Garcia-Hernando, Björn Stenger, Gyeongsik Moon, Ju Yong Chang, Kyoung Mu Lee, Pavlo Molchanov, Jan Kautz, Sina Honari, Liuhaog Ge, et al. Depth-based 3d hand pose estimation: From current achievements to future goals. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2636–2645, 2018. [6](#)
- [48] Shanxin Yuan, Qi Ye, Guillermo Garcia-Hernando, and Tae-Kyun Kim. The 2017 hands in the million challenge on 3d hand pose estimation. *arXiv preprint arXiv:1707.02237*, 2017. [2](#), [5](#), [6](#)
- [49] Shanxin Yuan, Qi Ye, Bjorn Stenger, Siddhant Jain, and Tae-Kyun Kim. Bighand2. 2m benchmark: Hand pose dataset and state of the art analysis. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4866–4874, 2017. [5](#)
- [50] Ho Yub Jung, Soochahn Lee, Yong Seok Heo, and Il Dong Yun. Random tree walk toward instantaneous 3d human pose estimation. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2467–2474, 2015. [1](#), [7](#)
- [51] Xingyi Zhou, Qingfu Wan, Wei Zhang, Xiangyang Xue, and Yichen Wei. Model-based deep hand pose estimation. In *Proc. International Joint Conference on Artificial Intelligence (IJCAI)*, 2016. [6](#)