# Multi-Person 3D Pose and Shape Estimation via Inverse Kinematics and Refinement – Supplementary Material

Junuk Cha<sup>1</sup>[0000-0003-2321-2797]</sup>, Muhammad Saqlain<sup>1,2†</sup>[0000-0001-5877-6432], GeonU Kim<sup>1‡</sup>, Mingyu Shin<sup>1,3‡</sup>, and Seungryul Baek<sup>1</sup>[0000-0002-0856-6880]</sup>

<sup>1</sup>UNIST, South Korea <sup>2</sup>eSmart Systems, Norway <sup>3</sup>Yeongnam Univ., South Korea

In this supplemental document, we try to involve more detailed explanation about our pipeline. For the purpose, we draw the schematic diagram for the entire framework in Fig. 1 and detail the overall sequence of the framework including the rearrange and concatenation method for  $\mathbf{F}_{ing}$  and  $\boldsymbol{\Theta}_{init}$  in Sec. 1. We further offer more details on our training process in Sec. 2, summarize how to calculate absolute 3D skeletons  $\mathbf{P}$  in Sec. 3 and discuss the additional quantitative results in Sec. 4. Finally, we present additional qualitative examples in Sec. 5.

# 1 Overall sequence of framework

In this section, we detail the overall sequence of framework.

Feature extraction. We use ResNet architecture [2] to extract image feature  $\mathbf{F}_{img} \in \mathbb{R}^{8 \times 8 \times 2048}$  from the cropped image  $\mathbf{X} \in \mathbb{R}^{256 \times 256 \times 3}$ . The cropped image  $\mathbf{X}$  goes through CONV, BN, ReLU, MAX Pool, Block 1, Block 2, Block 3, and Block 4 layers to produce the  $8 \times 8 \times 2,048$ -dimensional array. This feature array is used to further 1) estimate initial skeletons from the initial skeleton estimation network  $f^{P}$ , 2) estimate the twist angle  $\Phi$  and shape  $\beta_{init}$  from twist angle and shape estimator  $f^{TS}$  and 3) refine the initial mesh parameters via the relation-aware refiner  $f^{Ref}$ .

Initial skeleton estimation network  $f^{\mathbf{P}}$ . As denoted in Table 1, the obtained feature  $\mathbf{F}_{img}$  is fed to the 1 × 1 conv1 layer to generate 3D heatmaps whose dimensions are  $8 \times 8 \times 8 \times K$ . After that, 3D soft-argmax operation is applied on the 3D heatmaps to obtain the root-relative 3D skeletons  $\mathbf{P}_{rel}$ . Similarly,  $\mathbf{F}_{img}$  is fed to 1×1 conv2 layer to generate 2D heatmaps whose dimensions are  $8 \times 8 \times K$ . Then, 2D soft-argmax operation is applied on the 2D heatmaps to generate 2D skeletons  $\mathbf{P}_{img}$ . Finally, the absolute 3D skeletons  $\mathbf{P}$  are differentiably calculated by combining  $\mathbf{P}_{img}$  and  $\mathbf{P}_{rel}$  with camera intrinsic matrix as in [18].

**Global average pooling.** To be used for twist angle and shape estimator  $f^{\text{TS}}$  and relation-aware refiner  $f^{\text{Ref}}$ , the global average pooling (GAP) is applied on  $\mathbf{F}_{\text{img}}$  to reduce its dimension to  $1 \times 1 \times 2,048$ .

Twist and shape estimation network  $f^{\text{TS}}$ . The twist and shape estimation network  $f^{\text{TS}}$  is composed of multiple fully-connected layers as denoted in Table 2. It first changes the  $1 \times 1 \times 2$ , 048-dimensional feature array into the  $1 \times 1 \times 1$ , 024dimensional array via two fully-connected layers (i.e. FC1 and FC2). Then, it

