

BoIR: Box-Supervised Instance Representation for Multi-Person Pose Estimation

Uyoung Jeong¹
jeong.uyoung@unist.ac.kr

Seungryul Baek¹
srbaek@unist.ac.kr

Hyung Jin Chang²
h.j.chang@bham.ac.uk

Kwang In Kim³
kjmkin@postech.ac.kr

¹ Ulsan National Institute of Science and Technology
Ulsan, Republic of Korea

² University of Birmingham
Birmingham, United Kingdom

³ Pohang University of Science and Technology
Pohang, Republic of Korea

Abstract

Single-stage multi-person human pose estimation (MPPE) methods have shown great performance improvements, but existing methods fail to disentangle features by individual instances under crowded scenes. In this paper, we propose a bounding box-level instance representation learning called BoIR, which simultaneously solves instance detection, instance disentanglement, and instance-keypoint association problems. Our new instance embedding loss provides a learning signal on the entire area of the image with bounding box annotations, achieving globally consistent and disentangled instance representation. Our method exploits multi-task learning of bottom-up keypoint estimation, bounding box regression, and contrastive instance embedding learning, without additional computational cost during inference. BoIR is effective for crowded scenes, outperforming state-of-the-art on COCO val (0.8 AP), COCO test-dev (0.5 AP), CrowdPose (4.9 AP), and OCHuman (3.5 AP). Code will be available at <https://github.com/uyoung-jeong/BoIR>

1 Introduction

Multi-person human pose estimation (MPPE) localizes 2D keypoint locations of multiple human instances from an image. It is useful not only for 3D pose estimation and activity recognition [4], but also for human-robot interaction [5], autonomous driving [43], augmented/virtual reality and surveillance applications. In wild scenarios, where severe inter-person occlusion and background clutter frequently occur, the capability of multi-person pose estimation becomes even more crucial.

Recent advances in single-stage MPPE methods [20, 52, 40] have shown significant performance improvements. Compared to top-down methods [11, 16, 57], they do not require off-the-shelf person detectors and therefore robust to detection errors. Unlike bottom-up

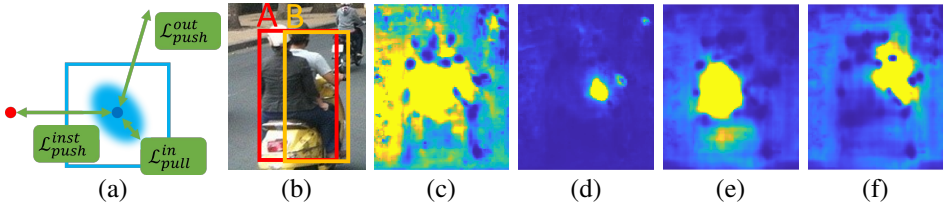


Figure 1: (a): Bbox Mask Loss framework. Blue dot is a query box(blue colour) center, while red dot is another box center. \mathcal{L}_{pull}^{in} pulls instance center and soft-masked mean embeddings inside the box, \mathcal{L}_{push}^{out} pushes pairwise instance-background embeddings, and $\mathcal{L}_{push}^{inst}$ pushes pairwise instance embeddings. (b)-(f): Visualization of feature similarities from the center features of bounding boxes in (b). (c) and (d) are CID feature similarities from A and B centers, respectively, while (e) and (f) are BoIR feature similarities.

methods [4, 6, 19, 35, 39], they solve instance-keypoint association problems by explicitly detecting instances, usually using instance center locations.

While single-stage methods showed promising results, they still suffer from instance-keypoint association under heavy inter-person occlusion, which often results in noisy predictions. We summarize the main reasons in two aspects. First, existing representation-based methods lack multi-task supervision to learn diverse aspects of instance representation. Even if they add multiple tasks, it would incur computational overhead during inference. Second, previous works have spatially sparse supervision. Many works apply losses only on ground-truth keypoint locations, which is too sparse for the model to holistically learn the entire image region, leading to noisy and globally inconsistent results, as illustrated in Fig. 1. Although heatmap-based approaches apply Gaussian kernel to generate ground-truth keypoint heatmaps, it is still more sparse than conventional segmentation level supervision.

In this paper, we focus on an effective instance representation learning method which can provide both rich spatial and multi-task supervision. First, we reformulate the existing MPPE pipeline to apply embedding loss on a separate embedding branch, which can effectively map nonlinear features of instances while the primary task branch’s performance is not degraded. Then, we design a new contrastive learning scheme, termed Bbox Mask Loss, using bounding box(bbox) supervision. It contrasts instance embeddings on both inside and outside of the ground-truth boxes, which provides learning signals on the entire image region. Combining with box regression and bottom-up keypoint heatmap regression as auxiliary tasks, we apply multi-task learning scheme to learn effective instance representation for multiple keypoint estimation.

In summary, we introduce a novel method for instance representation learning at the box level, named BoIR. BoIR adeptly addresses the challenges of instance disentanglement and instance detection simultaneously, without incurring any additional computational costs during inference. These are achieved through the following key contributions:

- Bbox Mask Loss effectively disentangles features by instances in the embedding space using a new embedding loss with spatially rich bounding box level supervision.
- Auxiliary task heads enrich instance representation by sharing multiple aspects of the instance, while no additional computational cost is induced during inference.

BoIR excels at challenging crowded scenes, surpassing comparative methods by 0.5 AP on COCO test-dev, 4.9 AP on CrowdPose test, and 3.5 AP on OCHuman test.

2 Related Works

2D Multi-person human pose estimation (MPPE). 2D MPPE methods can be roughly classified by instance handling approaches. Top-down methods use detectors [8, 26, 27] to get person boxes and use cropped images as input. Bottom-up methods first detect keypoints and group them into instances. Single-stage methods, on the other hand, detect instances first and then regress instance-wise keypoints. Single-stage methods eliminate the need to crop an image into multiple instance-wise images, and avoid the need for keypoint grouping.

SimpleBaseline [57] and HRNet [51] are top-down methods, and are generally used as backbone networks in various works. MIPNet [10] is one of the recent top-down approaches that considers multiple instances within a box by modulating the channel dimensions to regress individual keypoints.

OpenPose [1], PersonLab [24], and PifPaf [12] share a similar idea of estimating a vector field that associates keypoints with instances. HigherHRNet [9] and its subsequent works [6, 19, 65, 69] are another class of bottom-up methods using Associative Embedding [22]. From the pixel-wise one-dimensional embedding, they assign the detected keypoints to respective instances using off-the-shelf grouping algorithm [13]. These methods tend to lack the capability of instance detection since their training losses are mainly targeted for keypoint estimation.

There are several single-stage methods based on Transformers [33]. PETR [29] avoids using Hungarian algorithm for instance grouping by randomly initializing query embeddings to regress keypoints. In contrast, ED-Pose [40] extracts query embeddings via a human detection decoder, but it requires substantial computational cost due to the massive amount of learnable parameters, which is critical for real-time pose estimation. QueryPose [68] similarly performs box and keypoint regression via query embeddings and Transformers-based decoders, and its performance on CrowdPose test is inferior to CID by 0.2 AP with the same HRNet-W48 backbone.

