

Monocular, One-stage, Regression of Multiple 3D People

Supplementary Material

1. Introduction

This supplementary material first provides details about the network architecture design in Sec. 2, followed by hyper-parameter configurations and implementation details in Sec. 3. Finally, we present more qualitative results in Sec. 4, including an analysis of failure cases.

2. Ablation study on the architecture

The common architecture underlying most existing methods uses a ResNet-50 backbone followed by a single head that produces SMPL parameters. Most existing methods use an iterative approach, first introduced by HMR [9]. To arrive at the ROMP architecture in Fig. 1, we progressively add/change elements of this basic design and evaluate the effects in an ablation study. As shown in Fig. 2, two main design choices of the head architecture are explored, single head (SH) vs. separated multiple heads (MH), and the number of convolution blocks (NB) in each head. In this section, we first introduce the experiment settings of the ablation study and then report the results.

2.1. Experiment settings

Architecture details. 1) SH v.s. MH: The architecture of the single head model is presented as (b) in Fig. 2. It is a straightforward and lightweight design, which uses a single branch to jointly estimate the Camera and SMPL maps. In contrast, the separated multi-head design is adopted in Fig. 1. 2) NB: We adopt ResNet blocks (shown in Fig. 1) as the basic unit of the head networks. We conduct an ablation study to determine the proper number of blocks in each head. The architecture is presented as (c) in Fig. 2.

Datasets. For pretraining, we take three 2D pose datasets, COCO [15], CrowdPose [14], AICH [20] and one detection dataset, CrowdHuman [18]. For the formal training, we take four 3D pose datasets, Human3.6M [3], MPI-INF-3DHP [16], MuCo-3DHP [16], UP [13], and three 2D pose datasets, COCO [15], MPII [1], LSP [5]. Especially, we only use the 7 subjects (Subject 1, 2, 3, 4, 5, 6, and 7) of MPI-INF-3DHP for training. Subject 8 is used for validation. The test set of MPI-INF-3DHP is employed for evaluation.

Components				MPI-INF-3DHP	
SH	MH	NB	Backbone	MPJPE↓	PMPJPE↓
✓		1	HRNet-32	95.31	66.73
	✓	1	HRNet-32	95.11	65.89
✓		2	HRNet-32	97.19	65.50
	✓	2	HRNet-32	95.30	66.32
	✓	3	HRNet-32	97.19	65.94
✓		1	ResNet-50	105.31	69.62
	✓	1	ResNet-50	106.43	69.83
✓		2	ResNet-50	105.37	69.12
	✓	2	ResNet-50	104.92	70.20
	✓	3	ResNet-50	105.41	71.78

Table 1. Ablation study on the architecture design. SH denotes the single head design and MH is the separated multi-head design. NB is the number of blocks used in each head.

Training. Firstly, we pretrain the ResNet-50/HRNet-32 model on the pretraining 2D pose datasets for 120 epochs. The network architecture for pretraining is presented as (a) in Fig. 2. We follow the hyper-parameter settings of HigherHRNet [2] during this process.

During formal training, we set the learning rate to $5e^{-5}$, weight decay to $1e^{-6}$, batch size to 64. To achieve the best performance of each architecture as much as possible, the loss weights are adjusted, according to the visualization results on the validation set, to avoid the model falling into a local minimum. Each model has been trained at least 100 epochs. We observe that the model with fewer head layers often achieves its best performance at about 120 epochs after several rounds of adjusting hyperparameters, while the model with more head layers often achieves its best performance at about 50 epochs.

2.2. Result analysis

In Tab. 1, we present the results of different architecture designs. The experiments are performed with two kinds of backbone, HRNet-32 and ResNet-50. In this settings, HRNet-32 performs better than ResNet-50. Just like its superior performance of fine-tuning on 3DPW (Tab. 3 in the main paper), HRNet-32 once again proves its excellent ability to fit a specific data domain on MPI-INF-3DHP.

