

Occluded Human Mesh Recovery

Rawal Khirodkar¹ Shashank Tripathi² Kris Kitani¹
¹Carnegie Mellon University ²Max Planck Institute for Intelligent Systems
<https://rawalkhirodkar.github.io/ochmr>



Figure 1. To handle severe person-person occlusion, our proposed method OCHMR conditions the deep network on image spatial context using predicted body centermaps. OCHMR is trained using multi-person mesh interpenetration and depth ordering losses. In comparison to bottom-up ROMP [58], top-down OCHMR outputs pixel aligned mesh estimates for each individual under occlusion.

Abstract

Top-down methods for monocular human mesh recovery have two stages: (1) detect human bounding boxes; (2) treat each bounding box as an independent single-human mesh recovery task. Unfortunately, the single-human assumption does not hold in images with multi-human occlusion and crowding. Consequently, top-down methods have difficulties in recovering accurate 3D human meshes under severe person-person occlusion. To address this, we present Occluded Human Mesh Recovery (OCHMR) - a novel top-down mesh recovery approach that incorporates image spatial context to overcome the limitations of the single-human assumption. The approach is conceptually simple and can be applied to any existing top-down architecture. Along with the input image, we condition the top-down model on spatial context from the image in the form of body-center heatmaps. To reason from the predicted body centermaps, we introduce Contextual Normalization (CoNorm) blocks to adaptively modulate intermediate features of the top-down model. The contextual conditioning helps our model disambiguate between two severely overlapping human bounding-boxes, making it robust to multi-person occlusion. Compared with state-of-the-art methods, OCHMR achieves superior performance on challenging multi-person benchmarks like 3DPW, CrowdPose and OCHuman. Specifically, our proposed contextual reasoning architecture applied to the SPIN model with ResNet-50 backbone results in 75.2 PMPJPE on

3DPW-PC, 23.6 AP on CrowdPose and 37.7 AP on OCHuman datasets, a significant improvement of 6.9 mm, 6.4 AP and 20.8 AP respectively over the baseline. Code and models will be released.

1. Introduction

Estimating accurate 3D human meshes from single images has diverse applications in modeling human-scene interactions, understanding human behaviour, AR/VR and robotics. While recent approaches [4, 11, 25, 32, 43, 48, 50, 66] perform particularly well in images containing a single person, human mesh recovery for complex real-world scenes with multiple occluded people remains a challenging task. This can be attributed in part to simplifying assumptions made by existing methods. For instance, most top-down approaches expect a single subject in the input image, which affects robustness under in-the-wild scenarios containing severe person-person occlusion, such as crowding. In this paper, we address human mesh recovery in multi-person scenarios by mitigating the limitations of the single-person assumption of top-down approaches.

Current human mesh recovery methods can be categorized into *top-down* and *bottom-up* methods. Top-down methods [7, 10, 12, 24, 31, 35–37, 67] reduce the problem to a simpler task of single human mesh recovery by relying on a person detector to detect individual bounding box for each person in the image. Since each bounding box is

scaled to the same size, top-down methods are less sensitive to scale variations among subjects and can achieve pixel accurate mesh alignment [10, 67]. In contrast, bottom-up methods [58, 64, 68] simultaneously predict meshes for all subjects in the input image but are limited to a fixed input resolution due to computational constraints. *e.g.* ROMP [58], a bottom-up method, recovers a limited number of human meshes from a resized 512×512 input whereas SPIN [32], a top-down method, scales each bounding box to 224×224 , retaining higher input resolution per person (see Fig. 1). This observation has also been discussed by Cheng *et al.* [6] albeit in the context of 2D human pose estimation. Thus, top-down methods are currently the best performers on various multi-human benchmarks [18, 21, 23, 46, 55, 61]. Despite the advantages, due to the single-human assumption, when presented with multi-human inputs like crowded scenes, top-down methods are forced to select a single plausible mesh per detection bounding box. Bottom-up methods do not have this limitation and typically perform better under occlusion.

A general method should have both traits – be robust to scale variations and person-person occlusions. To this end, we rethink top-down human mesh recovery by predicting *multiple* meshes from the input bounding box. We condition the top-down model on image *spatial-context* in the form of body-center maps, refer Fig. 2. Our choice of using center maps for representing humans under occlusion is inspired by crowd-counting literature [44, 56, 62] and recent works in detection [8, 71, 72]. Our method, OCHMR, predicts the output mesh from the input image for the person of interest in the subject-specific *local* center-map. Similar to bottom-up methods, we also use information from the *global* center-map for understanding overall scene context, which is helpful for occlusion reasoning. With this strategy, we obtain the best of both worlds – OCHMR achieves pixel accurate mesh alignment similar to top-down methods and is robust to occlusions similar to bottom-up methods (See Fig. 1).

To design a top-down architecture capable of contextual conditioning using centermaps, we adopt the mechanism of feature normalization [9, 16, 51] and propose a novel Context Normalization (CoNorm) block to process the global and local centermaps. The CoNorm blocks are used to inject contextual information at multiple depths in the deep feature backbone network. The spatial context is necessary for 3D occlusion reasoning, and the CoNorm block allows for adaptive normalization of intermediate features of the network without changing the backbone. We show that unlike *early fusion* (*e.g.* channel-wise concatenation) of centermaps with input image I , CoNorm can effectively utilize the contextual information from the image. OCHMR is general and can be extended to other top-down human mesh recovery methods with minimal effort.

While the use of spatial-context allows our method to reason about occlusions, our method must also reason about the

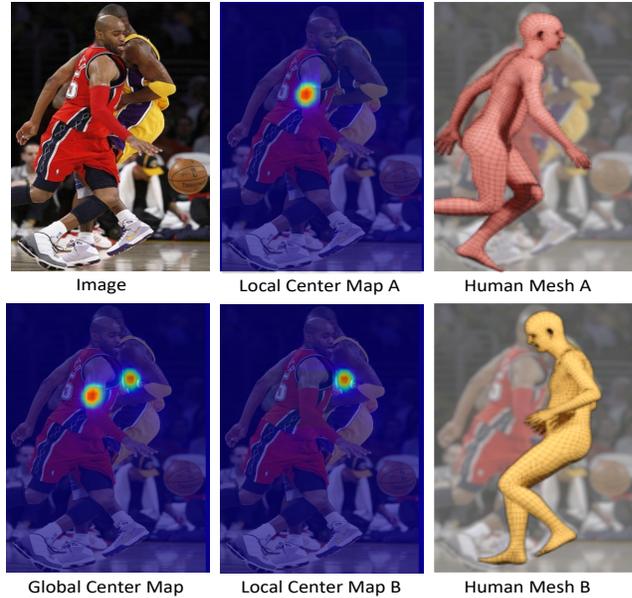


Figure 2. OCHMR leverages image spatial-context for occlusion reasoning by predicting body centermaps. The deep network predicts the mesh output using input image, the subject-specific local centermap and the scene-specific global centermap.

intersection of a set of 3D human meshes. To address this, following CRMH [19], we use an interpenetration loss to penalize intersections among reconstructed meshes and a differentiable depth-ordering loss for depth-consistent human mesh recovery. Furthermore, we make use of training-time data augmentation like scaling and cropping, which affords OCHMR the ability to predict meshes from a variety of body-center locations. We show that our proposed method is also robust to errors in estimated body-centers under severe occlusion. Our empirical results show that OCHMR does not require precise centermaps that correspond to actual body-centers but can also work with any point in its vicinity.

