

# Single-stage Multi-human Parsing via Point Sets and Center-based Offsets

Jiaming Chu<sup>1†</sup>, Lei Jin<sup>1\*†</sup>, Junliang Xing<sup>2</sup> and Jian Zhao<sup>3\*</sup>

<sup>1</sup>Beijing University of Posts and telecommunications

<sup>2</sup>Tsinghua University

<sup>3</sup>Institute of North Electronic Equipment

{chujiaming886, jinlei}@bupt.edu.cn, jlxing@tsinghua.edu.cn, zhaojian90@u.nus.edu

## Abstract

This work studies the multi-human parsing problem. Existing methods, either following top-down or bottom-up two-stage paradigms, usually involve expensive computational costs. We instead present a high-performance **Single-stage Multi-human Parsing (SMP)** deep architecture that decouples the multi-human parsing problem into two fine-grained sub-problems, *i.e.*, locating the human body and parts. SMP leverages the point features in the barycenter positions to obtain their segmentation and then generates a series of offsets from the barycenter of the human body to the barycenters of parts, thus performing human body and parts matching without the grouping process. Within the SMP architecture, we propose a *Refined Feature Retain* module to extract the global feature of instances through generated mask attention and a *Mask of Interest Reclassify* module as a trainable plug-in module to refine the classification results with the predicted segmentation. Extensive experiments on the *MHPv2.0* dataset demonstrate the best effectiveness and efficiency of the proposed method, surpassing the state-of-the-art method by 2.1% in  $AP_{50}^p$ , 1.0% in  $AP_{vol}^p$ , and 1.2% in  $PCP_{50}$ . In particular, the proposed method requires fewer training epochs and a less complex model architecture. We will release our source codes, pretrained models, and online demos to facilitate further studies.

## 1 Introduction

Instance-Aware Multi-Human Parsing (IAMHP) aims to cut apart the parts of human body according to semantics, and group them by human instance. Compared with semantic segmentation task and instance segmentation task, it is more challenging. Because for each pixel in the picture, it is necessary to judge not only the semantic label of the part level, but also the instance label of the human level. This paper attempts to propose a superior framework to conduct instance-aware multi-human parsing, which can significantly simplify the complicated processing pipelines of previous researches.

\*Lei Jin and Jian Zhao are corresponding authors. † Jiaming Chu and Lei Jin are equal contribution.

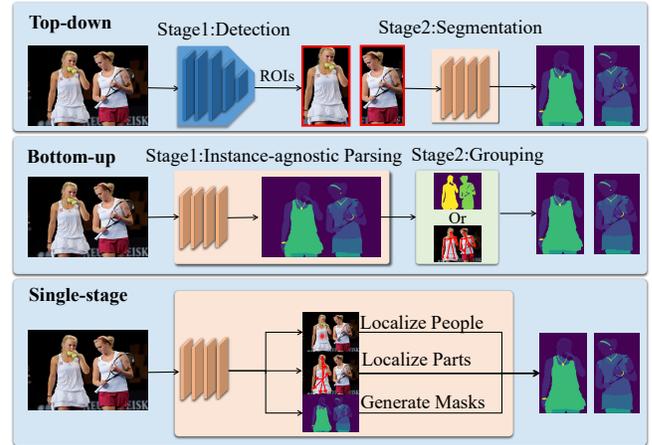


Figure 1: The differences between top-down, bottom-up, and single-stage methods in multi-human parsing task. The single-stage method is parallel computing, and neither the instance detection of the top-down method nor the grouping process of the bottom-up method are required. Best viewed in color.

A few efforts have been proposed to model the IAMHP problem. The existing multi-human parsing method could be roughly divided into two categories: bottom-up and top-down methods. The bottom-up methods [Zhao *et al.*, 2017] usually parse all the human parts through instance segmentation or semantic segmentation, and then group the human parts belonging to the same human instance through the different grouping methods, *e.g.*, clustering by human instance segmentation, instance edge prediction, *etc.* The top-down methods [Yang *et al.*, 2018; Yang *et al.*, 2020] usually detect human instances first, and then parse human parts for each human instance individually. Albeit encouraging results achieved, these methods still have the following limitations. Both bottom-up and top-down methods are involved two stages and complicated post-processings, *i.e.*, part parsing and grouping, instance detection and single-person parsing. Such sequential stages lead to redundant processes and high computational cost.

More recently, to eliminate redundant processes and improve model calculation efficiency, researchers proposed single-stage framework [Nie *et al.*, 2019; Geng *et al.*, 2021; Lei *et al.*, 2022; McNally *et al.*, 2021; Zhou *et al.*, 2021] for the human pose estimation task and have achieved impressive

results. The key idea of single-stage framework is a holistic representation for the human instance, *e.g.*, the hierarchical structure pose representation in SPM [Nie *et al.*, 2019] and the decoupled 3D pose representation in DRM [Lei *et al.*, 2022]. Overall, these methods only need a center point and center-based offsets to represent the human pose. The simple model structure and parallel inference schema greatly improve the efficiency of model inference. However, for multi-human parsing task, only center point can not encode fine-grained semantic information for human parts. To construct a single-stage multi-human parsing framework, how to unify human-instance-level information and human-part-level information is still a challenging and unexplored problem.

In this paper, we explore the possibility to understand human body with point sets and center-based offsets to conduct multi-human parsing. Specifically, the point sets are composed of human barycenter and the barycenters of parts, and the center-based offsets are the displacements from human barycenter to part barycenters. With such representation, we realize a Single-stage Multi-human Parsing (SMP) framework, which omits the time-consuming grouping process. In particular, we decouple the IAMHP task into four sub-tasks, termed, human instance localization, part instance localization, part instance segmentation and subordination mapping prediction of the two granularity instances. Additionally, we also propose the Refined Feature Retain (RFR) module and the Mask of Interest Reclassification (MIR) module. The former utilizes the instance correlation between mask features as an attention to make the model learn more relevant features. The latter can refine the classification results by taking the output of part segmentation as region of interest.