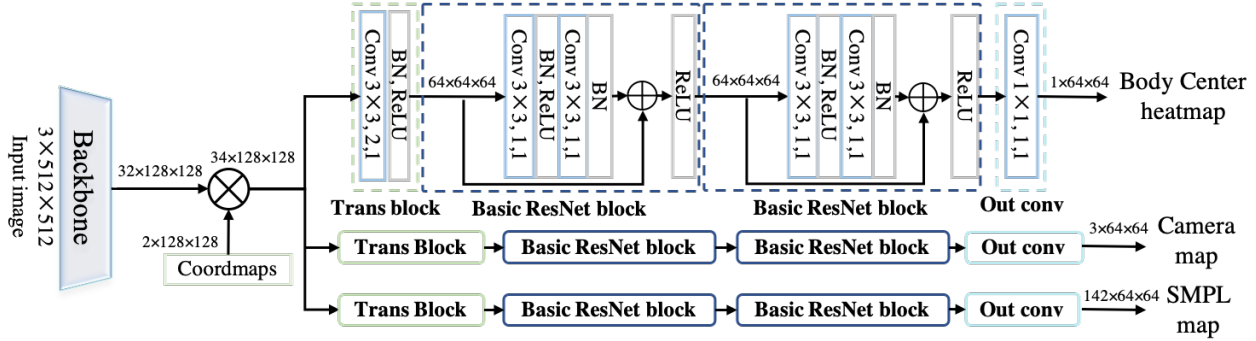


Figure 1. Architecture of the proposed ROMP.

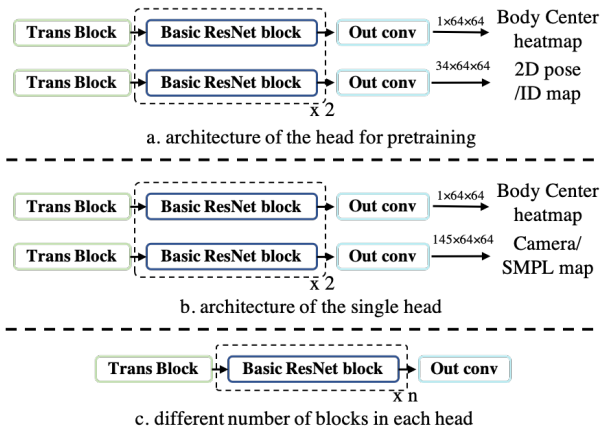


Figure 2. Architecture of the pre-train model.

Using the same backbone, the performance gap between different head architecture designs is relatively small. Among these designs, we find out that 1) compared with the SH, the disentangled multi-head design performs better in the most cases and is prone to train; 2) regarding the NB, setting $n = 2$ is a better choice to balance accuracy and training time.

Fig. 1 shows the details of the architecture. ROMP adopts a fully convolutional multi-head design. Compared with previous methods (like [9, 11, 12]), we do not need to use iterative regression. It is interesting because previous methods rely on this iterative updating to achieve good pose accuracy. ROMP has a harder task than these previous methods that are given a cropped image of the person. ROMP needs to detect people, sort out what image features belong to which person, and estimate the pose. By learning from a holistic view of the whole image, ROMP is forced to learn more about people and how they appear in images. For example, people at the edge of the image usually tend to be truncated. In addition, the holistic view provides the opportunity of learning the interaction between multiple people, which helps handle the crowded scenes. Since ROMP has to solve the pose-estimation problem given an image of the

whole scene, it must learn more distinguishable features to solve the task. We posit that these more powerful features enable it to estimate the body shape and pose without iteration.