Overall, OCHMR outperforms both top-down and bottom-up methods on various datasets. For challenging datasets such as 3DPW-PC [70], CrowdPose [34] and OCHuman [69], containing a larger proportion of cluttered scenes (with multiple overlapping people), OCHMR sets a new state-of-the-art for 3D reconstruction error (PMPJPE) and 2D keypoint average precision (AP) achieving 77.1 PMPJPE, 21.4 AP and 24.8 AP respectively on the *val* sets outperforming bottom-up methods (Tab. 1). Further, when evaluating using ground-truth bounding boxes, OCHMR dramatically improves SPIN [32] by 20.8 AP and 6.4 AP on the OCHuman and CrowdPose dataset respectively. In summary:

- OCHMR advances top-down human mesh recovery methods by addressing limitations caused by the single-human assumption. Our method leverages spatial-context in the form of centermaps to predict multiple

mesh outputs from an input image.

- We introduce novel Context Normalization (CoNorm) blocks to inject global and local centermap information at multiple depths of the top-down network.
- Our approach achieves state-of-the-art results on the occluded 3DPW-PC, CrowdPose and OCHuman datasets. Empirically, we also show that OCHMR is resilient to noisy body center estimates and demonstrates robust 3D reasoning using multi-person losses.

2. Related Work

Deep learning has significantly advanced 3D human mesh recovery [7, 10, 12, 24, 29–32, 35–37, 60, 67], facilitating the more challenging task of mesh recovery under severe multi-person occlusion [19, 34, 58, 69, 70], which is the main focus of this work.

Biased Human Mesh Recovery Benchmarks. Most benchmark datasets [18, 21, 23, 46, 55, 61, 63] used for learning human mesh recovery focus on a single person and do not accurately represent the distribution of possible occlusions present in the real world. Human3.6M [18], HumanEva [55] and TotalCapture [23] are popular datasets collected using motion capture (mocap) systems using optical markers. While providing accurate annotations, they only have a single subject in the image with limited image complexity due to the lack of background variation. In contrast, datasets like MPI-INF-3DHP [46], PanopticStudio [21] and 3DPW [61] contain multi-person annotations but have limited person-person occlusion – less than 27% of all annotations have crowding (at IoU 0.5). Although previous methods [24, 29, 30, 32, 35, 36, 67] leverage 2D keypoint annotations from datasets like COCO [39], MPII [1], LSP-Extended [20], the 2D datasets are also known to contain similar biases [28, 53, 69]. These biases have affected critical design decisions in state-of-the-art methods which lead to poor generalization under heavy occlusion [19, 58]. Recently, challenging datasets such as OCHuman [69], CrowdPose [34] and 3DPW-PC [70] containing heavy occlusion have been proposed to capture these biases. OCHMR shows a significant improvement over existing works under such challenging conditions.

Top-Down Human Mesh Recovery. Top-down methods [7, 10, 12, 24, 31, 35–37, 67] estimate 3D human mesh of a single person within a person bounding box. The bounding box is usually generated using a person detector [3, 5, 14, 38, 41, 52]. As the input bounding boxes are cropped and scaled to the same size, top-down methods are less sensitive to person scale variations in the image. In contrast, bottom-up methods have to deal with scale variations which compromises pixel alignment in the reconstruction results. For these reasons [6], most state-of-the-art 2D pose

estimation methods [40, 45, 57, 65] are also top-down. However, top-down methods inherently assume a single person in the input image and often fail under occlusions in multi-person scenarios. Recent works like [4, 7, 22, 47, 54, 60] use 2D/3D poses as input along with bounding boxes for human mesh recovery. However, obtaining accurate 2D poses under occlusion is difficult and pose errors like joint swaps [53] are magnified during the 3D reconstruction [22]. CRMH [19] handles multi-person scenarios by using RoI-aligned [14] features of each person to predict the SMPL [43] parameters. However, the reliance on bounding-box-level features makes it hard to effectively differentiate between two overlapping bounding boxes. OCHMR resolves these issues by conditioning the top-down model on image context in the form of *body-centers* – a representation which helps in resolving ambiguity under multi-person occlusion.

Bottom-Up Human Mesh Recovery. Unlike top-down methods, only few methods exist which use the bottom-up paradigm for human mesh recovery. Zanfir *et al.* [64] uses intermediate 3D poses to estimate the 3D mesh of each person in a bottom-up fashion. ROMP [58] uses a fixed resolution body-center map to disambiguate between multiple persons under occlusion. Due to the fixed input size, 512×512 , ROMP is limited to predicting a small numbers of meshes. In contrast, OCHMR is top-down and can leverage input resizing of subject bounding box to a higher resolution for pixel accurate shape estimation. Being top-down, OCHMR can be applied to all detected persons in the input image.

3. Method

OCHMR leverages the strengths of both top-down and bottom-up methods for multi-person mesh recovery under severe person-person occlusion/crowding. In this section, we briefly describe the top-down method used as baseline architecture in our approach. Then we provide details of our contextual representations i.e. local and global centermaps and the context estimation network. Finally, we describe the proposed architectural improvements in the form of Context Normalization (CoNorm) blocks and multi-person losses used in training.