Without the instance detection in top-down methods and the complex grouping process in bottom-up methods, our method completes the prediction of mapping relationship in parallel while predicting the location and segmentation of instances, thus improving the inference efficiency. Specifically, our method reduces the inference time by 102.5% and 91.2% with ResNet-50 [He *et al.*, 2016] and ResNet-101 as backbone, respectively. Finally, it is worth mentioning that, SMP also achieves state-of-the-art performance, outperforming the best result [Dong *et al.*, 2022b] by 2.1% in  $AP_{50}^p$ , 1.0% in  $AP_{vol}^p$ , and 1.2% in  $PCP_{50}$  on MHPv2.0 [Zhao *et al.*, 2018] dataset with the same settings. To show the generalization of SMP, we test SMP on DensePose COCO [Guler *et al.*, 2018] (Please refer to the appendix). Experiments show that SMP achieves superior performance on DensePose COCO as well.

Our contributions are summarized as follows.

- We make the first attempt to understand human body with point sets and center-based offsets. With this as a guideline, to the best knowledge of us, we introduce the first single-stage framework for instance-aware multi-human parsing.
- We propose a Refined Feature Retain (RFR) module which uses mask attention to guide the model to extract features, and further put forward a Mask of Interest Reclassification (MIR) module to optimize the classification results of the model.
- Our method significantly outperforms all the state-of-the-

art top-down and bottom-up methods in AP and PCP metrics on MHPv2.0 benchmark, obtaining fastest inference speed.

## 2 Related Work

### 2.1 Instance Segmentation

In general, existing instance segmentation techniques can be roughly categorized as top-down or bottom-up methods.

As a typical top-down method, Mask R-CNN [He *et al.*, 2017] first detects instances, and then performs semantic segmentation to achieve pixel-level predictions. Dai *et al.* [Dai *et al.*, 2016] utilize the location sensitivity of instances to segment instances based on the FCN network [Long *et al.*, 2014], and assign the prediction task of instance segmentation to different grids.

On the other hand, the bottom-up methods aim to learn better pixel-level feature representations. Some of them adopt the loss function about clustering [Neven *et al.*, 2019; Ying *et al.*, 2021], and others attempt to utilize grids to generate strongly related location features [Wang *et al.*, 2019; Bolya *et al.*, 2019]. But both of them output intensive predictions for each grid or pixel, which incurs redundant calculation, as not all pixels belong to specific instances.

In order to reduce unnecessary computation, a handful of works attempt to leverage conditional convolution [Yang *et al.*, 2019]. It replaces the conventional convolution kernel with the predicted features, further reducing considerable computation amount. So there are also many works [Tian *et al.*, 2020; Sofiiuk *et al.*, 2019; Karras *et al.*, 2018; Wang *et al.*, 2020] like ConDist, SOLOv2, utilize dynamically generated kernel to improve the capacity of the model, hence the mask head can be more compact, leading to faster inference speed. Such similiar kernel mechanism is also widely used in other fields [Cao *et al.*, 2022; Chen *et al.*, 2022; Ni *et al.*, 2020].

### 2.2 Multi-human Parsing

Compared with instance segmentation [Dong *et al.*, 2022b], semantic segmentation [Dong *et al.*, 2022a; He *et al.*, 2022; Li *et al.*, 2021] and pose estimation [Zhang *et al.*, 2021], multi-human parsing is a task that not only needs to distinguish the part instance label of each pixel, but also needs to determine the subordinate relationship between human parts and human instances. Similarly, existing multi-human parsing methods also can be categorized into top-down or bottom-up methods.

The top-down methods usually need detector and feature aligning. Parsing RCNN [Yang *et al.*, 2018] and RP-RCNN [Yang *et al.*, 2020] both generate bounding box first and utilize a geometric and context encoding (GCE) module or GSE-FPN to get finest semantic pixel features. In addition, some methods [Ruan *et al.*, 2019] utilize extra ground truth to enhance the ability of the model to extract instance context features. SNT [Ji *et al.*, 2019] which relies on Mask RCNN to distinguish human instances, predicts human parts by attributes in stages, reducing the difficulty of prediction.

The bottom-up methods do not give priority to distinguish human instances. They usually predict the instance label and semantic label of pixels and group them. Some works [Zhao *et al.*, 2020; Li *et al.*, 2017] apply the GAN-based network to improve the learning ability to instances and