3. Implementation Details

3.1. Training Strategy

The basic settings of pretraining and formal training are introduced in Sec. 2.1. Here, we introduce the detailed strategy we used for training ROMP. Our training uses 4 NVIDIA P40 GPUs with a batch size of 128. We adopt the Adam optimizer [10] for training. To avoid the multi-step training and adapt to people of diverse scales, we take both the cropped single-person images and the whole images as input. The ratio of loading the entire images is first set to 10% in the first 60 epochs, and then adjusted to 60% in the remaining 60 epochs. Especially, to accelerate the training and reduce the GPU memory usage, we use the automatic mixed precision (AMP) training of Pytorch [17].

3.2. Effect of hyper-parameter configurations during inference

ROMP has several hyper-parameters that can be adjusted during inference to adapt to different scenes, particularly the confidence threshold t_c of the Body Center map and the maximum number of people in an image, N . The confidence threshold t_c is used to filter out the detected people with the confidence value lower than t_c . Similarly, we set the max person number N to take the top N detected people (sorted by their confidence value on the Body Center map). Changing these parameters only affects the number of output bodies and does not require retraining the model.

We observe that a higher t_c filters out inaccurate/untrusted predictions, while a lower t_c leaves more detection results. The setting of the max person number N follows the same rule. The qualitative ablation study in Fig. 3 illustrates this conclusion. For evaluation on all benchmarks, we have set $t_c = 0.2$ and $N = 64$ for a fair comparison.



Figure 3. Qualitative ablation study of the confidence threshold t_c , on a crowded image.

Sequence Name	Frame Ranges
downtown_bus_00	1620-1900
courtyard_hug_00	100-500
courtyard_dancing_00	60-370
courtyard_dancing_01	60-270
courtyard_basketball_00	200-280
courtyard_captureSelfies_00	500-600

Table 2. Video sequences of the person-occluded subset, 3DPW-PC, in 3DPW.

3.3. Depth Ordering

For better visualization, we attempt to approximate the depth ordering between the estimated multi-person body meshes to render the meshes onto the original 2D images. In detail, we construct a depth ordering map to determine the visible meshes in front, using 2D body scale and center confidence as the cue. First, we sort by the body center confidence value from largest to smallest. Second, for each body mesh, we compute the 2D area of the body projected onto the image; this gives an approximate measure of its scale. Finally, we adjust the visible mesh at each position according to their 2D areas. Specifically, at a certain position, if the area value of the invisible body mesh is greater than the currently visible body mesh by a threshold, we swap their positions. In this way, we bring to the front, the body mesh that occupies a larger area on the image plane. This is an approximate solution and future work should explore an integrated solution for depth ordering during inference.

3.4. Datasets

3DPW-PC and 3DPW-OC are the subsets of the 3DPW [19] dataset that are used to evaluate the performance under person or object occlusion respectively. **3DPW-NC** is simply the rest images in the 3DPW. 3DPW-PC contains 1314 frames of 6 person-occluded video sequences. They contain severe person-person occlusion cases with at least two-people overlapping. Details are provided in Tab. 2. Following Zhang et al. [21], 3DPW-OC contains 23 object-occluded video sequences. Please refer to [21] for the details.

Crowdpose [14] is a crowded dataset with 2D pose annotations. It contains an abundant variety of person-

person and person-object occlusion. Specifically, it contains 20,000 images with about 80,000 persons. We employ their default splits with 9,963 samples for training, 7,991 test samples, and 1,997 validation samples.

3DOH50K [21] is a 3D human occlusion dataset. In the image, the human body is occluded by various objects, such as a laptop computer, box, chair, etc. It contains 50,310 images for training and 1,290 images for testing. It is used to evaluate the performance under object occlusion.

MPI-INF-3DHP [16] is a single-person multi-view 3D pose dataset. It contains 8 actors performing 8 activities. Over 1.3M frames are captured from all 14 cameras. Except for the indoor RGB videos of a single person, they also provide MATLAB code to generate a multi-person dataset, MuCo-3DHP, via mixing up segmented foreground human appearance.

COCO [15], **MPII** [1], **LSP** [5], **LSP Extended** [6], and **AICH** [20] are in-the-wild 2D pose datasets. We use them for training. Especially, we use the pseudo SMPL annotations of part images generated by [8, 12] for training.