Top-down Human Mesh Recovery. Top-down human mesh recovery aims to predict a 3D human mesh from an input image $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$. Most top-down methods transform this problem to estimating the parameters of a human body model like SMPL [43] and the camera parameters. We represent body pose, shape and camera parameters by $\Theta = [\theta_{\text{pose}}, \theta_{\text{shape}}, \theta_{\text{camera}}]$, $\theta_{\text{pose}} \in \mathbb{R}^{24 \times 6}$, $\theta_{\text{shape}} \in \mathbb{R}^{10}$, $\theta_{\text{camera}} \in \mathbb{R}^3$. The pose parameters θ_{pose} are the 6D representation of the joint rotations [73] of the 24 body joints and include the global root orientation of the SMPL body. The shape parameters θ_{shape} represent the first 10 coefficients of the PCA shape space. The camera parameters θ_{camera} describe the 2D scale s and translation $\mathbf{t} = (t_x, t_y)$. SMPL is

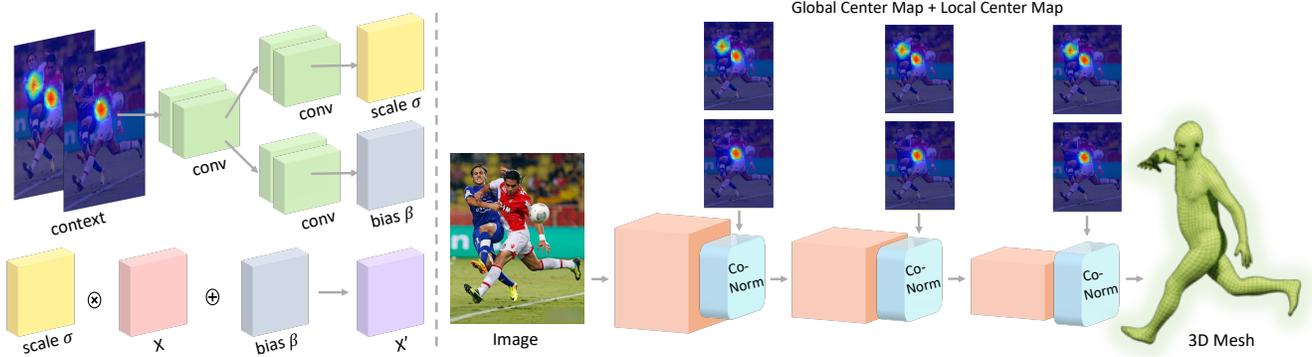


Figure 3. Context Normalization (CoNorm) Block (in blue) learns scale σ and bias β parameters for spatial affine transformation of intermediate features \mathbf{X} (in red) from the image context. *Left*: We concatenate 2D global and local center maps channel-wise to represent image context. *Right*: We insert multiple CoNorm blocks at various depths in the deep neural architecture - injection of high resolution contextual information throughout the network is critical for predicting accurate 3D meshes under occlusion.

linear and fully differentiable, making it a suitable representation for learning based methods.

Similar to [26, 32], we define a deep regression model P as our baseline top-down architecture for human mesh recovery. The bounding box at training and inference is scaled to $H \times W$ and is provided as an input to P . Let Θ^{gt} denote the ground-truth SMPL and camera parameters corresponding to the human in the input image \mathbf{I} . The deep regression model P transforms input \mathbf{I} to a single 3D mesh M , such that $\Theta = P(\mathbf{I})$. P is trained to minimize the sum of various 2D/3D pose and shape losses (using 2D pose annotations and segmentations masks if available) denoted by $\mathcal{L}(\Theta^{\text{gt}}, \Theta)$ [32].

3.1. Occluded Human Mesh Recovery

We propose to modify the top-down deep regression model P to predict multiple meshes as follows. Let N be the number of ground-truth subjects present in the image \mathbf{I} . N is set to the total number of subjects with atleast 5 visible 2D keypoints in the image. Let $\Theta_0^{\text{gt}}, \Theta_1^{\text{gt}}, \dots, \Theta_{N-1}^{\text{gt}}$ be the corresponding ground-truth mesh parameters. Our modified deep regression model P predicts N instances, $\Theta_0, \Theta_1, \dots, \Theta_{N-1}$ for an input \mathbf{I} . This is achieved by conditioning the network P on the *spatial-context* \mathbf{C}_i individually for each subject. P accepts both \mathbf{I} and \mathbf{C}_i as input and predicts $\Theta_i = P(\mathbf{I}, \mathbf{C}_i)$ where $i \in \{0, 1, \dots, N-1\}$. We define the OCHMR’s single person loss $\mathcal{L}_{\text{single}}$ as follows,

$$\mathcal{L}_{\text{single}} = \frac{1}{N} \sum_{i=0}^{N-1} \mathcal{L}(\Theta_i^{\text{gt}}, \Theta_i) \quad (1)$$

During inference, we vary the spatial-context \mathbf{C}_i to extract multiple mesh predictions from the same input image \mathbf{I} . In cases of severely overlapping people, it is hard for the baseline top-down method to estimate diverse body meshes M from similar image patches \mathbf{I} . OCHMR uses spatial-

context \mathbf{C} to resolve the implicit ambiguity of the bounding-box input representation in such multi-person cases.

3.2. Global and Local Center Map Estimation

Our top-down framework relies heavily on the representation of the spatial-context \mathbf{C} . It is crucial to define a representation which is explicit and robust to occlusion. Inspired by [8, 58], we choose body centers to encode the spatial context \mathbf{C} of the image. Specifically, we represent the contextual information of i^{th} instance as $\mathbf{C}_i = (\mathbf{C}_{\text{global}}, \mathbf{C}_{\text{local-}i})$ where $\mathbf{C}_{\text{global}}$ is the body-center heatmap of all the N instances present in the image \mathbf{I} and $\mathbf{C}_{\text{local-}i}$ is the body-center heatmap of the i^{th} instance (see Fig. 2). $\mathbf{C}_{\text{local-}i}$ is calculated by thresholding and iterating over pixel locations in $\mathbf{C}_{\text{global}}$. While $\mathbf{C}_{\text{local-}i}$ informs the network about the subject of interest, $\mathbf{C}_{\text{global}}$ places the subject in the context of its neighbors, thereby helping the network disambiguate between occluding persons.

The body center is defined as the center of visible torso joints (neck, left/right shoulders, pelvis, and left/right hips). When all torso joints are invisible, the center is the average of the visible joints. Following [58], we calculate the ground-truth body center from the ground-truth 2D pose. All the ground-truth 2D body center locations are converted into $\mathbf{C}_{\text{global}}^{\text{gt}}$ which is a heatmap of size $H \times W$ indicating the probability of the body centers at any spatial location [57]. At inference, we use a fully-convolutional [42] neural network F to predict $\mathbf{C}_{\text{global}}$ from the input image \mathbf{I} . F is trained to minimize the mean squared loss $\mathcal{L}_{\text{context}} = \text{MSE}(\mathbf{C}_{\text{global}}^{\text{gt}}, \mathbf{C}_{\text{global}})$. Finally, the context of the i^{th} instance \mathbf{C}_i is the channel-wise concatenation of $\mathbf{C}_{\text{global}}$ and $\mathbf{C}_{\text{local-}i}$ *i.e.* $\mathbf{C}_i \in \mathbb{R}^{H \times W \times 2}$.

3.3. Context Normalization Block

A key challenge is to design an architecture that incorporates spatial-context as a conditioning input. A naïve *early fusion* approach would be to simply concatenate the input im-

age \mathbf{I} with the spatial-context \mathbf{C} . Similarly, *late fusion* would concatenate feature maps from later layers within the network with appropriately down-sampled context \mathbf{C} . However, both of these approaches fail to improve performance.