FCPose [20] and CID [54] are single-stage methods using an instance center map. FCPose generates instance proposals from a single-stage person detector and employs instance-wise dynamic convolution on global features. Similarly, CID estimates instance center map to detect instances, and performs channel and spatial attention between sampled feature and global features, but it does not perform box regression. CID directly applies contrastive loss on the backbone network’s output feature, which does not effectively disentangle features by instances, as discussed in SimCLR [8]. Additionally, CID’s contrastive loss is spatially sparse since it is applied solely on instance center locations. Instead, we introduce a separate embedding branch that enhances learning keypoint features, providing richer spatial and multi-task guidance. KAPAO [21] is another method that reformulates keypoint regression task as an object detection task, jointly detecting persons and keypoints.

Representation learning with distance metrics. Deep metric learning aims to learn a distance metric in the embedding space for better representation, generally composed with a pull term for closing the distance among positive samples, and a push term for differentiating between different classes. Push loss term is crucial for effective representation learning, so many works are devoted to proposing various negative sampling strategies. Contrastive loss [7], triplet loss [28], N-pair loss [60] and InfoNCE loss [23] are some of the approaches. SimCLR [8], MoCo [9], and CLIP [25] are representative works using variants of InfoNCE loss. All of these methods use cosine similarity as a similarity metric.

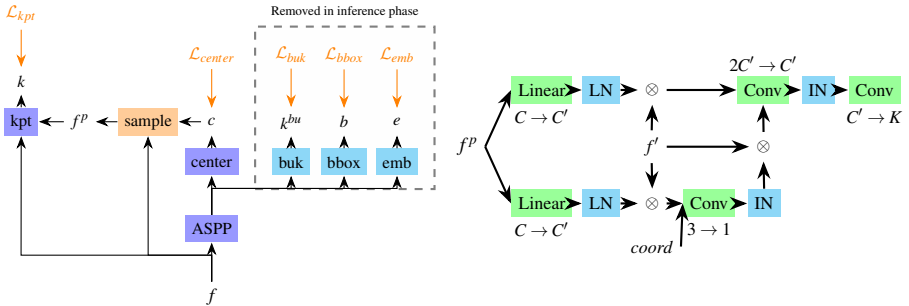


Figure 2: Left: Overview of our framework. Instance keypoint (kpt) head and center head are primary regression heads. bottom-up keypoint (buk) head, bounding box (bbox) head and embedding (emb) head are auxiliary task regressors which are not used during inference. Right: Layer composition of instance keypoint head. 'Linear': linear layer, 'Conv': convolution layer, 'LN': Layer Normalization, 'IN': Instance Normalization, ' \otimes ': Hadamard product. 'coord': relative coordinates of the heatmap pixel indices. $f^p \in \mathbb{R}^{C \times H \times W}$: projection of f by single convolution layer.

3 Method

3.1 Framework overview

Our framework comprises two main parts: auxiliary task branch and instance keypoint branch. Given an input image, backbone network outputs a feature $f \in \mathbb{R}^{C,H,W}$, where H is height and W is width. Task-specific heads produce instance center heatmaps $c \in \mathbb{R}^{1,H,W}$, box predictions $b \in \mathbb{R}^{4,H,W}$, bottom-up keypoint heatmaps $k^{bu} \in \mathbb{R}^{K,H,W}$ and instance embedding map $e \in \mathbb{R}^{D,H,W}$. During training, after detecting instances from the center map, instance features f^p are sampled from the backbone feature at the ground-truth center coordinates. f^p are used as conditions for regressing instance-wise keypoints k in the instance keypoint head, as proposed in [52]. In case of embedding branch, we sample instance embeddings p from e . During inference, f^p are sampled from predicted instance centers. We made several enhancements to the instance keypoint head, including Layer Normalization and Instance Normalization for stable learning, as illustrated in Fig. 2. Please note that b, k^{bu}, e are not estimated during inference.

3.2 Bbox Mask Loss

Existing instance representation learning methods such as Associative Embedding(AE) and CID's contrastive loss failed to handle multiple people in several aspects, often leading to noisy results. Firstly, they only compare instance embeddings with ground-truth(GT) instance locations, making it difficult to generate a push loss term when only one GT instance is present in an image. Secondly, there are unlabeled instances in training datasets, and existing works typically ignore these unlabeled instances, which induces additional noise during inference. Thirdly, the number of human instances per image in training datasets is insufficient for effective instance representation learning. For example, COCO train set has an average of 2.6 people per image, excluding labels with `iscrowd=1`. Similarly, CrowdPose trainval set has 4.2 people per image.

To alleviate aforementioned challenges, inspired by a weakly supervised instance segmentation method [46], we introduce spatially rich supervision via box annotations, termed Bbox Mask Loss. It disambiguates each instance embedding from outside of the box region, which can handle arbitrary unlabeled instances and background clutter. It applies soft masking on the inside of the box based on embedding similarity, which is effective for feature disentanglement under heavy cross-instance occlusions. Moreover, it can produce push loss term even when only a single GT instance is available in an image, serving as a simple but effective negative sampling method.

Bbox Mask Loss incorporates multitude of push and pull loss terms, including in-box pull \mathcal{L}_{pull}^{in} , out-box push \mathcal{L}_{push}^{out} , and cross-instance push $\mathcal{L}_{push}^{inst}$. First, given a GT instance and corresponding box with height h and width w , we compute pixel-wise embedding similarity between embedding map and the instance embedding as follows:

$$s_i^{(x,y)} = \psi(d(e^{(x,y)}, p_i)), \quad (x,y) \in \mathcal{B}_i, \quad (1)$$

where d is a distance metric, and ψ is an inversion operator to convert the distance to similarity with $[0,1]$ output range. From ablative experiment, as reported in Table 4, we find that L2 distance for d and Gaussian kernel for ψ outperforms cosine distance and cosine similarity. \mathcal{B}_i is a set of coordinates inside the box b_i , where $i = 1, 2, \dots, N$. As a pulling term inside the box, we want the model to produce similar embeddings on the foreground region of the same person. To realize the objective, we compare the instance center embedding with the mean instance embedding \bar{p}_i , as defined below:

$$\mathcal{L}_{pull}^{in} = \frac{1}{N} \sum_{i=1}^N d(p_i, \bar{p}_i), \quad \text{where} \quad \bar{p}_i = \frac{\sum_{(x,y) \in \mathcal{B}_i} e^{(x,y)} s_i^{(x,y)}}{\sum_{(x,y) \in \mathcal{B}_i} s_i^{(x,y)}}. \quad (2)$$

To decouple the instance embedding from the background, we define the out-box push loss using out-box mean embedding \bar{p}_i^c :

$$\mathcal{L}_{push}^{out} = \frac{1}{N} \sum_{i=1}^N \psi(d(p_i, \bar{p}_i^c)), \quad \text{where} \quad \bar{p}_i^c = \frac{\sum_{(x,y) \in \mathcal{B}_i^c} e^{(x,y)}}{|\mathcal{B}_i^c|}. \quad (3)$$

Note that \mathcal{B}_i^c is a set of coordinates outside the i th box, and \bar{p}_i^c is a mean embedding of the background except the i th box region. Lastly, cross-instance push term compares instance embeddings retrieved from ground-truths, which is the same as the existing losses.

$$\mathcal{L}_{push}^{inst} = \frac{1}{(N(N-1)/2)} \sum_{i=1}^N \sum_{j>i}^N \psi(d(p_i, p_j)). \quad (4)$$

3.3 Auxiliary tasks

In order to encourage the features to have richer and more disentangled information for MPPE, we incorporate multiple auxiliary tasks and instance representation learning in parallel. Our multi-task branch consists of shared layers and four separate regression heads, consisting of instance embedding, bottom-up keypoint, box, and instance center.