3.5. Evaluation Metrics

Per-vertex error (PVE) measures the average Euclidean distance from the 3D body mesh predictions to the ground truth after aligning the pelvis keypoints in millimeters.

The mean per joint position error (MPJPE) measures the average Euclidean distance from the 3D pose predictions to the ground truth in millimeters. The predictions are first translated to match the ground truth. Generally, they are aligned by the pelvis keypoint for comparison.

Procrustes-aligned MPJPE (PMPJPE) is MPJPE after rigid alignment of the predicted pose with ground truth in millimeters. It is also called as the reconstruction error. Through Procrustes alignment, the effects of translation, rotation, and scale are eliminated, thus PMPJPE focuses on evaluating the accuracy of the reconstructed 3D skeleton.

3D PCK & AUC. 3D PCK is the 3D version of the 2D PCK metric, which computes the Percentage of Correct Keypoints. Following the ECCV 2020 3DPW Challenge, the threshold of successful prediction is set to 50mm. Correspondingly, the AUC, which is the total area under the PCK-threshold curve, is calculated by computing PCKs by varying the threshold from 0 to 200mm.

The mean per joint angle error (MPJAE) measures



Figure 4. Qualitative results on 3DPW [19], 3DOH50K [21], and CMU Panoptic [7] from top to down.

the angle between the predicted joint orientation and the ground truth orientation in degrees. The orientation difference is measured as the geodesic distance in $SO(3)$. Specifically, only the angles of four limbs and the root are used for evaluation.

Procrustes-aligned MPJAE (PMPJAE) measures the MPJAE after applying the rotation matrix, obtained from the Procrustes alignment, on all predicted orientations. As above, it neglects the global mismatch.

Average Precision (AP) measures multi-person 2D pose accuracy. We employ it to measure the accuracy of the back-projected 2D body keypoints for evaluating the performance on crowded scenes. A detected keypoint candidate is considered to be correct (true positive) if it lies within a threshold of the ground-truth. Each keypoint separately calculates its correspondence with the ground-truth poses. AP correctly penalizes both missed detections and false positives. In [15], for multi-person pose estimation, AP is further designed by defining the object keypoint similarity (OKS) that is a similarity measure between the predictions and the ground truth. Analog to IoU in object detection, OKS is defined as

$$OKS = \frac{\sum_i \exp(-d_i^2/2s^2k_i^2)\delta(v_i > 0)}{\sum_i \delta(v_i > 0)}, \quad (1)$$

where d_i is the Euclidean distance between the detected keypoint and the corresponding ground truth, v_i is the ground-truth visibility flag, s is the person scale, and k_i is a per-keypoint constant that controls falloff. For each keypoint the OKS ranges between 0 and 1.

Given the OKS over all labeled keypoints, average precision (AP) and average recall (AR) can be computed. By

tuning OKS values, the precision-recall curve can be calculated, and AP and AR at different OKS can thoroughly reflect the performance of the testing algorithms. Here, we adopt $AP^{0.5}$ (AP at OKS = 0.50) and $AR^{0.5}$ for evaluation.

4. Qualitative Results

First, in Fig. 4, we present some qualitative results on evaluation benchmarks that are representative of our results. Next, we present more results on in-the-wild images in Fig. 5. Finally, in Fig. 6, we present some failure cases in estimating depth ordering, detection, and 3D pose. Additionally, we also present results of CRMH [4] on these failure cases for comparison.

Discussion of failure cases. Fig. 6 shows the performance of ROMP in estimating the depth order for complex scenes of overlapping people. It illustrates that our body-level depth ordering is limited in cases of extreme crowding with complex depth relationships. CRMH [4] is a Faster-RCNN-based multi-person state-of-the-art method that supervises mesh interpenetration and depth ordering at the vertex level. It has advantages for determining the multi-person depth ordering in crowded scenes. In contrast, ROMP produces more robust and accurate pose estimation in crowded scenes. In future work, we intend to develop a fine-grained depth estimation approach to tackle this problem.