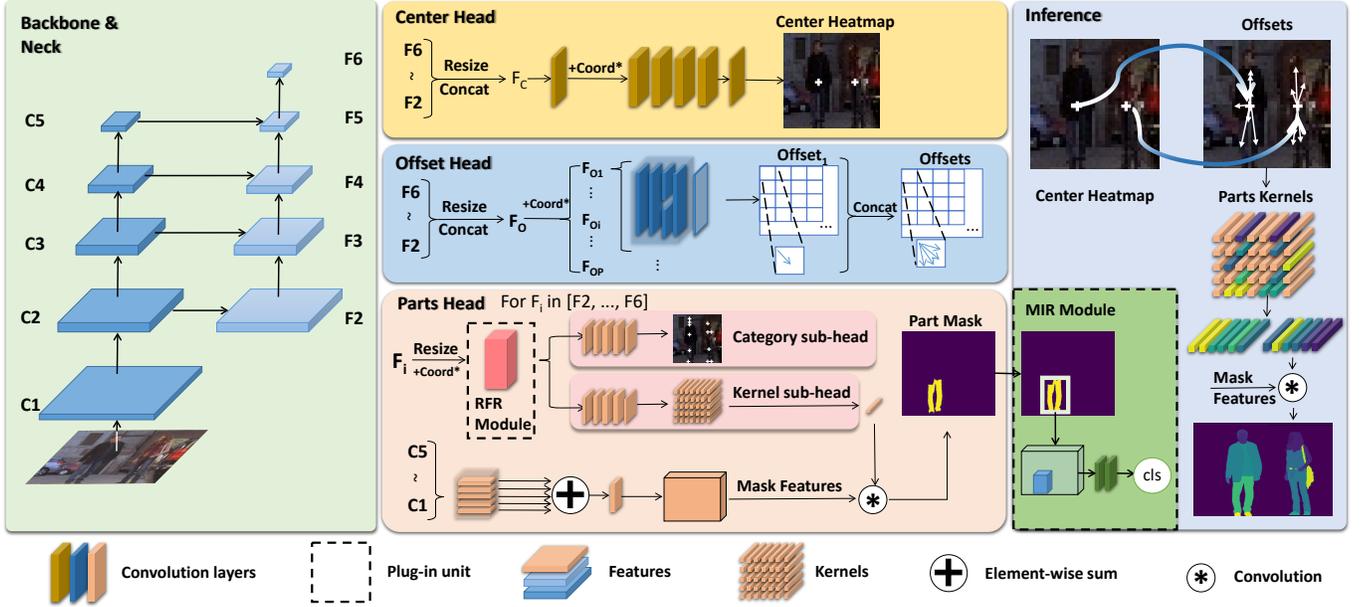


Figure 2: The Illustration of our Single-stage Multi-human Parsing (SMP) framework for instance-aware human semantic parsing. Coord\* refers to concatenating the relative coordinates to the feature. Details can be seen in Sec. 3.1 and Sec. 4.1. Best viewed in color.

semantic features of the model. Other works utilize additional groundtruth information such as pose [Zhou *et al.*, 2021] and edge [Gong *et al.*, 2018] to help the model learn instance characteristics.

### 3 Method

**Problem Formulation.** Given an image  $I$ , multi-human parsing aims to parse the human instances in  $I$ , and obtain their masks of body parts  $\mathcal{M} = \{M_m^{human}\}_{m=1}^N$ , where  $N$  denotes the total number of human instances in  $I$ . For each human instance, it is composed of part instance segmentation  $M_m^{human} = \sum_{i=1}^{C_P} M_i^{part}$ , where  $C_P$  is the number of part categories,  $M_i^{part}$  is the mask of  $i$ -th class part.

#### 3.1 Overview

The overview of our Single-stage Multi-human Parsing (SMP) framework is illustrated in Fig. 2. Firstly, we send an image to the Features Pyramid Network (FPN) [Lin *et al.*, 2016] to generate feature maps with different sizes. Then we process the feature maps with three dedicated heads, *i.e.*, center head, offset head, parts head, to predict human location and mask information. Finally, we can obtain instance-aware multi-human parsing results with the outputs of the three heads.

**Center Head.** The center head aims to predict the location of each individual human instance. To avoid the overlapped center problem, we utilize the barycenter of visible mask to represent each instance. Given an image  $I$ , the backbone and FPN extract its multi-scale semantic features  $F_i (i \in 2, \dots, 6)$ , as seen in Fig. 2. Next, We resize the multi-scale features to the same size  $S \times S$  with bilinear interpolation ( $S$  is set as 40), then concatenate them, and apply a  $3 \times 3$  convolution layer for feature fusion and compression. The fusion feature  $F_C$  can simultaneously perceive human instances of multiple scales.

Table 1: The settings of feature grids. Instances in different scales response in different feature maps. Human refers to the grids applied in center head and offset head.

	Parts					Human
Level	$F_2$	$F_3$	$F_4$	$F_5$	$F_6$	fusion
Scale	$\leq 96$	(48,192)	(96,384)	(192,768)	$\geq 384$	-
Grids	40	36	24	16	12	40

To further enhance the position information of the features, we adopt the coordinate convolution (CoordConv) [Liu *et al.*, 2018] to concatenate the relative position with the fusion features. Finally, we transfer the features by convolution layers to a one-channel heatmap  $H_{Center}^{S \times S \times 1}$ . The values in the  $H_{Center}^{S \times S \times 1}$  represent the existing confidence of the center.

**Offset Head.** In this paper, we predict the offsets from the barycenter of human body to the barycenters of its corresponding part instances to estimate the mapping relationship. To get features ( $F_O$  in the offset head of Fig. 2) from different scales, similar to the center head, we apply bilinear interpolation to resize the features to  $S \times S$  and concatenate them together ( $S$  is set as 40), and utilize CoordConv [Liu *et al.*, 2018] to obtain location information. Next, We employ the decoupled structure with AdaptBlock [Geng *et al.*, 2021] to predict the offset map for each human part. The input of each decoupled branch is only one part of the fusion features, which greatly reduces computation cost. All decoupled branch outputs will eventually be concatenated together to obtain offset maps  $M_{Offset}^{S \times S \times C_P \times 2}$ , where  $C_P$  is the number of parts class. The pixel around the human barycenter in offset maps contains  $C_P$  regression vectors, pointing to the barycenter positions of the  $C_P$  parts from the pixel.