We concurrently reduce dimensionality of the backbone feature and incorporate multi-resolution shared feature representation based on ASPPv2 [2]. It resolves the problem of regressing globally consistent instance features. Original ASPPv2 module incurs heavy

| Method | Backbone | Input size | AP | AP ⁵⁰ | AP ⁷⁵ | AP ^M | AP ^L | AR |
|----------------------|----------------|------------|-------------|------------------|------------------|-----------------|-----------------|-------------|
| Top-down methods | | | | | | | | |
| SBL [37] | ResNet-152 | 384×288 | 73.7 | 91.9 | 81.1 | 70.3 | 80.0 | - |
| HRNet [62] | HRNet-W32 | 384×288 | 74.9 | 92.5 | 82.8 | 71.3 | 80.9 | - |
| Bottom-up methods | | | | | | | | |
| HrHRNet [9] | HrHRNet-W32 | 512 | 66.4 | 87.5 | 72.8 | 61.2 | 74.2 | - |
| DEKR [6] | HRNet-W32 | 512 | 67.3 | 87.9 | 74.1 | 61.5 | 76.1 | 72.4 |
| SWAHR [19] | HrHRNet-W32 | 512 | 67.9 | 88.9 | 74.5 | 62.4 | 75.5 | - |
| Single stage methods | | | | | | | | |
| FCPose [20] | ResNet-101+FPN | 800 | 65.6 | 87.9 | 72.6 | 62.1 | 72.3 | - |
| PETR [29] | ResNet-101 | 800 | 68.5 | 90.3 | 76.5 | 62.5 | 77.0 | - |
| ED-Pose [40] | ResNet-50 | 800 | 69.8 | 90.2 | 77.2 | 64.3 | 77.4 | - |
| CID [34] | HRNet-W32 | 512 | 68.9 | 89.9 | 76.0 | 63.2 | 77.7 | 74.6 |
| CID [34] | HRNet-W48 | 640 | 70.7 | 90.3 | 77.9 | 66.3 | 77.8 | 76.4 |
| BoIR | HRNet-W32 | 512 | 69.5 | 90.4 | 76.9 | 64.2 | 77.3 | 75.3 |
| BoIR | HRNet-W48 | 640 | 71.2 | 90.8 | 78.6 | 67.0 | 77.6 | 77.1 |

Table 1: Comparison with state-of-the-art methods on COCO *test-dev* set. Best scores are marked as bold for small (e.g. HRNet-W32) and large (e.g. HRNet-W48) models respectively.

computational cost when fusing multiple resolution features. We alleviate this by further squeezing the output channel size of each multi-resolution feature to 128, and then apply a fusion layer to obtain a final feature with 256 channel size. This design reduces the number of trainable parameters of ASPP by 50%. This shared bottleneck module design helps to prevent auxiliary tasks from dominating over the primary task, by restricting the amount of information flow to the auxiliary tasks.

Each regression head comprises with one residual block and one output convolution layer for sufficient capability of learning nonlinear feature transformation. In case of box regression, we adopt anchor free method [15] for efficient training. For clarity, we do not use the bbox head outputs during inference, and the box head serves as an efficient and informative auxiliary task head.

3.4 Training losses

We employ five loss functions: instance-wise keypoint heatmap loss \mathcal{L}_{kpt} , center heatmap loss \mathcal{L}_{center} , bottom-up keypoint heatmap loss \mathcal{L}_{buk} , bbox loss \mathcal{L}_{bbox} , and embedding loss \mathcal{L}_{emb} .

$$\mathcal{L} = \mathcal{L}_{kpt} + \mathcal{L}_{center} + \mathcal{L}_{buk} + \mathcal{L}_{bbox} + \mathcal{L}_{emb}. \quad (5)$$

Focal loss [14, 45] is used for \mathcal{L}_{kpt} , \mathcal{L}_{center} and \mathcal{L}_{buk} , while CIoU loss [44] is used for \mathcal{L}_{bbox} . For embedding loss, we use three loss terms as defined in Equation 2,3,4. We use AE loss for calculating respective terms:

$$\mathcal{L}_{emb} = \mathcal{L}_{pull}^{in} + \mathcal{L}_{push}^{out} + \mathcal{L}_{push}^{inst}. \quad (6)$$

| Method | Backbone | Input size | AP | AP ⁵⁰ | AP ⁷⁵ | AP ^E | AP ^M | AP ^H |
|----------------------|-------------|------------|-------------|------------------|------------------|-----------------|-----------------|-----------------|
| Top-down methods | | | | | | | | |
| SBL [57] | ResNet-101 | - | 60.8 | 81.4 | 65.7 | 71.4 | 61.2 | 51.2 |
| SPPE [16] | ResNet-101 | 320×256 | 66.0 | 84.2 | 71.5 | 75.5 | 66.3 | 57.4 |
| Bottom-up methods | | | | | | | | |
| HrHRNet [9] | HrHRNet-W48 | 640 | 65.9 | 86.4 | 70.6 | 73.3 | 66.5 | 57.9 |
| DEKR [9] | HrHRNet-W32 | 512 | 65.7 | 85.7 | 70.4 | 73.0 | 66.4 | 57.5 |
| SWAHR [19] | HrHRNet-W48 | 640 | 71.6 | 88.5 | 77.6 | 78.9 | 72.4 | 63.0 |
| Single stage methods | | | | | | | | |
| PETR [29] | Swin-L | 800 | 71.6 | 90.4 | 78.3 | 77.3 | 72.0 | 65.8 |
| ED-Pose [40] | ResNet-50 | 800 | 69.9 | 88.6 | 75.8 | 77.7 | 70.6 | 60.9 |
| CID [54] | HRNet-W32 | 512 | 71.3 | 90.6 | 76.6 | 77.4 | 72.1 | 63.9 |
| CID [54] | HRNet-W48 | 640 | 72.3 | 90.8 | 77.9 | 78.7 | 73.0 | 64.8 |
| CID* [54] | HRNet-W32 | 512 | 74.9 | 91.8 | 81.0 | 82.0 | 75.8 | 66.3 |
| BoIR | HRNet-W32 | 512 | 70.6 | 89.9 | 76.5 | 77.1 | 71.2 | 63.0 |
| BoIR | HRNet-W48 | 640 | 71.2 | 90.3 | 76.7 | 77.8 | 71.8 | 63.5 |
| BoIR* | HRNet-W32 | 512 | 75.8 | 92.2 | 82.3 | 82.3 | 76.5 | 67.5 |
| BoIR* | HRNet-W48 | 640 | 77.2 | 92.4 | 83.5 | 82.7 | 78.1 | 69.8 |

Table 2: Comparison with state-of-the-art methods on CrowdPose test set. Best scores are marked as bold for small(e.g. HRNet-W32) and large(e.g. HRNet-W48) models respectively. Models with * are trained on COCO and finetuned on CrowdPose.