#### 2 Cha et al.



Fig. 1: The detailed diagram of our entire framework. From the cropped image **X**, feature array  $\mathbf{F}_{img}$  is extracted and it is used for 1) initial skeleton estimation via  $f^{P}$ , for 2) estimating initial mesh parameter  $\boldsymbol{\Theta}_{init} = [\boldsymbol{\theta}_{init}; \boldsymbol{\beta}_{init}; \mathbf{C}_{init}]$  via  $f^{TS}$  and for 3) refining the initial mesh parameter  $\boldsymbol{\Theta}_{init}$  by  $\boldsymbol{\Theta}_{ref} = \boldsymbol{\Theta}_{init} + \Delta \boldsymbol{\Theta}_{ref}$  where  $\Delta \boldsymbol{\Theta}_{ref} = [\Delta \boldsymbol{\theta}_{ref}; \Delta \boldsymbol{\beta}_{ref}; \Delta \mathbf{C}_{ref}]$  via  $f^{Ref}$ . Overall sequences are described in Sec. 1 and detailed architecture for  $f^{P}$ ,  $f^{TS}$  and  $f^{Ref}$  are shown in Tables 1, 2 and 3, respectively.

maps it to the twist angle vector  $\boldsymbol{\Phi} \in \mathbb{R}^{K \times 2}$  and shape vector  $\boldsymbol{\beta}_{\text{init}} \in \mathbb{R}^{1 \times 10}$  via respective fully-connected layers (i.e. FC-twist, FC-shape).

**Relation-aware refinement network**  $f^{\text{Ref}}$ . While both initial skeleton estimation and twist and shape estimation have been performed individually for each person in each image, the relation-aware refinement network  $f^{\text{Ref}}$  needs to deal with multi-persons contained in the image altogether. For the purpose, we rearrange and concatenate the intermediate output vectors as follows:

Concatenating responses for N persons. The intermediate feature obtained by applying GAP on the  $\mathbf{F}_{img}$  is the  $1 \times 1 \times 2$ , 048-dimensional array. We concatenate it for N sampled persons and reshape it towards the  $N \times 1 \times 2$ , 048-dimensional array. In parallel, the pose parameter  $\boldsymbol{\theta}_{init}$ , shape parameter  $\boldsymbol{\beta}_{init}$  and camera parameter  $\mathbf{C}_{init}$  are concatenated for N persons and reshaped into  $N \times K \times 6$ -dimensional array,  $N \times 1 \times 10$  and  $N \times 1 \times 3$ -dimensional array, respectively.

Re-arranging vectors for K joints. Afterwards, we need to align the dimension of the shape  $\beta_{init}$ , camera parameters  $C_{init}$  and feature  $F_{img}$  to that of the pose parameter  $\theta_{init}$  obtained for each joint. For the purpose, we apply fullyconnected layers (i.e. FC-shape-rearrange1, FC-cam-rearrange1 and FC-imgFrearrange1 layers) to make their second dimensions as K. As a result, we obtain the  $N \times K \times 10$ -dimensional array,  $N \times K \times 6$ -dimensional array and  $N \times K \times 2,048$ dimensional array. Finally, we obtain the  $N \times K \times 2,067$ -dimensional array by concatenating resultant shape, camera, pose and image feature arrays. This is the input to the relation-aware refinement network  $f^{\text{Ref}}$  whose architecture is denoted in Table 3.

Transformer outputs. From the Transformer, three arrays (i.e.  $N \times K \times 6$ dimensional array,  $N \times K \times 10$ -dimensional array and  $N \times K \times 3$ -dimensional array) are obtained. The first array is the residual pose vector  $\Delta \theta_{\rm ref} \in \mathbb{R}^{N \times K \times 6}$ . The second and third arrays are transformed via FC-shape-rearrange2 and FCcam-rearrange2 layers towards  $N \times 1 \times 10$  and  $N \times 1 \times 3$ -dimensional arrays, respectively to obtain the shape  $\Delta \beta_{\rm ref} \in \mathbb{R}^{N \times 1 \times 10}$ , and camera  $\Delta \mathbf{C}_{\rm ref} \in \mathbb{R}^{N \times 1 \times 3}$ residual vectors.

**Pose and shape discriminators**  $D_{\theta}$  and  $D_{\beta}$ . We used the same architecture of [8] as pose discriminator  $D_{\theta}$  and shape discriminator  $D_{\beta}$ .

Table 1: Architecture of initial 3D skeleton estimation network  $f^{\rm P}$ . Input is  $\mathbf{F}_{\rm img}$  from ResNet.