Fig. 6 also shows some extremely challenging images in terms of pose and occlusion. Images like them clearly challenge the state of the art but may also be challenging for humans to perceive.

References

- [1] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2D human pose estimation: New benchmark and state of the art analysis. In *CVPR*, 2014. 1, 3
- [2] Bowen Cheng, Bin Xiao, Jingdong Wang, Honghui Shi, Thomas S Huang, and Lei Zhang. Higherhrnet: Scale-aware representation learning for bottom-up human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5386–5395, 2020. 1
- [3] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6M: Large scale datasets and predictive methods for 3D human sensing in natural environments. *TPAMI*, 2014. 1
- [4] Wen Jiang, Nikos Kolotouros, Georgios Pavlakos, Xiaowei Zhou, and Kostas Daniilidis. Coherent reconstruction of multiple humans from a single image. In *CVPR*, 2020. 4
- [5] Sam Johnson and Mark Everingham. Clustered pose and nonlinear appearance models for human pose estimation. In *BMVC*, 2010. 1, 3
- [6] Sam Johnson and Mark Everingham. Learning effective human pose estimation from inaccurate annotation. In *CVPR*, 2011. 3
- [7] Hanbyul Joo, Hao Liu, Lei Tan, Lin Gui, Bart Nabbe, Iain Matthews, Takeo Kanade, Shohei Nobuhara, and Yaser Sheikh. Panoptic studio: A massively multiview system for social motion capture. In *ICCV*, 2015. 4
- [8] Hanbyul Joo, Natalia Neverova, and Andrea Vedaldi. Exemplar fine-tuning for 3D human pose fitting towards in-the-wild 3D human pose estimation. In *ECCV*, 2020. 3
- [9] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018. 1, 2
- [10] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv*, 2014. 2
- [11] Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. VIBE: Video inference for human body pose and shape estimation. In *CVPR*, 2020. 2
- [12] 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 *ICCV*, 2019. 2, 3
- [13] Christoph Lassner, Javier Romero, Martin Kiefel, Federica Bogo, Michael J Black, and Peter V Gehler. Unite the people: Closing the loop between 3D and 2D human representations. In *CVPR*, 2017. 1
- [14] Jiefeng Li, Can Wang, Hao Zhu, Yihuan Mao, Hao-Shu Fang, and Cewu Lu. CrowdPose: Efficient crowded scenes pose estimation and a new benchmark. In *CVPR*, 2019. 1, 3
- [15] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 1, 3, 4
- [16] 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, 3
- [17] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019. 2
- [18] Shuai Shao, Zijian Zhao, Boxun Li, Tete Xiao, Gang Yu, Xiangyu Zhang, and Jian Sun. Crowdhuman: A benchmark for detecting human in a crowd. *arXiv preprint arXiv:1805.00123*, 2018. 1
- [19] Timo von Marcard, Roberto Henschel, Michael Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3D human pose in the wild using imus and a moving camera. In *ECCV*, 2018. 3, 4
- [20] Jiahong Wu, He Zheng, Bo Zhao, Yixin Li, Baoming Yan, Rui Liang, Wenjia Wang, Shipei Zhou, Guosen Lin, Yanwei Fu, et al. Ai challenger: A large-scale dataset for going deeper in image understanding. *arXiv*, 2017. 1, 3
- [21] Tianshu Zhang, Buzhen Huang, and Yangang Wang. Object-occluded human shape and pose estimation from a single color image. In *CVPR*, 2020. 3, 4



Figure 5. Qualitative results on in-the-wild images.



(a) Input images

(b) Predicted Center maps

(c) Our results

(d) Results of CRMH

Figure 6. Failure cases on estimating depth ordering, 3D pose, and detection results in extremely challenging scenes.