4 Experiments

4.1 Datasets and evaluation metrics

We evaluated the performance of our approach on four benchmark datasets.

COCO Keypoint 2017 [17] comprises train (57K images), val (5K images), and test-dev (20K images) splits, annotated with 17 keypoints. We use train split for training, and val split for hyperparameter tuning.

CrowdPose [16] consists of 20K images and 80K instances, annotated with 14 keypoints. Following the evaluation protocol of [54], we use trainval split (12K images, 43.4K instances) for training and test split (8K images, 29K instances) for evaluation.

OCHuman [42] is targeted for evaluation on crowded scenes with extreme conditions. 2,500 images are for val set, and 2,231 images are for test set. We evaluate our method following [10, 54].

Evaluation metrics. We follow COCO evaluation protocol, where AP(Average Precision) and AR(Average Recall) are computed based on OKS(Object Keypoint Similarity) with varying thresholds, including AP (averaged AP), AP⁵⁰ (AP at OKS=0.5), and AP⁷⁵ (AP at OKS=0.75). In case of CrowdPose, we additionally report metrics based on crowd index, including AP^E (easy), AP^M (medium), and AP^H (hard).

4.2 Implementation details

Our implementation is based on [54]. We use HRNet-W32 and HRNet-W48 as backbone networks and perform hyperparameter tuning with COCO val set results. We apply

| Method | Backbone | COCO val | | OCHuman val | | OCHuman test | |
|----------|-----------|-------------|-------------|-------------|-------------|--------------|-------------|
| | | AP | AR | AP | AR | AP | AR |
| DEKR [6] | HRNet-W32 | 68.0 | 73.0 | 37.9 | - | 36.5 | - |
| DEKR [6] | HRNet-W48 | 71.0 | 76.0 | - | - | - | - |
| CID [52] | HRNet-W32 | 69.8 | 75.4 | 44.9 | - | 44.0 | - |
| CID [52] | HRNet-W48 | - | - | 46.1 | - | 45.0 | - |
| BoIR | HRNet-W32 | 70.6 | 76.3 | 47.4 | 80.1 | 47.0 | 80.3 |
| BoIR | HRNet-W48 | 72.5 | 78.3 | 49.4 | 80.8 | 48.5 | 80.7 |

Table 3: Comparison with state-of-the-art methods on COCO val and OCHuman val, test set. OCHuman performance is evaluated with COCO pretrained model without fine-tuning.

AdamW optimizer with initial learning rate 1.0e-3, weight decay 2.5e-2, and cosine learning rate scheduler with 10 warmup epochs, following [18]. For COCO, we train the model for 140 epochs on 4 GPUs (RTX 3090 for HRNet-W32, A6000 for HRNet-W48) with AMP, with 20 batch size for each device. For CrowdPose, similar to [34], we train the model for 310 epochs when training from scratch, while 100 epochs with 1 warmup epoch are applied for transfer learning. Following [4, 6, 52], we apply single scale test with flipping.

4.3 Comparison with state-of-the-arts

Results on COCO dataset. We report COCO val results in Table 3, and test-dev results in Table 1. Our method outperforms existing state-of-the-art under the same or similar backbone. Our method with HRNet-W32 backbone outperforms CID by 0.8 AP on val and 0.6 AP on test-dev. Similarly, we achieve 0.5 AP improvement on test-dev with HRNet-W48 backbone. Furthermore, we conducted *t*-test on five COCO val set results from respective methods, and our method achieved statistically significant and consistent improvements over CID with *p*-value 6.4×10^{-5} .

Results on CrowdPose dataset. We compare other methods on CrowdPose test in Table 2. BoIR is the second best among state-of-the-art methods. Nonetheless, our method suffers from performance drop of 0.7 AP on the HRNet-W32 backbone and 1.1 AP on HRNet-W48 backbone. We speculate that as the model size increases, the model suffers from insufficient amount of training data on CrowdPose, as the performance difference between CID and ED-Pose on CrowdPose is also reversed on COCO. To validate the hypothesis, we introduce finetuning on CrowdPose using the model weights trained on COCO train set. Finetuning strategy is proven to be far more effective, surpassing existing state-of-the-art by 4.5 AP with the HRNet-W32 backbone, and 4.9 AP with HRNet-W48 backbone. For a fair comparison, we additionally conducted the same finetuning strategy on CID, and our method also outperforms the baseline by 0.9 AP.

OCHuman results. Comparison on OCHuman is summarized in Table 3. Following the protocol in [10], we evaluate the model trained on COCO without finetuning on OCHuman. BoIR outperforms comparative methods on both val and test set by large margin. Therefore, our instance representation learning is effective especially for crowded scenes.

| Bbox Mask Loss | Bbox Head | AP | Emb. Loss | Dist. Metric | AP |
|----------------|-----------|------|-------------|--------------|------|
| | | 69.6 | Contrastive | cosine | 70.3 |
| ✓ | | 70.2 | Contrastive | L2 | 70.2 |
| | ✓ | 70.4 | AE | L2 | 70.6 |
| ✓ | ✓ | 70.6 | | | |

Table 4: Left: Ablation study of Bbox Mask Loss and bbox regression head on COCO val set. Right: Ablation study of embedding loss function and distance metric on COCO val set, where Bbox Mask Loss and bbox head are used.

| Method | Backbone | # params. (M) | Time (ms) | AP |
|---------|-----------|---------------|-----------|------|
| CID | HRNet-W32 | 29.3 | 86.7 | 69.8 |
| CID | HRNet-W48 | 65.4 | - | - |
| ED-Pose | ResNet-50 | 47.9 | 113.9 | 71.6 |
| ED-Pose | Swin-L | 218.8 | 272.1 | 74.3 |
| BoIR | HRNet-W32 | 31.8 | 110.6 | 70.6 |
| BoIR | HRNet-W48 | 68.9 | 167.3 | 72.5 |

Table 5: Computational cost comparison on COCO val set. Inference time is measured with single RTX 3090 and 1 batch size.

4.4 Ablation study

We have performed ablative experiments, as illustrated in Table 4. The effectiveness of Bbox Mask Loss and Bbox Head has been validated by assessing four possible combinations, and the result shows that our proposed methods are useful. We additionally conduct ablative experiments on embedding losses and distance metrics. AE loss turns out to be superior to Contrastive loss. We hypothesize that L2 distance with Gaussian kernel used in AE loss is better suited for keypoint evaluation criteria, as claimed in [22]. We also extensively compare computational cost in Table 5. Our method manages to keep the computational cost within a reasonable extent, compared to ED-Pose.