Layer	Operation	Kernel	Dimensionality
	Input: $\mathbf{F}_{img}$	-	$8\times8\times2048$
$1 \times 1 \text{ conv1}$	Conv.	$1 \times 1$	$8\times8\times8\times K$
3D Soft-argmax	3D Soft-argmax	-	$K \times 3$
$1 \times 1 \text{ conv}2$	Conv.	$1 \times 1$	$8 \times 8 \times K$
2D Soft-argmax	2D Soft-argmax	-	$K \times 2$
Absolute pose recovery	-	-	$K \times 3$

Table 2: Architecture of twist angle and shape estimation network  $f^{\text{TS}}$ . Input is global average pooled  $\mathbf{F}_{\text{img}}$  from ResNet.

Layer	Operation	Kernel Dimensionality		
	Input: global average pooled $\mathbf{F}_{img}$	-	2048	
FC1	Linear + dropout(0.5)	-	1024	
FC2	Linear + dropout(0.5)	-	1024	
FC-twist	Linear	-	$K \times 2$	
FC-shape	Linear	-	10	

Table 3: Architecture of relation-aware refiner  $f^{\text{Ref}}$ . Input is  $\boldsymbol{\Theta}_{\text{init}}$  concatenated with global average pooled  $\mathbf{F}_{\text{img}}$  from ResNet.

Layer	Operation	Kernel	Dimensionality
	Input: Input patches	-	$N\times K\times 2067$
FC-input	Linear	-	$N\times K\times 1024$
Norm	LayerNorm	-	$N\times K\times 1024$
Multi-Head Attention	Attention	-	$N\times K\times 1024$
Norm	LayerNorm	-	$N\times K\times 1024$
MID	Linear + GELU	-	$N\times K\times 2048$
MILI	Linear	-	$N\times K\times 1024$
Norm	LayerNorm	-	$N\times K\times 1024$
Multi-Head Attention	Attention	-	$N\times K\times 1024$
Norm	LayerNorm	-	$N\times K\times 1024$
MLD	Linear + GELU	-	$N\times K\times 2048$
MILI	Linear	-	$N\times K\times 1024$
MLP-pose	Linear	-	$N \times K \times 6$
MI D shape	Linear	-	$N \times K \times 512$
MLF-snape	Linear	-	$N\times K\times 10$
MI P. com	Linear	-	$N\times K\times 512$
MLF-cam	Linear	-	$N\times K\times 3$

4 Cha et al.

Algorithm 1: The summary of our entire training process
Input: Image I
$\textbf{Output: P, } \beta_{\textbf{init}},  \phi,  \Theta_{\textbf{init}},  \Theta_{\textbf{ref}}$
$\mathbf{for}  t = 1, \dots, T  \mathbf{do}$
Crop the image to produce $\mathbf{X}$ using $M$ bounding boxes.
Obtain initial 3D skeletons ( <b>P</b> ) of each person from $f^{\rm P}$ .
Obtain SMPL shape parameters $(\boldsymbol{\beta}_{\text{init}})$ and twist angles $(\boldsymbol{\Phi})$ from $f^{\text{TS}}$ .
Obtain initial SMPL pose parameters $(\boldsymbol{\Theta}_{init})$ by inverse kinematics.
Sample $\boldsymbol{\Theta}_{\text{init}}$ of N persons.
Feed them to $f^{\text{Ref}}$ concatenated with image features ( $\mathbf{F}_{\text{img}}$ ) and refine
them to get $\boldsymbol{\Theta}_{\mathrm{ref}}$ .
Feed $\boldsymbol{\Theta}_{\text{ref}}$ to discriminators $D_{\theta}$ and $D_{\beta}$ .
Calculate gradient $\nabla L$ (Eq. 1) and update $f^{\rm P}$ , $f^{\rm TS}$ , and $f^{\rm Ref}$
end

# 2 More details on training process

We obtain the initial skeleton estimation network  $f^{\rm P}$ , twist angle and shape estimation network  $f^{\rm TS}$  and relation-aware refinement network  $f^{\rm Ref}$  via optimizing the following loss functions:

$$L(f^{\mathrm{P}}, f^{\mathrm{TS}}, f^{\mathrm{Ref}}) = L_{\mathrm{P}}(f^{\mathrm{P}}) + L_{\mathrm{TS}}(f^{\mathrm{TS}}) + L_{\mathrm{Ref}}(f^{\mathrm{Ref}}).$$
 (1)

where individual terms are defined in the main paper. The overall training procedure is summarized in Algorithm 1.