We describe the Context Normalization (CoNorm) block that can be easily introduced in any existing feature extraction backbone to overcome this issue (see Fig. 3). The key intuition is that CoNorm allows normalization of intermediate feature maps using the conditioning input \mathbf{C} . The deep regression model P uses CoNorm blocks to leverage contextual information for predicting multiple meshes from the input image \mathbf{I} . Similar to Batch Normalization [17], CoNorm learns to influence the output of the neural network by applying an affine transformation to the network’s intermediate features based on \mathbf{C} .

Let $\mathbf{X} \in \mathbb{R}^{H' \times W' \times D}$ be an intermediate feature in the deep network P . The CoNorm block consists of operations Φ_{latent} , Φ_{scale} and Φ_{bias} on the context \mathbf{C} . \mathbf{C} is spatially downsampled to the same 2D resolution $H' \times W'$ as \mathbf{X} .

$$\boldsymbol{\lambda} = \Phi_{\text{latent}}(\mathbf{C}), \quad (2)$$

$$\boldsymbol{\sigma} = \Phi_{\text{scale}}(\boldsymbol{\lambda}), \quad (3)$$

$$\boldsymbol{\beta} = \Phi_{\text{bias}}(\boldsymbol{\lambda}), \quad (4)$$

$$\mathbf{X}' = \boldsymbol{\sigma} * \mathbf{X} + \boldsymbol{\beta}. \quad (5)$$

Φ_{latent} maps \mathbf{C}_i to $\boldsymbol{\lambda}$ which is in a V dimensional latent space *i.e.* $\boldsymbol{\lambda} \in \mathbb{R}^{H' \times W' \times V}$. Φ_{scale} and Φ_{bias} use the latent vector $\boldsymbol{\lambda}$ to predict $\boldsymbol{\sigma}$ and $\boldsymbol{\beta}$ respectively. $\boldsymbol{\sigma}, \boldsymbol{\beta} \in \mathbb{R}^{H' \times W' \times D}$. We use the predicted $\boldsymbol{\sigma}$ and $\boldsymbol{\beta}$ to normalize the intermediate feature \mathbf{X} using element-wise operations to output \mathbf{X}' .

3.4. Multi-Person Losses

In multi-person scenarios, the regression model P can often predict meshes that are intersecting and have incoherent depth ordering. Following [19], we adopt two multi-person losses – i) interpenetration and ii) depth-ordering loss, refer Fig. 4. We briefly describe the losses here for completeness but refer to [19] for more details.

Interpenetration Loss. Let Ω be the modified Signed Distance Field (SDF) [13] over the 3D space. Ω takes a

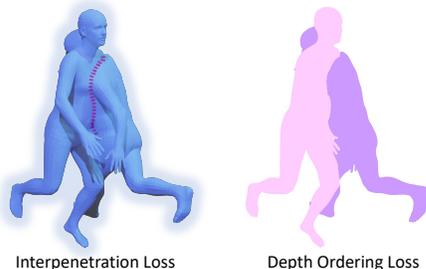


Figure 4. Interpenetration loss prevents mesh intersections. Depth ordering loss is useful for depth-consistent mesh outputs.

positive value for all the points inside the 3D human mesh \mathbf{M} , proportional to the distance from the mesh surface and is 0 everywhere else. We compute a separate distance field Ω_i for each human mesh $\mathbf{M}_i \in \{0, 1, \dots, N\}$ in the image \mathbf{I} . We define the pairwise interpenetration loss $\mathcal{L}_{\text{collision}}^{ij}$ between mesh \mathbf{M}_i and mesh \mathbf{M}_j as follows,

$$\Omega(x, y, z) = -\min(\text{SDF}(x, y, z), 0), \quad (6)$$

$$\mathcal{L}_{\text{collision}}^{ij} = \sum_{v \in \mathbf{M}_j} \Omega_i(v), \quad (7)$$

$$\mathcal{L}_{\text{collision}} = \sum_{i=1}^N \sum_{\substack{j=1 \\ i \neq j}}^N \mathcal{L}_{\text{collision}}^{ij}. \quad (8)$$

$\mathcal{L}_{\text{collision}}$ is the sum of valid pairwise mesh collisions (Fig.4).

Depth-ordering Loss. We now define the depth-ordering loss $\mathcal{L}_{\text{depth}}$. The key idea is to leverage the ground-truth instance segmentation maps available in the COCO datasets [39]. We render all the meshes and the corresponding depth maps onto the image plane using a differentiable renderer [27] and optimize the vertex locations based on the agreement with the ground-truth instance segmentation map of the image \mathbf{I} (Fig.4).

Finally, we train the network P to minimize the loss \mathcal{L} where w_{single} , $w_{\text{collision}}$ and w_{depth} are loss weights,

$$\mathcal{L} = w_{\text{single}} \mathcal{L}_{\text{single}} + w_{\text{collision}} \mathcal{L}_{\text{collision}} + w_{\text{depth}} \mathcal{L}_{\text{depth}} \quad (9)$$

4. Experiments

4.1. Implementation Details

OCHMR. For a fair comparison with other approaches [32, 58], we use ResNet-50 [15] as the default backbone for the mesh regression model P and HRNet-W32 [57] as the backbone for the context estimator F . We insert CoNorm blocks after each of the 4 ResNet block in the backbone. We set the CoNorm’s latent space dimensionality K as 128 for all our experiments. The input images are resized to 224×224 , keeping the same aspect ratio and padding with zeros. Following [57], gaussians of size 6 pixels is used to generate the local/global centermaps. The train-time data-augmentation, training schedule and all other hyper-parameters are set similar to [32]. The loss weights are set to $w_{\text{single}} = 1$, $w_{\text{collision}} = 0.2$, $w_{\text{depth}} = 0.4$ to ensure that the weighted loss items are of the same magnitude. The threshold of the local/global center heatmaps is set to 0.3.

Training Datasets. Similar to [32], we use MPI-INF-3DHP [46], COCO [39], MPII [1], LSP-Extended [20] for training (we do not use Human3.6M [18] due to licensing issues). Only the training sets are used, following the standard split protocols. We use ground-truth SMPL annotations from MPI-INF-3DHP and 2D annotations from COCO, MPII and LSP-Extended. The instance segmentation masks from COCO are used to compute $\mathcal{L}_{\text{depth}}$.