For qualitative and visual analysis, we compare our method with CID in Fig. 3. The color coding in the figure represents the t-SNE results of the learned features (3 dim). For skateboarding (left), CID missed the border, explaining the inconsistent and less disentangled features. This was evident from the similarity of the color of the border to the background, which also appeared overall noisy. In contrast, BoIR demonstrated clear, consistent, and distinct t-SNE colors for the border, effectively separating it from the background. For right, BoIR further exhibited distinct colors for different individuals, successfully distinguishing between two closely interacting people, where CID failed.

We also visualize the behaviour of our method for overlapping instances in Fig. 1 (b-f). Two boxes A and B in (b), (c), and (d) respectively show the feature similarities of CID from the respective box centers. Similarly, (e) and (f) show the corresponding similarities of BoIR. Our method demonstrates a notably effective separation of features for closely interacting individuals.

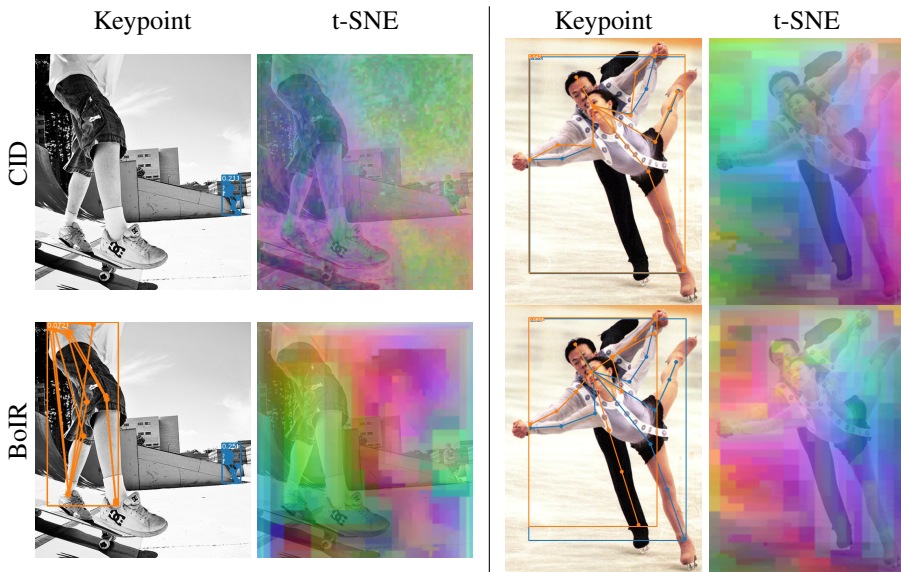


Figure 3: Example outcomes using our approach. The image on the left is from the COCO val set, while the image on the right is from the CrowdPose test set. We employed t-SNE, running it for 250 iterations, on the output backbone feature, with three output dimensions per pixel, corresponding directly to normalized RGB values.

5 Conclusion

This paper proposes a new multi-person pose estimation method using bounding box-supervised instance representation learning, called BoIR. It provides rich spatial supervision, utilizing embedding similarity as a soft mask for positive sampling, and the background region as a negative sample. It also incorporates auxiliary tasks for richer multi-task learning, without additional computation cost during inference. Our instance embedding can effectively disentangle instances in crowded scenes, surpassing comparable state-of-the-art methods on multiple human pose estimation benchmarks. Despite notable performance improvement with transfer learning, effective representation learning on small training data is a remaining issue, and we plan to mitigate the limitation as future work. Potential future work also involves enhancing auxiliary supervision by incorporating additional tasks, *e.g.* action recognition, and leveraging image captions within the framework of multi-modal contrastive training.

Acknowledgements

This work was supported by the National Research Foundation of Korea (NRF) grant (No. 2021R1A2C2012195), Institute of Information & Communications Technology Planning & Evaluation (IITP) grant (2021-0-00537, Visual Common Sense Through Self-supervised Learning for Restoration of Invisible Parts in Images), and IITP grant (2020-0-01336, Artificial Intelligence Graduate School Program, UNIST), all funded by the Korea government (MSIT). We thank Yunpyo An and Jihun Lee for their assistance in conducting a qualitative analysis of the results.

Appendix

In this supplementary material, we provide further details about 1) Bbox Mask Loss 2) architecture composition 3) additional experiments 4) visualization and comparative analysis with CID.

6 Bbox Mask Loss

From the instance embedding map e , we first apply L2 normalization on e . Then, we sample instance embedding p from e and compute respective loss terms as following:

$$\mathcal{L}_{pull}^{in} = \frac{1}{DN} \sum_{i=1}^N \sum_{d=1}^D (p_{(i,d)} - \bar{p}_{(i,d)})^2 \quad (7)$$

$$\mathcal{L}_{push}^{out} = \frac{1}{N} \sum_{i=1}^N \exp \left\{ -\frac{\beta}{D} \sum_{d=1}^D (p_{(i,d)} - \bar{p}_{(i,d)}^c)^2 \right\} \quad (8)$$

$$\mathcal{L}_{push}^{inst} = \frac{1}{\frac{N(N-1)}{2}} \sum_{i=1}^N \sum_{j>i}^N \exp \left\{ -\frac{\beta}{D} \sum_{d=1}^D (p_{(i,d)} - p_{(j,d)})^2 \right\} \quad (9)$$

D is the dimension of the instance embedding, N is the number of ground-truth instances in an image, and i, j represent the instance indices. $\beta = \frac{1}{2\sigma^2}$ is a scaling coefficient for the Gaussian kernel proposed in Associative Embedding [2].

7 Architecture Composition

In the case of 512x512 input size, output heatmap size is set to 128x128. In the case of 640x640 input size, output heatmap size is 160x160. HRNet-W32 backbone outputs 480 channels, due to concatenation of all block outputs. Similarly, HRNet-W48 backbone outputs 720 channels.

In the case of auxiliary task heads, 1 residual block and 1 convolution layer are applied. Residual block receives 256 input channels and outputs 128 channels. The final convolution layer outputs task-specific output channels. In case of bottom-up keypoint head, it is the number of keypoints (17 in COCO, 14 in CrowdPose). In case of the bounding box head, it outputs 4 channels (left, top, right, bottom distance). In case of embedding head, it is D . All convolution layers in the auxiliary task head have a 3x3 kernel size.

In case of instance-wise keypoint regression head, 64 hidden channel size is applied for HRNet-W32 backbone. 96 hidden channel size is used for HRNet-W48 backbone.