**Part Head.** The goal of part head is to localize the part in-

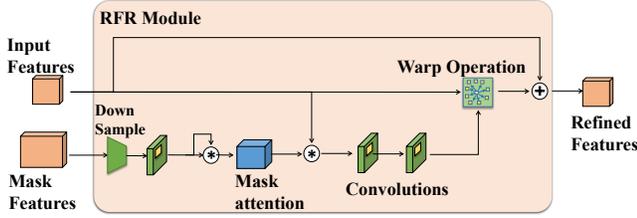


Figure 3: The detailed architecture about RFR module. It is placed after downsampling of category and kernel branch as shown in Fig. 2.

stances and generate their fine-grained masks. Inspired by SOLOv2 [Wang *et al.*, 2020] and YOLOv3 [Redmon and Farhadi, 2018], we utilize Features Pyramid Network (FPN) to predict part instances with different scales. The outputs of each level are resized to different grid sizes and compressed to the same number of channels. Instances with different scales are corresponding to different levels as shown in Tab. 1. The outputs are utilized as input for each prediction head: parts category sub-head and parts kernel sub-head. The parameters of these heads are shared across different levels. We adopt a similar model like center head, and we also use the barycenter of visible part mask to represent each part. A convolution layer with  $C_P$  output channels is at the end of the **part category sub-head**, and we could obtain the parts heatmap  $H_{Part}^{S \times S \times C_P}$  finally.

After localizing part instance, we refer to condition convolution [Tian *et al.*, 2020] and SOLOv2 to generate the refined mask. The input feature is also the single-layer feature of FPN, and the size is aligned with the  $H_{Part}$ . This sub-head also concatenates the relative coordinates by CoordConv. Its structure is the same as that of the  $H_{Part}$  sub-head, but the number of channels in the last convolution layer is  $C_K$  which is set as 256, then the final output of **part kernel sub-head** is  $F_{Kernel}^{S \times S \times C_K}$ .

To ensure the fineness of the generated mask, the feature map will not be downsampled to  $S \times S$ . The features employ the multi-scale output features of the backbone (C1,...,C5 in Fig. 2). The features at different levels are all upsampled to 1/4 the original resolution through bilinear interpolation. With element-wise summation and a convolution layer, the convoluted feature map  $F_{Mask}^{h \times w \times 256}$  is obtained, where  $h, w$  are 1/4 height and width of the original image.

### 3.2 Refined Feature Retain Module

To guide the model to better extract the features around the barycenter, we propose Refined Feature Retain (RFR) module. The position of RFR module in our framework SMP is shown in Fig. 2.

The main idea of RFR module is utilizing the mask feature as an attention to guide the learning of the category branch. Part head completes the instance segmentation by conditional convolution. For each value on output segmentation results, it is actually the inner product of the convolution kernel and the corresponding feature on the feature map. When the vectors are nearly normalized, the inner product can be regarded as the cosine similarity between them. During training, the

projection between the positive convolution kernel and the convolution features on the positive position will gradually increase. Since the prediction label is limited to  $[0, 1]$ , the convolution result will be closer to the cosine similarity. As seen in Fig. 3, through the self-correlation calculation of the convoluted feature map, we could get the attention map of the instance on the corresponding position. The self-attention map, namely mask attention, has superior instance guidance. By multiplying the category feature with the mask attention of each position, we could obtain a new refined feature map. We utilize the new feature as offset input and conduct warp operation to guide model adaptively obtain more instance information.

### 3.3 Mask of Interest Reclassify Module

The RFR module enhances the ability of feature extraction in the category sub-branch. However, the model is still difficult to classify the fashion classes with strong inter-class similarity [Zhang *et al.*, 2022b]. For the classification task, the model prefers to learn the target with relatively stable scale. Therefore, we propose a Mask of Interest Reclassify (MIR) module that can be trained quickly, and it can be used as a plug-in module to improve the model performance. The position of MIR module in our framework SMP is shown in Fig. 2.

Our model can output fine-grained segmentation as the Region of Interest (ROI) to achieve secondary classification. MIR module is separate and could utilize the output results of other branches. We choose the fusion feature of FPN as the input feature, extract the feature through consecutive convolution layers, and use semantic segmentation labels to supervise the feature. The mask generated by the part head is used as ROI to obtain local features. The features are reduced to a fixed size through the ROIalign with the size of 14, and the convolution layer with the kernel size of 14 is employed to compress the features again. Finally, two consecutive full connection layers are applied to output the classification results.

### 3.4 Loss Computation

The output of our model is divided into three parts: the heatmaps of the barycenters of both the human body and the parts, the mapping relationships between these two types of barycenter, and the fine masks of the parts. For each type of outputs, we apply distinct labels for supervision.

**Center and Part Localization Loss.** For the heatmaps localizing the barycenter of two granularity instances, we employ the traditional Focal Loss [Lin *et al.*, 2017]. Given the barycenter  $(c_x, c_y)$  and the height  $h$  and width  $w$  of the instance, the barycenter is extended to a rectangular area called center region  $(c_x, c_y, \epsilon h, \epsilon w)$ , where  $\epsilon$  is a hyperparameter to control the scale of region ( $\epsilon$  is set as 0.2). Any grid that falls into the center region is regarded as a positive sample.