Method	Extra Data	3DPW-PC ↓			AP	OCHuman ↑				CrowdPose ↑		
		MPJPE	PMPJPE	PVE		AP ⁵⁰	AP ⁷⁵	AP ^M	AP ^L	AP	AP ⁵⁰	AP ⁷⁵
SPIN [32]	✓	129.6	82.6	157.6	12.7	46.8	19.4	17.8	26.2	16.4	40.1	10.6
PyMaf [67]	✓	126.7	81.3	154.3	14.3	48.7	21.5	18.0	28.7	17.4	42.7	13.0
ROMP* [58]	✓	115.6	75.8	147.5	19.8	56.2	25.0	19.3	32.9	28.5	58.8	24.7
SPIN [32]	✗	132.7	83.7	162.3	11.1	41.4	18.6	15.6	25.9	14.8	38.5	9.5
ROMP [58]	✗	119.7	79.7	152.8	15.6	55.0	23.6	18.7	30.0	18.9	44.6	13.8
OCHMR (Ours)	✗	117.5 (-2.2)	77.1	149.6	24.8 (+9.2)	60.7	28.6	22.3	34.2	21.4 (+2.5)	48.3	16.5
Using ground-truth bounding boxes												
SPIN [32]	✗	128.4	82.1	155.7	16.9	56.1	25.4	20.0	31.4	17.2	42.4	11.2
OCHMR (Ours)	✗	112.2 (-16.2)	75.2	145.9	37.7 (+20.8)	76.4	33.0	25.0	37.7	23.6 (+6.4)	51.1	18.7

Table 1. Comparisons to the state-of-the-art methods under severe occlusion using FasterRCNN [52] and ground-truth bounding boxes. OCHMR significantly outperforms top-down as well as bottom-up approaches across all benchmarks. ROMP* trains on CrowdPose.

Evaluation Benchmarks. 3DPW-PC [70] is employed as the main benchmark for evaluating 3D mesh/joint error since it contains in-the-wild multi-person videos with abundant 2D/3D annotations. 3DPW-PC is the *person-occluded* subset of 3DPW [61]. We also evaluate OCHMR under severe occlusion on Crowdpose [34] and OCHuman [69] which are crowded-in-the-wild 2D pose benchmarks. For completeness, we also benchmark our approach on the general datasets like 3DPW and COCO.

Evaluation Metrics. We report mean per joint position error (MPJPE), Procrustes-aligned MPJPE (PMPJPE) and per-vertex error (PVE) on the 3D datasets. MPJPE and PMPJPE evaluates the 3D joint rotation accuracy and PVE evaluates the 3D surface error. Also, to evaluate the pose accuracy under occlusion, we report standard metrics such as AP, AP⁵⁰, AP⁷⁵, AP^M, AP^L, AR at various Object Keypoint Similarity [34, 39]. We also report results using bounding boxes obtained via Faster R-CNN [52] detector.

4.2. Comparison to the State-of-the-Art

Occlusion benchmarks. To validate the stability under occlusion, we evaluate OCHMR on multiple occlusion benchmarks. Firstly, on the person-occluded 3DPW-PC, OCHuman and Crowdpose, results in Tab. 1 show that OCHMR significantly outperforms previous state-of-the-art methods [32, 58, 67]. Additionally, in Fig. 5, we qualitatively demonstrate the robustness of OCHMR under severe occlusion in comparison to top-down SPIN [32] and bottom-up ROMP [58]. Further, when using ground-truth bounding boxes, the gains of OCHMR are significant in comparison to baselines. These results show that using high-resolution input images along with global/local centermaps is key for occlusion reasoning.

General benchmarks. We also compare OCHMR with other approaches on general benchmarks like 3DPW (Tab. 2) and COCO (Tab. 3). OCHMR undergoes no performance degradation on non-occlusion cases. In fact, OCHMR improves baseline SPIN’s MPJPE error by 5mm on 3DPW.

Without using extra supervision, our method achieves comparable performance to ROMP with ResNet-50 backbone. We also outperform other methods on the COCO dataset.

Method	H3.6M	MPJPE ↓	PMPJPE ↓	PVE ↓
HMR [25]	✓	130.0	76.7	-
Kanazawa et al. [26]	✓	116.5	72.6	139.3
Arnab et al. [2]	✓	-	72.2	-
GCMR [33]	✓	-	70.2	-
DSD-SATN [59]	✓	-	69.5	-
SPIN [32]	✓	96.9	59.2	116.4
ROMP (ResNet-50) [58]	✓	91.3	54.9	108.3
I2L-MeshNet* [47]	✓	93.2	58.6	-
EFT* [22]	✓	-	54.2	-
VIBE* [30]	✓	93.5	56.5	113.4
PyMaf* [67]	✓	92.8	58.9	110.1
ROMP (ResNet-50)* [58]	✓	89.3	53.5	105.6
SPIN [32]	✗	94.7	60.2	111.4
OCHMR (Ours)	✗	89.7 (-5.0)	58.3 (-1.9)	107.1 (-4.3)

Table 2. Comparisons to the state-of-the-art methods on 3DPW test set using Protocol 2 [58]. * denotes extra training data in comparison to SPIN [32]. OCHMR does not use Human3.6M [18] for training and achieves comparable results to prior art that uses extra supervision.

Method	AP ↑	AP ⁵⁰ ↑	AP ⁷⁵ ↑	AP ^M ↑	AP ^L ↑	AR ↑
SPIN [32]	11.3	28.6	5.8	10.2	11.4	22.8
CRMH* [19]	12.6	33.8	7.6	13.2	12.8	25.0
PyMaf* [67]	13.8	35.8	9.7	14.8	14.2	28.9
ROMP [58]	14.7	36.7	9.8	15.3	14.8	29.0
OCHMR (ours)	15.3 (+0.6)	38.7	10.2	16.7	15.9	29.4
Using ground-truth bounding boxes						
SPIN [32]	13.0	33.8	7.0	13.6	12.9	26.8
OCHMR (Ours)	17.4 (+4.4)	41.9	11.8	18.2	17.4	32.4

Table 3. Comparisons to the state-of-the-art methods on COCO val set evaluated for 2D keypoint projection. * denotes extra training data compared to OCHMR.

4.3. Analysis

We perform all our analysis on the 3DPW-PC dataset with ground-truth boxes for evaluations.

CoNorm Block Architecture. We compare CoNorm blocks against *early* and *late* fusion in Tab. 4. In *early* fusion, we perform channel-wise concatenation of input image, global centermap and local centermap. In *late* fusion we concatenate the intermediate feature after the third ResNet block with the downsampled context information. We observe that injection of high-resolution context information at multiple-depths in the form of CoNorm blocks is important for accurate human mesh recovery under occlusion. Further, we vary the dimension K of the latent space of the four CoNorm blocks in the OCHMR backbone. We show that increasing K improves performance under occlusion in comparison to baseline SPIN, with $K = 128$ achieving the optimal balance between the parameter overhead and human recovery performance.