8 Additional Experiments

We report full comparative evaluation results on COCO val set on Table 6. CID paper does not report full results, so we report the scores using the provided trained model weights. Since HRNet-W48 backbone model is not available, CID's HRNet-W48 results are not reported. BoIR outperforms all comparative state-of-the-arts except for HRNet-W48 backbone

| Method | Backbone | Input size | AP | AP ⁵⁰ | AP ⁷⁵ | AP ^M | AP ^L | AR |
|----------------------|-------------|------------|-------------|------------------|------------------|-----------------|-----------------|-------------|
| Top-down methods | | | | | | | | |
| SBL [37] | ResNet-152 | 384×288 | 74.3 | 89.6 | 81.1 | 70.5 | 81.6 | 79.7 |
| HRNet [62] | HRNet-W32 | 384×288 | 75.8 | 90.6 | 82.5 | 72.0 | 82.7 | 80.9 |
| Bottom-up methods | | | | | | | | |
| HrHRNet [9] | HrHRNet-W32 | 512 | 67.1 | 86.2 | 73.0 | 61.5 | 76.1 | - |
| HrHRNet [9] | HrHRNet-W48 | 640 | 69.9 | 87.2 | 76.1 | 65.4 | 76.4 | - |
| DEKR [6] | HRNet-W32 | 512 | 68.0 | 86.7 | 74.5 | 62.1 | 77.7 | 73.0 |
| DEKR [6] | HRNet-W48 | 512 | 71.0 | 88.3 | 77.4 | 66.7 | 78.5 | 76.0 |
| SWAHR [19] | HrHRNet-W32 | 512 | 67.1 | 86.2 | 73.0 | 61.5 | 76.1 | - |
| SWAHR [19] | HrHRNet-W48 | 640 | 69.9 | 87.2 | 76.1 | 65.4 | 76.4 | - |
| Single stage methods | | | | | | | | |
| PETR [29] | ResNet-101 | 800 | 70.0 | 88.5 | 77.5 | 63.6 | 79.4 | - |
| ED-Pose [40] | ResNet-50 | 800 | 71.6 | 89.6 | 78.1 | 65.9 | 79.8 | - |
| CID [54] | HRNet-W32 | 512 | 69.8 | 88.5 | 76.6 | 64.0 | 78.9 | 75.4 |
| BoIR | HRNet-W32 | 512 | 70.6 | 89.2 | 77.4 | 65.1 | 79.0 | 76.3 |
| BoIR | HRNet-W48 | 640 | 72.5 | 89.9 | 79.1 | 68.2 | 79.4 | 78.3 |

Table 6: Comparison with state-of-the-art methods on COCO val set. Best scores are marked as bold for small(e.g. HRNet-W32) and large(e.g. HRNet-W48) models respectively.

| β | AP | AR | Emb. Dim | AP | AR |
|---------|------|------|----------|------|------|
| 0.5 | 69.8 | 75.7 | 1 | 70.4 | 76.3 |
| 1 | 70.3 | 76.2 | 8 | 70.4 | 76.3 |
| 5 | 70.2 | 76.0 | 16 | 70.4 | 76.2 |
| 10 | 70.5 | 76.2 | 32 | 70.4 | 76.3 |
| 15 | 70.2 | 76.0 | 64 | 70.1 | 75.8 |
| | | | 128 | 70.5 | 76.2 |

Table 7: Left: Ablation experiment of β on COCO val set. $D = 128$ by default. Right: Ablation experiment of embedding dimension D on COCO val set. $\beta = 10$ by default.

on AP^L. Our method with HRNet-W32 backbone outperforms CID by 0.8 AP. Our method with HRNet-W48 even outperforms ED-Pose by 0.9 AP.

We report ablation experiments on AE’s Gaussian kernel scaling coefficient β and embedding dimension D on COCO val set on Table 7. For fast training and simplicity, we use Bbox Mask Loss and do not use bbox head during experiment. In case of β , 10 was the best among the options. In case of D , changing the embedding dimension shows little performance difference. We conjecture that AE loss for higher dimensions needs more refinement to benefit high dimensional representation. The primary cause is L2 normalization over embedding dimension before loss computation, which significantly drops the loss scale and floating point precision compared to the original AE loss formulation. Simply removing the normalization would cause unstable training with AMP, so further research is required to improve the current framework.

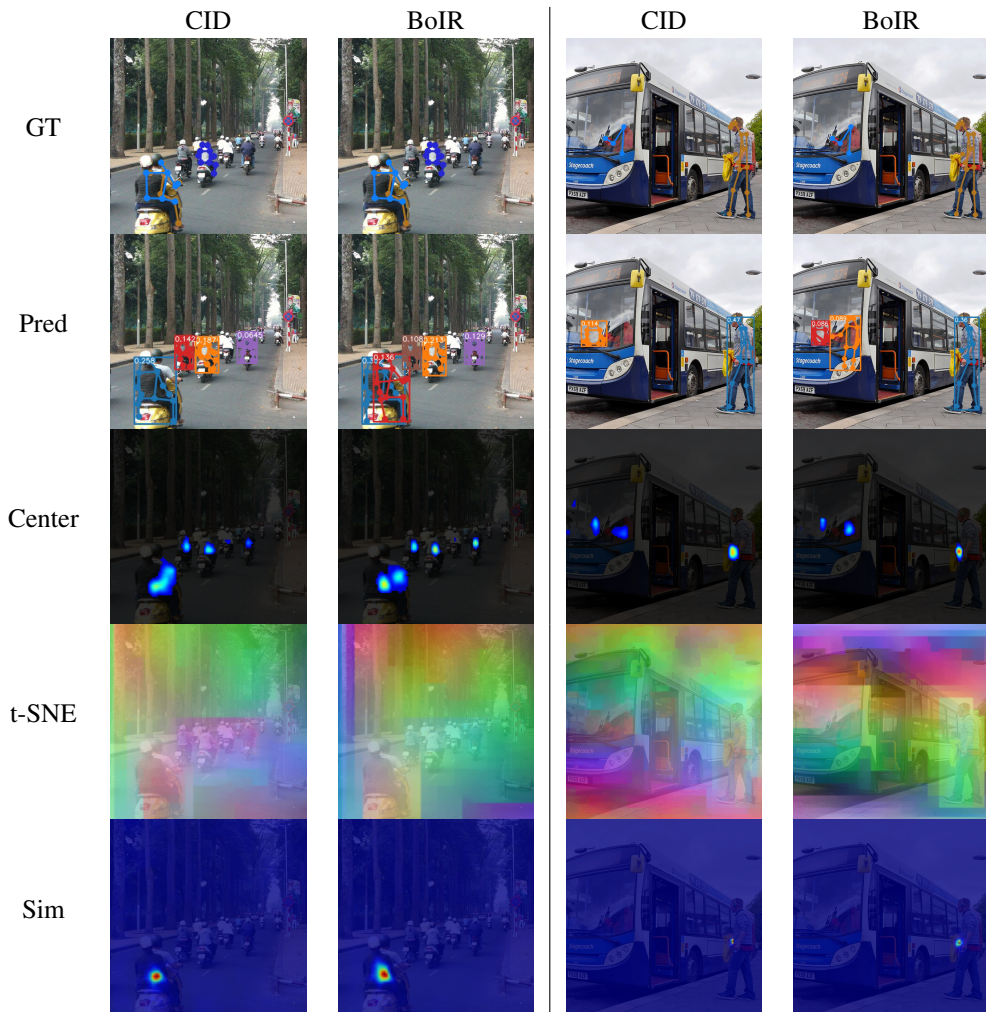


Figure 4: Comparative visualization on COCO val set.