**Offset Loss.** The offset prediction can be viewed as a regression problem. The supervision labels of offset head are the offset vectors. The offset vectors on each grid are obtained by subtracting the grid coordinates from the barycenter coordinates of the corresponding parts.

$$Offset_{i,j,p} = (c_{px} - i, c_{py} - j), \quad (1)$$

where  $i, j \in [0, S-1]$  are the grid coordinates,  $c_{px}$  and  $c_{py}$  are the barycenter coordinates of the person's part,  $p \in [0, C_p - 1]$

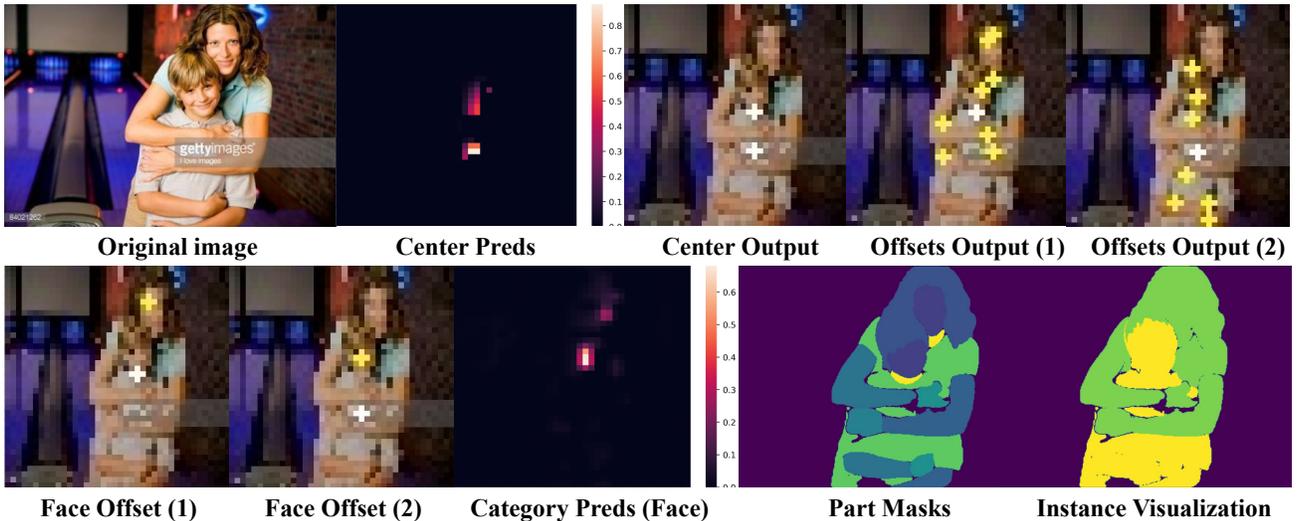


Figure 4: The visualization results of intermediate outputs for SMP. Center preds and Category preds (Face) are the visual confidence maps (“preds” is short for “predictions”). Center output is generated by maxpooling layer based on center preds. Yellow crosses in offsets output are the part instance locations pointed by offsets and white crosses are human centers. (1), (2) denote different human instances. Except original image, parts mask and instance visualization, other visualizations are resized to  $S \times S$  ( $S=40$  during inference) as description in Sec. 3.1 and Sec. 4.1. Best viewed in color.

is the label of part category. We adopted Smooth  $\ell_1$  Loss [Girshick, 2015] for the regression problem. When calculating the loss, we will multiply each grid with a scale-related weight to balance the impact of scale.

**Mask Loss.** For the segmentation of part instances, the label supervision method is the same as the general instance segmentation. We apply Dice Loss [Milletari *et al.*, 2016] for optimization.

**Total.** The loss of SMP is the sum of the loss for each head.

$$\mathcal{L}_{\text{total}} = \lambda_c \mathcal{L}_{\text{center}} + \lambda_p \mathcal{L}_{\text{part}} + \lambda_d \mathcal{L}_{\text{dice}} + \lambda_o \mathcal{L}_{\text{offset}}, \quad (2)$$

where  $\lambda_c$  and  $\lambda_p$  are both used to balance the center and part loss, with a value of 1.0.  $\lambda_d$  is used to balance the mask loss, with value of 3.0.  $\lambda_o$  is used to balance the Smooth  $\ell_1$  Loss for offsets, with value of 10.0. The above parameters are verified by the experiment search and SMP gets better results.

## 4 Experiment

### 4.1 Experiment Setup

**Dataset.** MHPv2.0 [Zhao *et al.*, 2018] is the largest and most challenging dataset in the field of multi-human parsing. It contains 15,403 training images, 5,000 validation images. Each picture contains 2 - 26 human instances, with an average of 3. The categories of 58 parts include 11 human body parts labels and 47 clothing and accessory labels. Meanwhile, the data set also provides labels for 16 keypoints of human pose. To show the generalization of SMP, we also conducted experiments on the DensePose COCO dataset [Guler *et al.*, 2018]. See the appendix for the relevant results.

**Metrics.** For instance-aware human parsing performance, we employ the Average Precision based on part ( $AP^p$ ) [Zhao *et al.*, 2018] for multi-human parsing evaluation, which calculates mean part instance level pixel IoU of different semantic part categories within a human instance to determine if the human

instance is a true positive. We choose the  $AP^p_{50}$  and  $AP^p_{vol}$  as the evaluation metrics. The former defines the instance whose IoU is larger than the threshold of 0.5 as the positive, and the latter is the average of  $AP^p$  in the IoU threshold ranging from 0.1 to 0.9 in increments of 0.1. In addition, we also report the official metric, Percentage of Correctly parsed semantic Parts (PCP) [Zhao *et al.*, 2018].

**Implementation Details.** We implement our SMP based on mmdetection [Chen *et al.*, 2019] on a server with 8 NVIDIA Tesla V100 GPUs and 32GB memory per card. We adopt ResNet-50, ResNet-101 and FPN as backbone and neck in all architectures, each of which is trained end-to-end. A mini-batch involves 32 images. We use scale jitter where the shorter image side is randomly sampled from 640 to 800 pixels, following SOLO [Wang *et al.*, 2019]. We regard 12 epochs as  $1 \times$  training with an initial learning rate of 0.02, which is then divided by 10 at 9-th and 11-th epoch. We use  $3 \times$  training schedule whose learning rate is divided by 10 at 27-th and 33-th epoch for performance comparison.