Effect of Multi-Person losses. To understand the effect of multi-person losses like interpenetration loss $\mathcal{L}_{\text{collision}}$ and depth-ordering loss $\mathcal{L}_{\text{depth}}$, we perform an ablative study using the loss weights in the OCHMR framework in Tab. 5. We achieve the best performance when using both losses, however the use of supervised $\mathcal{L}_{\text{depth}}$ loss gives better gains than self-supervised $\mathcal{L}_{\text{collision}}$. Note, OCHMR still significantly outperforms baseline SPIN when only using $\mathcal{L}_{\text{single}}$.

Choice of Context. CoNorm blocks allow conditioning the network P with various representations of the spatial-context C . Tab. 6 shows the effect of using ground-truth and predicted (using F) Local and Local + Global Centermaps along with 2D keypoints as C . We use offshelf pose-estimation network HRNet-W48 [57] trained on COCO dataset as our F . In case of 2D keypoints, C is a 17-channel heatmap corresponding to keypoint locations. In comparison to local centermaps, the addition of global centermaps helps improve performance under occlusion. Interestingly, conditioning us-

Method	MPJPE ↓	PMPJPE ↓	PVE ↓
SPIN	128.4	82.1	155.7
OCHMR, <i>early-fusion</i>	115.8	76.4	150.1
OCHMR, <i>late-fusion</i>	119.8	80.2	151.8
OCHMR, $K = 16$	116.8	76.9	150.2
OCHMR, $K = 32$	114.2	76.2	148.6
OCHMR, $K = 64$	113.0	75.0	146.4
OCHMR, $K = 128$	112.2	75.2	145.9
OCHMR, $K = 256$	113.1	74.7	146.1

Table 4. Comparison of CoNorm block with *early* and *late* fusion of context along with variation of CoNorm block’s latent space dimension K . Injection of contextual information at multiple depths outperforms *early/late* fusion. Increase in K results in better context normalization with better performance under occlusion.

$\mathcal{L}_{\text{single}}$	$\mathcal{L}_{\text{collision}}$	$\mathcal{L}_{\text{depth}}$	MPJPE ↓	PMPJPE ↓	PVE ↓
✓	✗	✗	116.9	77.1	149.2
✓	✓	✗	115.3	76.2	148.7
✓	✗	✓	113.6	75.6	147.0
✓	✓	✓	112.2	75.2	145.9

Table 5. Ablation of multi-person losses in OCHMR. We default w_{single} to 1 to ensure model convergence. We found the relative importance of $\mathcal{L}_{\text{depth}}$ to be greater than $\mathcal{L}_{\text{collision}}$.

ing ground-truth 2D keypoints outperforms all other choices. However, when ground-truth keypoints are unavailable, body centers outperform estimated 2D keypoints as estimating accurate 2D pose under occlusion is more challenging than estimating body centers.

Context C	Ground-Truth		Estimated by F	
	MPJPE ↓	PMPJPE ↓	MPJPE ↓	PMPJPE ↓
Local Center	113.0	76.4	114.8	77.1
Local + Global Center	111.4	74.7	112.2	75.2
2D Keypoints	109.5	73.9	116.8	78.9

Table 6. Comparison of various choices of contexts in OCHMR conditioning. Local + Global centermaps performs better than other conditioning choices when being estimated by the network F .

Limitations. OCHMR is a multi-stage top-down method and is, therefore, not real-time during inference. Though OCHMR improves performance under multi-person occlusion, it is still susceptible to failure under truncation and extreme cropping due to object occlusion. Moreover, OCHMR fails to handle extreme poses and shapes due to the lack of training data, as shown in the Sup. Mat. In the future, OCHMR can be extended and incorporated with the recent progress to handle various kinds of occlusions [31, 49, 70].

5. Conclusion

Most top-down methods for human mesh recovery assume a single subject in the input, causing them to fail under severe person-person occlusion. In this work, we introduce OCHMR, a novel top-down method to handle multiple occluded people in crowded scenes. Our key idea is conditioning top-down models on spatial-context from the image, in the form of local and global centermaps which allows OCHMR to effectively disambiguate between overlapping humans. We propose Contextual Normalization (CoNorm) blocks, a novel architectural improvement which can be easily extended to any existing top-down method. While OCHMR draws inspiration from bottom-up methods, we retain the advantages of both bottom-up and top-down methods, resulting in a method that can handle multi-person occlusion and achieve pixel-aligned reconstruction results.