9 Visualization

We provide extensive outputs of our model in Fig. 4 and Fig. 5. We visualize keypoint prediction outputs along with instance center heatmap, t-SNE of backbone output feature, and feature similarity between top-1 confident instance parameter and the entire feature map. t-SNE is applied on the output backbone feature for 250 iterations with 3 output dimensions per pixel, which directly corresponds to normalized RGB values. Instance similarity is measured by computing the L2 distance between the top-1 confident instance’s parameter and the feature map, and then applying a Gaussian Kernel.

We additionally report failure cases of BoIR in Fig. 6. In case of the left images, BoIR produces duplicated predictions on the same person, due to wide activation area of the center heatmap. In case of the right images, BoIR places occluded joints on implausible positions, while this does not affect evaluation performance. However, BoIR generally produces dis-

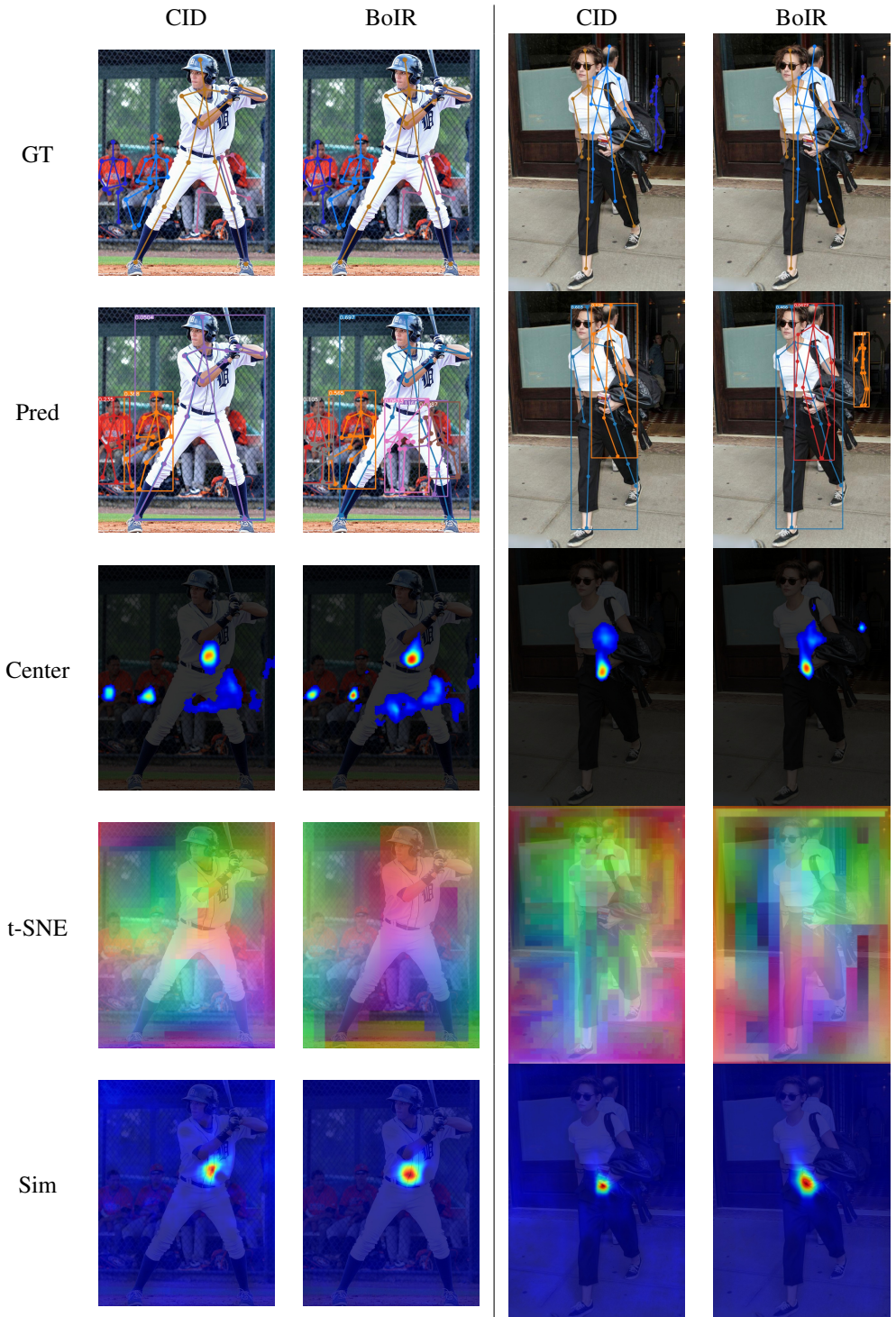


Figure 5: Comparative visualization on CrowdPose test set.

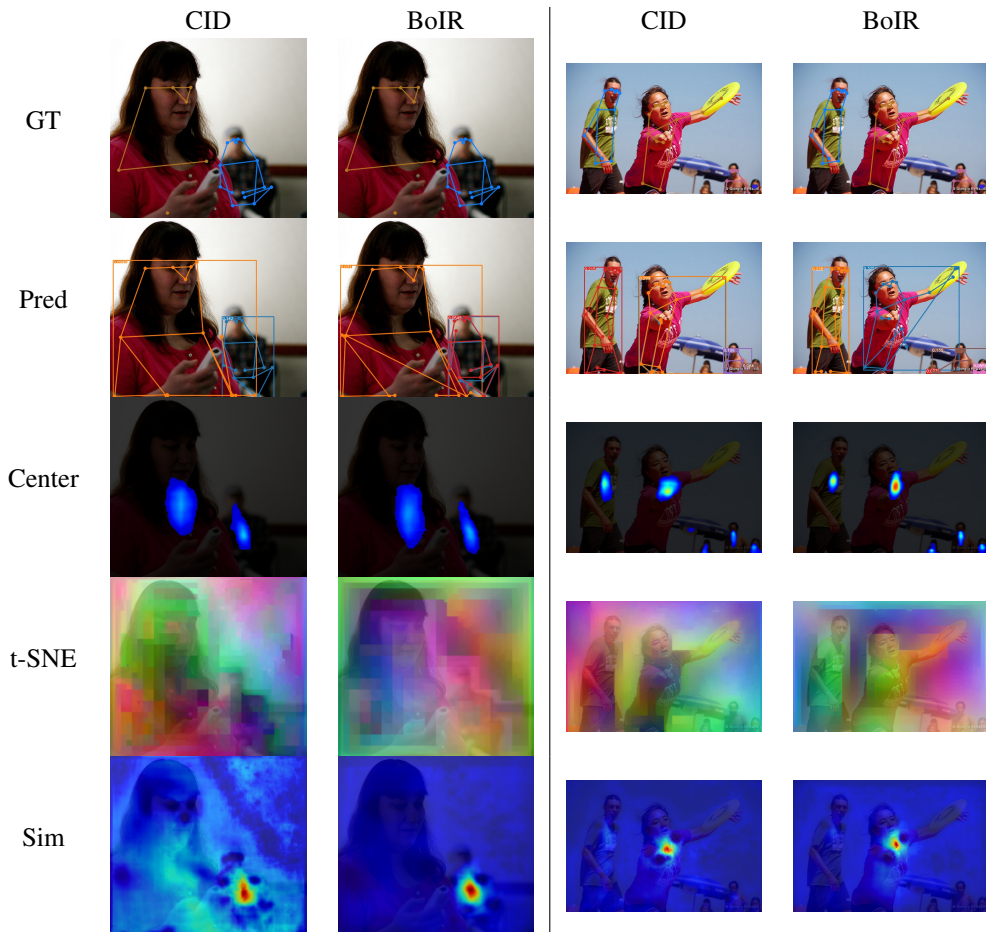


Figure 6: Failure cases on COCO val set.

entangled instance features and detects people better than CID.