**Inference.** Since this work is the first attempt that the single-stage method is applied to the instance-aware multi-human parsing task, we design an inference scheme dedicated to the single-stage method. Fig.4 shows the visual intermediate outputs of SMP. Since SMP represents the instance parsing as point sets and center-based offsets, we utilize maxpooling to remove the redundant human barycenters in  $H_{\text{Center}}$ , and obtain the offsets to get the positions  $(C_x, C_y)$  of the part instances corresponding to the human instance. Before getting the positions of part instances by offset, we resize  $H_{\text{Part}}$  and  $F_{\text{Kernel}}$  to the same size as  $M_{\text{Offset}}$ . We judge whether the human instance contains the part instance by comparing the confidence value in  $H_{\text{Part}}$  on the position pointed by the offset vectors with the threshold. And the redundant part instances will be filtered out. After localizing the part instances, We

Table 2: The comparable results on MHPv2.0 dataset. \* denotes the backbone with DCNv2 [Zhu *et al.*, 2018]. For better comparison, we mark the state-of-the-art methods under the same settings in gray. RP-RCNN with 150 epochs is the best available pretrained model which is also used in the qualitative comparison (RP-RCNN is inferior than AIParsing with ResNet-50 and 75-epoch training strategy).

Methods	Backbone	Extra data	Epoch	$AP_{50}^p$	$AP_{vol}^p$	PCP <sub>50</sub>
<b>Top-Down</b>						
M-RCNN (He et al. TPAMI'17.)	ResNet-50	-	-	14.9	33.9	25.1
P-RCNN (Yang et al. CVPR'18.)	ResNeXt-101	Pose	75	30.2	41.8	44.2
M-CE2P (Ruan et al. AAAI'19.)	ResNet-101	-	150	34.5	42.7	43.7
SNT (Ji et al. ECCV'19.)	ResNet-101	-	-	34.4	42.5	43.5
RP-RCNN (Yang et al. ECCV'20.)	ResNet-50	Pose	75	40.5	45.2	39.2
RP-RCNN (Yang et al. ECCV'20.)	ResNet-50	Pose	150	45.3	46.8	43.8
AIParsing (Zhang et al. TIP'22.)	ResNet-50	-	75	41.1	45.9	45.3
AIParsing (Zhang et al. TIP'22.)	ResNet-101	-	75	43.2	46.6	47.3
<b>Bottom-up</b>						
PGN (Gong et al. ECCV'18.)	ResNet-101	-	-	17.6	35.5	26.9
MHParse (Li et al. TOMM'21.)	ResNet-101	-	-	18.0	36.1	27.0
NAN (Zhao et al. IJCV'20.)	-	-	~80	25.1	41.8	32.3
DSPF (Zhou et al. CVPR'21.)	ResNet-101	Pose	150	39.0	44.3	42.3
<b>Single-stage</b>						
SMP (Ours)	ResNet-101	-	12	38.4	44.8	43.0
SMP (Ours)	ResNet-101*	-	12	43.0	46.4	45.5
SMP (Ours)	ResNet-50	-	36	42.0	46.1	46.1
SMP (Ours)	ResNet-50*	-	36	46.9	48.0	50.1
SMP (Ours)	ResNet-101	-	36	45.3	47.6	48.5
SMP (Ours)	ResNet-101*	-	36	<b>47.1</b>	<b>48.2</b>	<b>51.5</b>

could be obtained their corresponding convolution kernel from  $F_{Kernel}$ , and the final parsing result can be obtained by convolution operation.

## 4.2 Main Results

**Comparison with State-of-the-art Models.** Tab. 2 reports the results of comparison with the state-of-the-art top-down and bottom-up methods on MHPv2.0 val set. Our method is significantly superior to the existing top-down and bottom-up model in all instance-aware metrics with less training epochs and without the aid of additional ground truth other than mask. For other methods, the longer learning schedule is set up to 150 epochs while our method sets 36 epochs as  $3 \times$  learning schedule. With deformable convolution in backbone, SMP achieves 47.1%, 48.2% and 51.5% on  $AP_{50}^p$ ,  $AP_{vol}^p$  and PCP<sub>50</sub>. Based on ResNet-101, ours outperforms by +2.1, +1.0 and + 1.2 than AIParsing [Zhang *et al.*, 2022a] with less training time. Compared with DSPF [Zhou *et al.*, 2021], which is the SOTA model in bottom-up methods, SMP outperforms on all metrics by +6.3 on  $AP_{50}^p$ , +3.3 on  $AP_{vol}^p$  and +6.2 on PCP<sub>50</sub>.

**Qualitative Results.** As shown in Fig. 5, we visualize some comparison results with the top-down method RP-RCNN [Yang *et al.*, 2020] (RP-RCNN performs best in the open-source works). We present four different images belonging to different scenarios: normal, occlusion, crowd, and multi-scales. In the normal and the multi-scales images, we can find that SMP performs better in the parsing of small parts, *e.g.*, half of the left foot, and the edge processing of the instance. For the occlusion image, RP-RCNN has an error detection, in which two people are detected as one. Crowd pictures are challenging for the multi-human parsing task. Too many people make it difficult for the part head respond and easy to arise miss detection. However, compared with RP-

Table 3: The results of ablative study on MHPv2.0 val set.