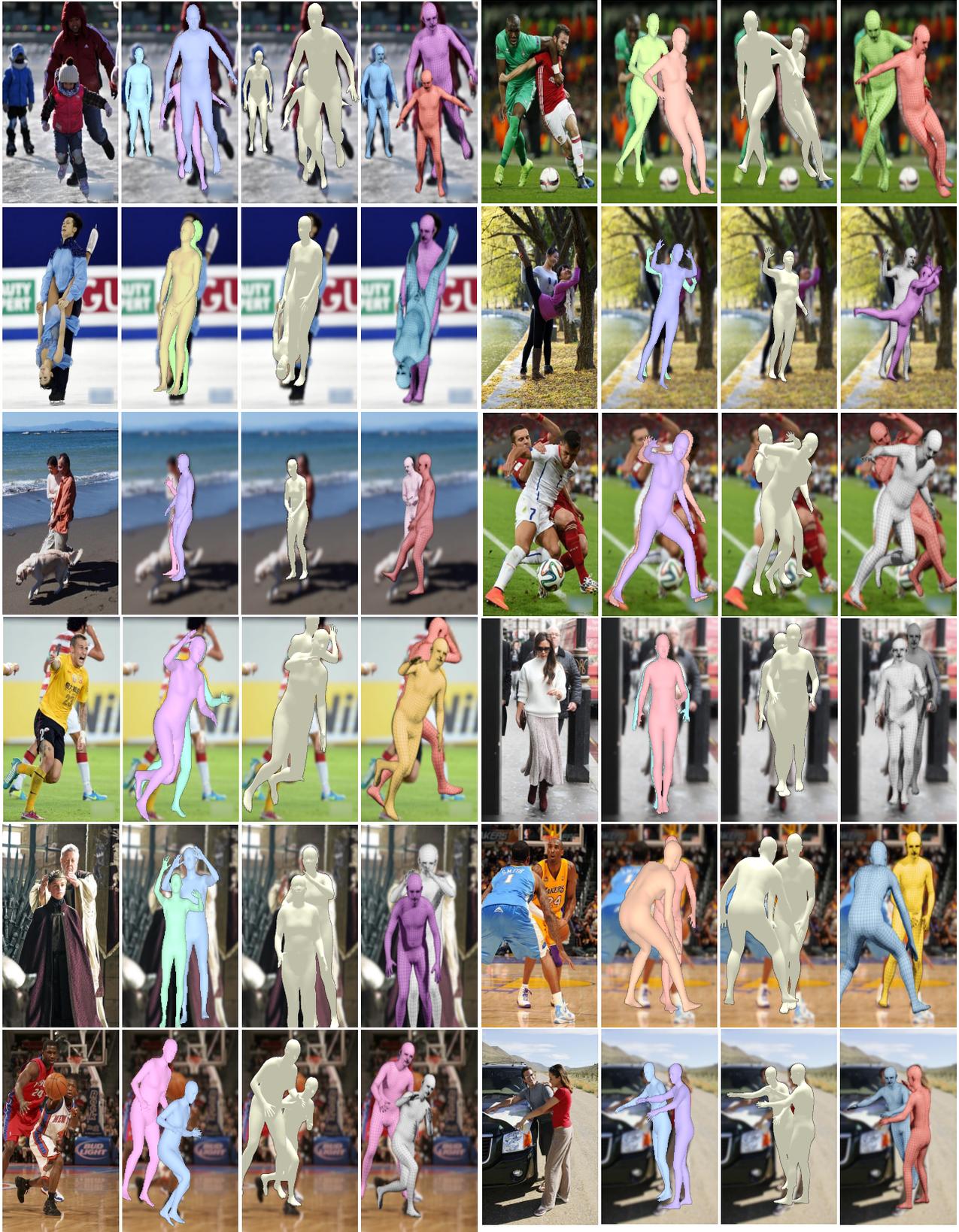


Figure 5. Qualitative results on the OCHuman val set. Each image (left to right) shows RGB image, SPIN [32] predictions, ROMP [58] predictions and OCHMR predictions. Due to occlusions, SPIN often misses the person in the background which is recovered by OCHMR. In comparison to ROMP, OCHMR outputs pixel aligned meshes with correct depth-ordering. Please see additional results in Sup. Mat.

References

- [1] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3686–3693, 2014. 3, 5
- [2] Anurag Arnab, Carl Doersch, and Andrew Zisserman. Exploiting temporal context for 3d human pose estimation in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3395–3404, 2019. 6
- [3] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*, 2020. 3
- [4] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In *European conference on computer vision*, pages 561–578. Springer, 2016. 1, 3
- [5] Bowen Cheng, Yunchao Wei, Honghui Shi, Rogerio Feris, Jinjun Xiong, and Thomas Huang. Revisiting rcnn: On awakening the classification power of faster rcnn. In *Proceedings of the European conference on computer vision (ECCV)*, pages 453–468, 2018. 3
- [6] 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. 2, 3
- [7] Hongsuk Choi, Gyeongsik Moon, and Kyoung Mu Lee. Pose2mesh: Graph convolutional network for 3d human pose and mesh recovery from a 2d human pose. In *European Conference on Computer Vision*, pages 769–787. Springer, 2020. 1, 3
- [8] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6569–6578, 2019. 2, 4
- [9] Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. *arXiv preprint arXiv:1610.07629*, 2016. 2
- [10] Sai Kumar Dwivedi, Nikos Athanasiou, Muhammed Kocabas, and Michael J Black. Learning to regress bodies from images using differentiable semantic rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11250–11259, 2021. 1, 2, 3
- [11] Peng Guan, Alexander Weiss, Alexandru O Balan, and Michael J Black. Estimating human shape and pose from a single image. In *2009 IEEE 12th International Conference on Computer Vision*, pages 1381–1388. IEEE, 2009. 1
- [12] Shanyan Guan, Jingwei Xu, Yunbo Wang, Bingbing Ni, and Xiaokang Yang. Bilevel online adaptation for out-of-domain human mesh reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10472–10481, 2021. 1, 3
- [13] Mohamed Hassan, Vasileios Choutas, Dimitrios Tzionas, and Michael J Black. Resolving 3d human pose ambiguities with 3d scene constraints. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2282–2292, 2019. 5
- [14] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 3
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5
- [16] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1501–1510, 2017. 2
- [17] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015. 5
- [18] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence*, 36(7):1325–1339, 2013. 2, 3, 5, 6
- [19] Wen Jiang, Nikos Kolotouros, Georgios Pavlakos, Xiaowei Zhou, and Kostas Daniilidis. Coherent reconstruction of multiple humans from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5579–5588, 2020. 2, 3, 5, 6
- [20] Sam Johnson and Mark Everingham. Learning effective human pose estimation from inaccurate annotation. In *CVPR 2011*, pages 1465–1472. IEEE, 2011. 3, 5
- [21] 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 *Proceedings of the IEEE International Conference on Computer Vision*, pages 3334–3342, 2015. 2, 3
- [22] Hanbyul Joo, Natalia Neverova, and Andrea Vedaldi. Exemplar fine-tuning for 3d human model fitting towards in-the-wild 3d human pose estimation. *arXiv preprint arXiv:2004.03686*, 2020. 3, 6
- [23] Hanbyul Joo, Tomas Simon, and Yaser Sheikh. Total capture: A 3d deformation model for tracking faces, hands, and bodies. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8320–8329, 2018. 2, 3
- [24] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. corr abs/1712.06584 (2017). *arXiv preprint arXiv:1712.06584*, 2017. 1, 3
- [25] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7122–7131, 2018. 1, 6
- [26] Angjoo Kanazawa, Jason Y Zhang, Panna Felsen, and Jitendra Malik. Learning 3d human dynamics from video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5614–5623, 2019. 4, 6
- [27] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neu-