References

- [1] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017.
- [2] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation, 2017. URL <https://arxiv.org/abs/1706.05587>.
- [3] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [4] 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 (CVPR)*, June 2020.
- [5] Qing Gao, Jinguo Liu, Zhaojie Ju, and Xin Zhang. Dual-hand detection for human-robot interaction by a parallel network based on hand detection and body pose estimation. *IEEE Transactions on Industrial Electronics*, 66(12):9663–9672, 2019. doi: 10.1109/TIE.2019.2898624.
- [6] Zigang Geng, Ke Sun, Bin Xiao, Zhaoxiang Zhang, and Jingdong Wang. Bottom-up human pose estimation via disentangled keypoint regression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14676–14686, June 2021.
- [7] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006.
- [8] 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.
- [9] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020.
- [10] Sheng Jin, Wentao Liu, Enze Xie, Wenhai Wang, Chen Qian, Wanli Ouyang, and Ping Luo. Differentiable hierarchical graph grouping for multi-person pose estimation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VII 16*, pages 718–734. Springer, 2020.
- [11] Rawal Khirodkar, Visesh Chari, Amit Agrawal, and Amrith Tyagi. Multi-instance pose networks: Rethinking top-down pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3122–3131, October 2021.
- [12] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi. Pifpaf: Composite fields for human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [13] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
- [14] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.
- [15] Chuyi Li, Lulu Li, Hongliang Jiang, Kaiheng Weng, Yifei Geng, Liang Li, Zaidan Ke, Qingyuan Li, Meng Cheng, Weiqiang Nie, Yiduo Li, Bo Zhang, Yufei Liang, Linyuan Zhou, Xiaoming Xu, Xiangxiang Chu, Xiaoming Wei, and Xiaolin Wei. Yolov6: A single-stage object detection framework for industrial applications, 2022.

- [16] 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 (CVPR)*, June 2019.
- [17] 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 David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 740–755, Cham, 2014. Springer International Publishing. ISBN 978-3-319-10602-1.
- [18] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11976–11986, 2022.
- [19] Zhengxiong Luo, Zhicheng Wang, Yan Huang, Liang Wang, Tieniu Tan, and Erjin Zhou. Rethinking the heatmap regression for bottom-up human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13264–13273, June 2021.
- [20] Weian Mao, Zhi Tian, Xinlong Wang, and Chunhua Shen. Fcpose: Fully convolutional multi-person pose estimation with dynamic instance-aware convolutions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9034–9043, June 2021.
- [21] William McNally, Kanav Vats, Alexander Wong, and John McPhee. Rethinking keypoint representations: Modeling keypoints and poses as objects for multi-person human pose estimation. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Computer Vision – ECCV 2022*, pages 37–54, Cham, 2022. Springer Nature Switzerland. ISBN 978-3-031-20068-7.
- [22] Alejandro Newell, Zhiao Huang, and Jia Deng. Associative embedding: End-to-end learning for joint detection and grouping. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/8edd72158ccd2a879f79cb2538568fdc-Paper.pdf>.
- [23] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [24] George Papandreou, Tyler Zhu, Liang-Chieh Chen, Spyros Gidaris, Jonathan Tompson, and Kevin Murphy. Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geometric embedding model. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 269–286, 2018.
- [25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

- [26] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- [27] 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, 2015.
- [28] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.
- [29] Dahu Shi, Xing Wei, Liangqi Li, Ye Ren, and Wenming Tan. End-to-end multi-person pose estimation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11069–11078, June 2022.
- [30] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL <https://proceedings.neurips.cc/paper/2016/file/6b180037abbebea991d8b1232f8a8ca9-Paper.pdf>.
- [31] 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 (CVPR)*, June 2019.
- [32] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, 2019.
- [33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [34] Dongkai Wang and Shiliang Zhang. Contextual instance decoupling for robust multi-person pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11060–11068, June 2022.
- [35] Haixin Wang, Lu Zhou, Yingying Chen, Ming Tang, and Jinqiao Wang. Regularizing vector embedding in bottom-up human pose estimation. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Computer Vision – ECCV 2022*, pages 107–122, Cham, 2022. Springer Nature Switzerland. ISBN 978-3-031-20068-7.
- [36] Adrian Wolny, Qin Yu, Constantin Pape, and Anna Kreshuk. Sparse object-level supervision for instance segmentation with pixel embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4402–4411, 2022.
- [37] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *European Conference on Computer Vision (ECCV)*, 2018.
- [38] Yabo Xiao, Kai Su, Xiaojuan Wang, Dongdong Yu, Lei Jin, Mingshu He, and Zehuan Yuan. Querypose: Sparse multi-person pose regression via spatial-aware part-level query. *Advances in Neural Information Processing Systems*, 35:12464–12477, 2022.

- [39] Nan Xue, Tianfu Wu, Gui-Song Xia, and Liangpei Zhang. Learning local-global contextual adaptation for multi-person pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13065–13074, June 2022.
- [40] Jie Yang, Ailing Zeng, Shilong Liu, Feng Li, Ruimao Zhang, and Lei Zhang. Explicit box detection unifies end-to-end multi-person pose estimation. In *International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=s4WVupnJjmX>.
- [41] Pengfei Zhang, Cuiling Lan, Wenjun Zeng, Junliang Xing, Jianru Xue, and Nanning Zheng. Semantics-guided neural networks for efficient skeleton-based human action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [42] 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.
- [43] Jingxiao Zheng, Xinwei Shi, Alexander Gorban, Junhua Mao, Yang Song, Charles R. Qi, Ting Liu, Visesh Chari, Andre Cornman, Yin Zhou, Congcong Li, and Dragomir Anguelov. Multi-modal 3d human pose estimation with 2d weak supervision in autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 4478–4487, June 2022.
- [44] Zhaohui Zheng, Ping Wang, Wei Liu, Jinze Li, Rongguang Ye, and Dongwei Ren. Distance-iou loss: Faster and better learning for bounding box regression. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 12993–13000, 2020.
- [45] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points, 2019. URL <https://arxiv.org/abs/1904.07850>.