RFR	MIR	$AP_{50}^p$	$AP_{vol}^p$	PCP <sub>50</sub>
		42.5	45.8	44.1
	✓	42.8 (+0.3)	46.1 (+0.3)	44.7 (+0.6)
✓		42.6 (+0.1)	46.3 (+0.5)	44.3 (+0.2)
✓	✓	43.0 (+0.5)	46.4 (+0.6)	45.5 (+1.4)

RCNN, SMP has no problem of ROI detection error, and its ability to distinguish categories is also superior.

## 4.3 Empirical Experiments

**Ablative Experiments.** In this section, we performed ablative study on Refined Feature Retain module and Mask of Interest Reclassify module. The model in the ablative study was trained by the short learning schedule (the number of epochs is 12), and the backbone of the model is ResNet-101 with the deformable convolution [Zhu *et al.*, 2018]. From Tab. 3, we can see that the performance of the model is enhanced with the addition of RFR module and MIR module. From the results of ablation experiments, it can be seen that after adding RFR module to the model, all metrics have been improved. As a separate module added after the model, MIR module can be quickly trained which is trained with 6 epochs. We utilize the Focal Loss [Lin *et al.*, 2017] for supervision.

**Running Time Analysis.** In addition, we also conducted runtime analysis with the open-source reproducible methods. The experiment is conducted on one V100 GPU. The existing top-down and bottom-up methods all take multi-stage paradigms, leading to computational redundancy. From Tab. 4, we can see that the speed of our method is far faster than the bottom-up methods, and is about 1 time faster than the fastest top-down method Parsing R-CNN [Yang *et al.*, 2018]. (SMP reduces

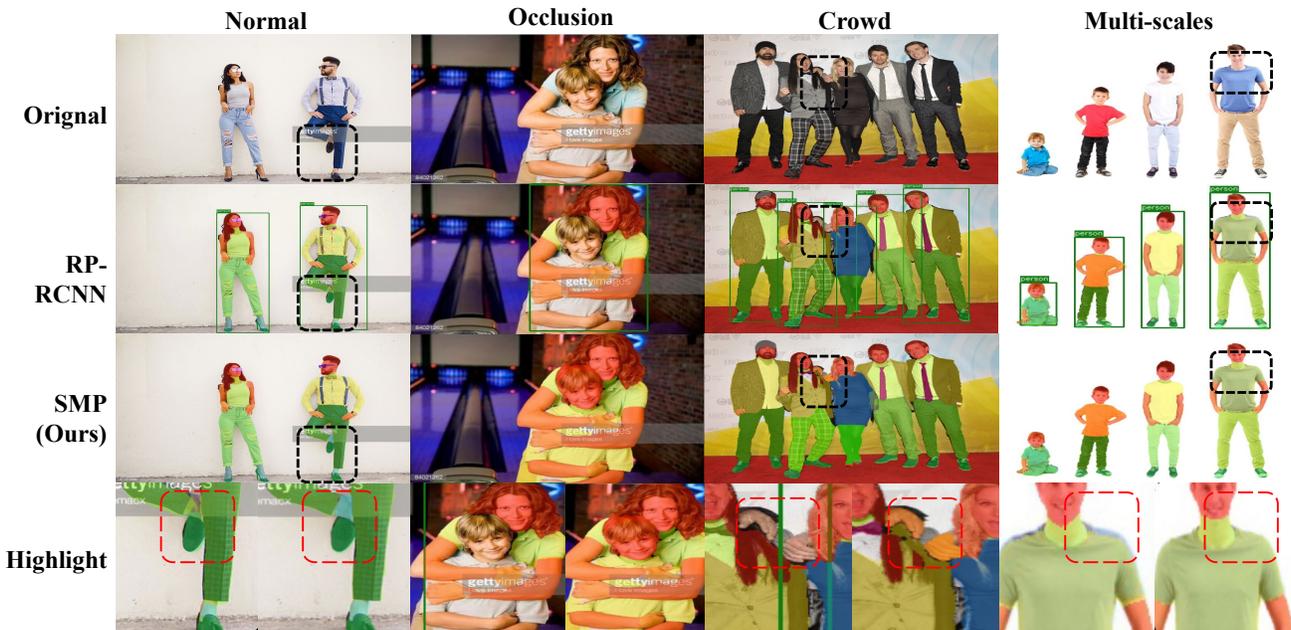


Figure 5: The visualization results from SMP including instances with different scales and quantity. In the highlight row, the left is RP-RCNN, the right is SMP. Best viewed in color.

Table 4: The runtime analysis on MHPv2.0 val set with average 2.6 people per image. The blue values are the decreases in inference time compared with the fastest open-source reproducible model. \* denotes the backbone with DCNv2 [Zhu *et al.*, 2018].

	Methods	Inference time (ms)
Top-down	SNT [Ji <i>et al.</i> , 2019]	3546
	M-CE2P [Ruan <i>et al.</i> , 2019]	1023
	P-RCNN [Yang <i>et al.</i> , 2018]	256
	RP-RCNN [Yang <i>et al.</i> , 2020]	341
Bottom-up	MHPParser [Yang <i>et al.</i> , 2020]	1224
	NAN [Zhao <i>et al.</i> , 2018]	997
	PGN [Gong <i>et al.</i> , 2018]	524
Single-stage	SMP (Ours, ResNet-50*)	124 (↓ 102.5%)
	SMP (Ours, ResNet-101*)	132 (↓ 91.2%)

the inference time by 102.5% and 91.2% with ResNet-50 and ResNet-101 respectively).