Datasets. For training, we involved multiple datasets to train our framework. We used the full or partial losses for each datasets according to their groundtruth types: Human3.6M [4] dataset is used for calculating the full losses (i.e.  $L_{\rm P}(f^{\rm P})$ ,  $L_{\rm TS}(f^{\rm TS})$ , and  $L_{\rm Ref}(f^{\rm Ref})$ ), as it provides the SMPL pose and shape parameters, 3D skeleton, 2D skeleton ground truths. MPI-INF-3DHP [15] is used for calculating the partial losses  $L_{\rm P}(f^{\rm P})$  and  $L_{\rm pose}(f^{\rm Ref})$ , as it provides only 2D and 3D skeleton ground-truths. LSP [6], MSCOCO [13] and MPII [1] datasets are used to calculate only the 2D losses in  $L_{\rm P}(f^{\rm P})$  and  $L_{\rm pose}(f^{\rm Ref})$ , as they provide only 2D skeleton ground-truth. Additionally, MuCo-3DHP [16], CMU-Panoptic [7], SAIL-VOS [3], SURREAL [21], AIST++ [12] are used to calculate the  $L_{\rm P}(f^{\rm P})$ .

# 3 Details on calculating P

We followed [18] for calculating the absolute 3D skeletons **P**. Our aim is obtaining absolute 3D skeletons **P** from root-relative 3D skeletons  $\mathbf{P}_{rel}$  and 2D skeletons  $\mathbf{P}_{img}$ . The root-relative 3D skeletons  $\mathbf{P}_{rel}$ , 2D skeletons  $\mathbf{P}_{img}$ , and absolute 3D skeletons can be expressed as  $\{(\mathbf{X}_k, \mathbf{Y}_k, \mathbf{Z}_k)\}_{k=1}^K$ ,  $\{(x_k, y_k)\}_{k=1}^K$ , and  $\{(\mathbf{X}_k + \mathbf{X}_o, \mathbf{Y}_k + \mathbf{Y}_o, \mathbf{Z}_k + \mathbf{Z}_o)\}_{k=1}^K$ , respectively, where K is the number of joints

Table 4: Ablation study on loss Table 5: The effectiveness of re-<br/>on 3DPW.Table 6: Inference frame rate.on 3DPW.finer  $f^{Ref}$  on 3DPW.Table 6: Inference frame rate.

on ode w.			inter j on obriv.		Method	FPS
Method	$\mathrm{MPJPE}(\downarrow)$	$\mathrm{dOrder}(\uparrow)$	Method N	$MPJPE(\downarrow)$	BOMP [11]	22.7
Ours-Lpose-Ladv	68.8	47.9	HybrIK [11]	80.0	PARE [11]	22.1
$Ours-L_{mesh}$	66.3	94.7	HybrIK [11] w/ Ref	78.4	SPEC [11]	21.4
$Ours-L_{pose}$	67.0	51.2	Ours w/o Rei	07.3 66.0	METRO [11]	16.6
Ours- $L_{adv}$	67.0	95.3	Ouis	00.0	Hubrik [11]	24.8
Ours	66.0	96.5				24.0 91 5

and  $(\mathbf{X}_o, \mathbf{Y}_o, \mathbf{Z}_o)$  is the offset. We should recover the offset  $(\mathbf{X}_o, \mathbf{Y}_o, \mathbf{Z}_o)$ . A normalized image coordinates can be calculated as  $(\tilde{x}_k, \tilde{y}_k)^T = K^{-1}(x_k, y_k)^T$ . It can be expressed as follows:

$$\begin{bmatrix} \tilde{x}_k \\ \tilde{y}_k \end{bmatrix} = \begin{bmatrix} (\mathbf{X}_k + \mathbf{X}_o) / (\mathbf{Z}_k + \mathbf{Z}_o) \\ (\mathbf{Y}_k + \mathbf{Y}_o) / (\mathbf{Z}_k + \mathbf{Z}_o) \end{bmatrix}$$
(2)

where  $\tilde{x}_k$ ,  $\tilde{y}_k$ ,  $\mathbf{X}_k$ ,  $\mathbf{Y}_k$ , and  $\mathbf{Z}_k$  are estimated values. This equation can be arranged to

$$\begin{bmatrix} \mathbf{X}_o - \tilde{x}_k \mathbf{Z}_o \\ \mathbf{Y}_o - \tilde{y}_k \mathbf{Z}_o \end{bmatrix} = \begin{bmatrix} \tilde{x}_j \mathbf{Z}_k - \mathbf{X}_k \\ \tilde{y}_j \mathbf{Z}_k - \mathbf{Y}_k \end{bmatrix}$$
(3)