- ral 3d mesh renderer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3907–3916, 2018. 5
- [28] Rawal Khirodkar, Visesh Chari, Amit Agrawal, and Ambrish Tyagi. Multi-hypothesis pose networks: Rethinking top-down pose estimation. *arXiv preprint arXiv:2101.11223*, 2021. 3
- [29] Imry Kissos, Lior Fritz, Matan Goldman, Omer Meir, Eduard Oks, and Mark Kliger. Beyond weak perspective for monocular 3d human pose estimation. In *European Conference on Computer Vision*, pages 541–554. Springer, 2020. 3
- [30] Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. Vibe: Video inference for human body pose and shape estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5253–5263, 2020. 3, 6
- [31] Muhammed Kocabas, Chun-Hao P Huang, Otmar Hilliges, and Michael J Black. Pare: Part attention regressor for 3d human body estimation. *arXiv preprint arXiv:2104.08527*, 2021. 1, 3, 7
- [32] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2252–2261, 2019. 1, 2, 3, 4, 5, 6, 8
- [33] Nikos Kolotouros, Georgios Pavlakos, and Kostas Daniilidis. Convolutional mesh regression for single-image human shape reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4501–4510, 2019. 6
- [34] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10863–10872, 2019. 2, 3, 6
- [35] Jiefeng Li, Chao Xu, Zhicun Chen, Siyuan Bian, Lixin Yang, and Cewu Lu. Hybrik: A hybrid analytical-neural inverse kinematics solution for 3d human pose and shape estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3383–3393, 2021. 1, 3
- [36] Kevin Lin, Lijuan Wang, and Zicheng Liu. End-to-end human pose and mesh reconstruction with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1954–1963, 2021. 3
- [37] Kevin Lin, Lijuan Wang, and Zicheng Liu. Mesh graphormer. In *ICCV*, 2021. 1, 3
- [38] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 3
- [39] 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 *European conference on computer vision*, pages 740–755. Springer, 2014. 3, 5, 6
- [40] HuaJun Liu, Fuqiang Liu, Xinyi Fan, and Dong Huang. Polarized self-attention: Towards high-quality pixel-wise regression. *arXiv preprint arXiv:2107.00782*, 2021. 3
- [41] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016. 3
- [42] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 4
- [43] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):1–16, 2015. 1, 3
- [44] Zhiheng Ma, Xiaopeng Hong, Xing Wei, Yunfeng Qiu, and Yihong Gong. Towards a universal model for cross-dataset crowd counting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3205–3214, 2021. 2
- [45] William McNally, Kanav Vats, Alexander Wong, and John McPhee. Evopose2d: Pushing the boundaries of 2d human pose estimation using accelerated neuroevolution with weight transfer. *IEEE Access*, 2021. 3
- [46] 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 *2017 international conference on 3D vision (3DV)*, pages 506–516. IEEE, 2017. 2, 3, 5
- [47] Gyeongsik Moon and Kyoung Mu Lee. I2l-meshnet: Image-to-lixel prediction network for accurate 3d human pose and mesh estimation from a single rgb image. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VII 16*, pages 752–768. Springer, 2020. 3, 6
- [48] Mohamed Omran, Christoph Lassner, Gerard Pons-Moll, Peter Gehler, and Bernt Schiele. Neural body fitting: Unifying deep learning and model based human pose and shape estimation. In *2018 international conference on 3D vision (3DV)*, pages 484–494. IEEE, 2018. 1
- [49] Georgios Pavlakos, Jitendra Malik, and Angjoo Kanazawa. Human mesh recovery from multiple shots. *arXiv preprint arXiv:2012.09843*, 2020. 7
- [50] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3d human pose and shape from a single color image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 459–468, 2018. 1
- [51] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018. 2
- [52] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28:91–99, 2015. 3, 6
- [53] Matteo Ruggero Ronchi and Pietro Perona. Benchmarking and error diagnosis in multi-instance pose estimation. In *Proceedings of the IEEE international conference on computer vision*, pages 369–378, 2017. 3
- [54] Akash Sengupta, Ignas Budvytis, and Roberto Cipolla. Synthetic training for accurate 3d human pose and shape estimation in the wild. *arXiv preprint arXiv:2009.10013*, 2020.

- 3
- [55] Leonid Sigal, Alexandru O Balan, and Michael J Black. Humaneva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion. *International journal of computer vision*, 87(1-2):4, 2010. 2, 3
- [56] Qingyu Song, Changan Wang, Zhengkai Jiang, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, and Yang Wu. Rethinking counting and localization in crowds: A purely point-based framework. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3365–3374, 2021. 2
- [57] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5693–5703, 2019. 3, 4, 5, 7
- [58] Yu Sun, Qian Bao, Wu Liu, Yili Fu, Michael J Black, and Tao Mei. Monocular, one-stage, regression of multiple 3d people. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11179–11188, 2021. 1, 2, 3, 4, 5, 6, 8
- [59] Yu Sun, Yun Ye, Wu Liu, Wenpeng Gao, Yili Fu, and Tao Mei. Human mesh recovery from monocular images via a skeleton-disentangled representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5349–5358, 2019. 6
- [60] Shashank Tripathi, Siddhant Ranade, Amrith Tyagi, and Amit Agrawal. Posenet3d: Learning temporally consistent 3d human pose via knowledge distillation. In *2020 International Conference on 3D Vision (3DV)*, pages 311–321. IEEE, 2020. 3
- [61] Timo von Marcard, Roberto Henschel, Michael J Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 601–617, 2018. 2, 3, 6
- [62] Changan Wang, Qingyu Song, Boshen Zhang, Yabiao Wang, Ying Tai, Xuyi Hu, Chengjie Wang, Jilin Li, Jiayi Ma, and Yang Wu. Uniformity in heterogeneity: Diving deep into count interval partition for crowd counting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3234–3242, 2021. 2
- [63] Hongwei Yi, Chun-Hao P. Huang, Dimitrios Tzionas, Muhammed Kocabas, Mohamed Hassan, Siyu Tang, Justus Thies, and Michael J. Black. Human-aware object placement for visual environment reconstruction. In *Computer Vision and Pattern Recognition (CVPR)*, June 2022. 3
- [64] Andrei Zanfir, Elisabeta Marinoiu, Mihai Zanfir, Alin-Ionut Popa, and Cristian Sminchisescu. Deep network for the integrated 3d sensing of multiple people in natural images. *Advances in Neural Information Processing Systems*, 31:8410–8419, 2018. 2, 3
- [65] Feng Zhang, Xiatian Zhu, Hanbin Dai, Mao Ye, and Ce Zhu. Distribution-aware coordinate representation for human pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7093–7102, 2020. 3
- [66] Hongwen Zhang, Jie Cao, Guo Lu, Wanli Ouyang, and Zhenan Sun. Danet: Decompose-and-aggregate network for 3d human shape and pose estimation. In *Proceedings of the 27th ACM International Conference on Multimedia*, pages 935–944, 2019. 1
- [67] Hongwen Zhang, Yating Tian, Xinchu Zhou, Wanli Ouyang, Yebin Liu, Limin Wang, and Zhenan Sun. Pymaf: 3d human pose and shape regression with pyramidal mesh alignment feedback loop. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11446–11456, 2021. 1, 2, 3, 6
- [68] Jianfeng Zhang, Dongdong Yu, Jun Hao Liew, Xuecheng Nie, and Jiashi Feng. Body meshes as points. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 546–556, 2021. 2
- [69] Song-Hai Zhang, Ruilong Li, Xin Dong, Paul Rosin, Zixi Cai, Xi Han, Dingcheng Yang, Haozhi Huang, and Shi-Min Hu. Pose2seg: Detection free human instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 889–898, 2019. 2, 3, 6
- [70] Tianshu Zhang, Buzhen Huang, and Yangang Wang. Object-occluded human shape and pose estimation from a single color image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7376–7385, 2020. 2, 3, 6, 7
- [71] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019. 2
- [72] Xingyi Zhou, Jiacheng Zhuo, and Philipp Krahenbuhl. Bottom-up object detection by grouping extreme and center points. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 850–859, 2019. 2
- [73] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5745–5753, 2019. 3