There is K joints, thus we can obtain 2K linear equations with the variables  $(\mathbf{X}_o, \mathbf{Y}_o, \mathbf{Z}_o)$ . We used Cholesky decomposition to solve them. After calculating the offset  $(\mathbf{X}_o, \mathbf{Y}_o, \mathbf{Z}_o)$ , the absolute 3D skeletons **P** whose 2D skeleton lies outside the image are calculated as  $(\mathbf{X}_k + \mathbf{X}_o, \mathbf{Y}_k + \mathbf{Y}_o, \mathbf{Z}_k + \mathbf{Z}_o)^T$  and the others are calculated as  $(\tilde{x}_k, \tilde{y}_k, 1)^T \cdot (\mathbf{Z}_k + \mathbf{Z}_o)$ .

### 4 More quantitative results

#### 4.1 Ablation study on losses

In this section, We studied the effectiveness of losses  $L_{\text{pose}}$ ,  $L_{\text{mesh}}$ , and  $L_{\text{adv}}$  for training  $f^{\text{Ref}}$ . The results are shown in Table 4. When using only  $L_{\text{mesh}}$ , the  $f^{\text{ref}}$  is trained on Human3.6M which has SMPL parameters ground-truth. Its generalization performance is worse than using a variety of dataset. When we do not involve  $L_{\text{mesh}}$ , it is trained on various datasets (MPI-INF-3DHP, LSP, COCO, etc.). Its performance is better than other methods except 'Ours'. When we do not involve  $L_{\text{pose}}$ , it is trained on Human3.6M by  $L_{\text{mesh}}$  and various datasets by  $L_{\text{adv}}$ . It has better performance compared to 'Ours- $L_{\text{pose}}$ - $L_{\text{adv}}$ '.  $L_{\text{adv}}$  punishes it for wrong estimated pose and shape parameters. When we do not involve  $L_{\text{pose}}$ . Finally, the method using three losses (ie.  $L_{\text{pose}}$ ,  $L_{\text{mesh}}$  and  $L_{\text{adv}}$ ) has the best performance. Further, we extend our framework to use the original camera parameter and measure the depth-order accuracy similar to [5, 20, 22] and its result is shown in 'dOrder' column of the Table. We could see that results are gradually increasing as more loss functions are involved.

6 Cha et al.

# 4.2 Effectiveness of refiner $f^{\text{Ref}}$

We conducted the experiment to find out whether the refiner module  $f^{\text{Ref}}$  is really improve the final mesh quality. Table 5 shows that mesh obtained from HybrIK [11] can be improved by our refiner module  $f^{\text{Ref}}$  as well as our coarse mesh. However, its gain is limited because HybrIK uses less competitive 3D poses.

### 4.3 Inference time

We compare the inference frame rate of our method with other state-of-thearts. In Table 6, it is confirmed that the inference speed is comparable to other methods.



Fig. 2: Qualitative comparisons with HybrIK [11]. Red circles highlight wrongly estimated parts.

# 5 More qualitative results

Fig. 2 shows our qualitative results compared to HybrIK [11]. In Figs. 3, 4 and 5, we present more qualitative results compared to three state-of-the-art methods [10, 19, 9] on three multi-person pose estimation benchmark datasets (i.e. 3DPW [14], AGORA [16] and MuPoTS [17]), respectively. Fig. 6 illustrates the intermediate outputs (i.e. initial skeletons, initial meshes obtained by inverse kinematics and refined meshes) obtained from our pipeline. Figs. 7 and 8 show the top-view and second-view results and failure cases, respectively.



Fig. 3: Qualitative comparisons on 3DPW [14]. Red circles highlight wrongly estimated parts.



Fig. 4: Qualitative comparisons on AGORA [17]. Red circles highlight wrongly estimated parts.



Fig. 5: Qualitative comparisons on MuPoTS [16]. Red circles highlight wrongly estimated parts.



Fig. 6: Example outputs from our pipeline: (a) input RGB image, (b) initial skeleton estimation results obtained from the input image, (c) initial meshes obtained from the inverse kinematics process, (d) refined meshes obtained from the refinement Transformer.



11

Fig. 7: Example outputs from our pipeline: (a) input RGB image, (b) refined meshes overlaid on input RGB image, (c) top view, (d) side view.



Fig. 8: Failure cases: (a) input RGB image, (b) initial skeleton estimation results obtained from the input image, (c) initial meshes obtained from the inverse kinematics process, (d) refined meshes obtained from the refinement Transformer.

# References

- 1. Andriluka, M., Pishchulin, L., Gehler, P., Schiele, B.: 2d human pose estimation: New benchmark and state of the art analysis. In: CVPR (2014)
- He, K., Zhang, X., Ren, S., Sun, J.: Identity mappings in deep residual networks. In: ECCV (2016)
- Hu, Y.T., Chen, H.S., Hui, K., Huang, J.B., Schwing, A.G.: Sail-vos: Semantic amodal instance level video object segmentation-a synthetic dataset and baselines. In: CVPR (2019)
- Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C.: Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. TPAMI (2013)
- Jiang, W., Kolotouros, N., Pavlakos, G., Zhou, X., Daniilidis, K.: Coherent reconstruction of multiple humans from a single image. In: CVPR (2020)
- Johnson, S., Everingham, M.: Clustered pose and nonlinear appearance models for human pose estimation. In: BMVC (2010)
- Joo, H., Liu, H., Tan, L., Gui, L., Nabbe, B., Matthews, I., Kanade, T., Nobuhara, S., Sheikh, Y.: Panoptic studio: A massively multiview system for social motion capture. In: ICCV (2015)
- 8. Kanazawa, A., Black, M.J., Jacobs, D.W., Malik, J.: End-to-end recovery of human shape and pose. In: CVPR (2018)
- Kocabas, M., Huang, C.H.P., Hilliges, O., Black, M.J.: Pare: Part attention regressor for 3d human body estimation. In: ICCV (2021)
- Kocabas, M., Huang, C.H.P., Tesch, J., Müller, L., Hilliges, O., Black, M.J.: Spec: Seeing people in the wild with an estimated camera. In: ICCV (2021)
- Li, J., Xu, C., Chen, Z., Bian, S., Yang, L., Lu, C.: Hybrik: A hybrid analyticalneural inverse kinematics solution for 3d human pose and shape estimation. In: CVPR (2021)
- Li, R., Yang, S., Ross, D.A., Kanazawa, A.: Learn to dance with aist++: Music conditioned 3d dance generation. arXiv:2101.08779 (2021)
- Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: ECCV (2014)
- von Marcard, T., Henschel, R., Black, M.J., Rosenhahn, B., Pons-Moll, G.: Recovering accurate 3d human pose in the wild using imus and a moving camera. In: ECCV (2018)
- Mehta, D., Rhodin, H., Casas, D., Fua, P., Sotnychenko, O., Xu, W., Theobalt, C.: Monocular 3d human pose estimation in the wild using improved cnn supervision. In: 3DV (2017)
- Mehta, D., Sotnychenko, O., Mueller, F., Xu, W., Sridhar, S., Pons-Moll, G., Theobalt, C.: Single-shot multi-person 3d pose estimation from monocular rgb. In: 3DV (2018)
- 17. Patel, P., Huang, C.H.P., Tesch, J., Hoffmann, D.T., Tripathi, S., Black, M.J.: Agora: Avatars in geography optimized for regression analysis. In: CVPR (2021)
- Sárándi, I., Linder, T., Arras, K.O., Leibe, B.: Metrabs: Metric-scale truncationrobust heatmaps for absolute 3d human pose estimation. IEEE Transactions on Biometrics, Behavior, and Identity Science (2020)
- Sun, Y., Bao, Q., Liu, W., Fu, Y., Black, M.J., Mei, T.: Monocular, one-stage, regression of multiple 3d people. In: ICCV (2021)
- Ugrinovic, N., Ruiz, A., Agudo, A., Sanfeliu, A., Moreno-Noguer, F.: Body size and depth disambiguation in multi-person reconstruction from single images. In: 3DV (2021)

- 14 Cha et al.
- Varol, G., Romero, J., Martin, X., Mahmood, N., Black, M.J., Laptev, I., Schmid, C.: Learning from synthetic humans. In: CVPR (2017)
- 22. Zhang, J., Yu, D., Liew, J.H., Nie, X., Feng, J.: Body meshes as points. In: CVPR (2021)