**Impact of Grid Numbers.** We also conduct additional experiments about grids number and hyperparameter  $\epsilon$  under the 12-epoch training schedule. The setting of grids number in Tab. 1 and  $\epsilon = 0.2$  can stably generate 2-6 positive labels for each instance. The above setting is effective for learning about semantic features of part head due to the consistency of instance scales. As the results in Tab. 5, under the grid numbers with the  $2\times$  size, the performance of SMP is improved by about 1% in all indicators, but the running speed is reduced by about 20%. While smaller grid number produces too rough features, resulting in poorer effect. For the hyperparameter  $\epsilon$ , its value is proportional to the number of positive grids generated by each instance. Too large  $\epsilon$  is prone to overlap and occlusion of center regions between instances. Too small  $\epsilon$  will result in too small response area on the feature map, and the generalization ability of the model is limited. Finally,

Table 5: The impact of grid numbers and center region on MHPv2.0 val set.  $1\times$  denotes the grids setting which is the same as Tab. 1.  $0.5\times$ ,  $2\times$  denote the grid numbers multiplied by the corresponding coefficients.  $\epsilon$  is the hyperparameter controlling the scale of center region. FPS indicates the inference speed.

Grids	$\epsilon$	$AP_{50}^p$	$AP_{vol}^p$	PCP <sub>50</sub>	FPS
$0.5\times$	0.2	26.6	36.7	31.8	8.45
$0.5\times$	0.1	25.7	36.3	29.9	8.45
$1\times$	0.2	43.0	46.4	45.5	7.47
$1\times$	0.1	41.9	45.8	44.0	7.47
$2\times$	0.2	44.1	47.1	47.3	5.94
$2\times$	0.1	44.0	47.1	46.2	5.94

we choose  $1\times$  grid numbers and  $\epsilon = 0.2$  to obtain the best trade-off between accuracy and speed.

## 5 Conclusion

This paper proposes to understand humans with point sets and center-based offsets, which results in a new framework, namely SMP, to solve the instance-aware multi-human parsing task in a single stage. Specifically, the feature of points in barycenters of human parts is utilized to generate the masks of the part instances. The offsets from the human center to part barycenters are used to unify the human instance. In order to enhance the representation of instance features for classification, we propose the Refined Feature Retain (RFR) module, which can utilize mask features to generate mask attention to guide feature extraction. For the problem of fashion classification errors with high inter-class similarity, we propose the Mask of Interest Reclassify (MIR) module, which employs the generated mask as the region of interest to refine the classification result. SMP has the advantages of fast inference, high accuracy and concise pipeline, which contribute to the human-centered research fields.

## References

- [Bolya *et al.*, 2019] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. Yolact: Real-time instance segmentation. *ICCV*, 2019.
- [Cao *et al.*, 2022] Gaofeng Cao, Fei Zhou, Han Yan, Anjie Wang, and Leidong Fan. Kpn-mfi: A kernel prediction network with multi-frame interaction for video inverse tone mapping. *IJCAI*, 2022.
- [Chen *et al.*, 2019] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyi Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.
- [Chen *et al.*, 2022] Zehui Chen, Zhenyu Li, Shiquan Zhang, Liangji Fang, Qinghong Jiang, Feng Zhao, Bolei Zhou, and Hang Zhao. Autoalign: Pixel-instance feature aggregation for multi-modal 3d object detection. *IJCAI*, 2022.
- [Dai *et al.*, 2016] Jifeng Dai, Kaiming He, Yi Li, Shaoqing Ren, and Jian Sun. Instance-sensitive fully convolutional networks. *ECCV*, 2016.
- [Dong *et al.*, 2022a] Hexin Dong, Zifan Chen, Mingze Yuan, Yutong Xie, Jie Zhao, Fei Yu, Bin Dong, and Li Zhang. Region-aware metric learning for open world semantic segmentation via meta-channel aggregation. *IJCAI*, 2022.
- [Dong *et al.*, 2022b] Zhangfu Dong, Yuting He, Xiaoming Qi, Yang Chen, Huazhong Shu, Jean-Louis Coatrieux, Guanyu Yang, and Shuo Li. Mnet: Rethinking 2d/3d networks for anisotropic medical image segmentation. *IJCAI*, 2022.
- [Geng *et al.*, 2021] Zigang Geng, Ke Sun, Bin Xiao, Zhaoxiang Zhang, and Jingdong Wang. Bottom-up human pose estimation via disentangled keypoint regression. *CVPR*, 2021.
- [Girshick, 2015] Ross Girshick. Fast r-cnn. *arXiv: Computer Vision and Pattern Recognition*, 2015.
- [Gong *et al.*, 2018] Ke Gong, Xiaodan Liang, Yicheng Li, Yimin Chen, Ming Yang, and Liang Lin. Instance-level human parsing via part grouping network. *ECCV*, 2018.
- [Guler *et al.*, 2018] Riza Alp Guler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. *CVPR*, 2018.
- [He *et al.*, 2016] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CVPR*, 2016.
- [He *et al.*, 2017] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. *TPAMI*, 2017.
- [He *et al.*, 2022] Wenbin He, William Surmeier, Arvind Kumar Shekar, Liang Gou, and Liu Ren. Self-supervised semantic segmentation grounded in visual concepts. *IJCAI*, 2022.
- [Ji *et al.*, 2019] Ruyi Ji, Dawei Du, Libo Zhang, Longyin Wen, Yanjun Wu, Chen Zhao, Feiyue Huang, and Siwei Lyu. Learning semantic neural tree for human parsing. *ECCV*, 2019.
- [Karras *et al.*, 2018] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. *TPAMI*, 2018.
- [Lei *et al.*, 2022] Jin Lei, Chenyang Xu, Xiaojuan Wang, Yabo Xiao, Yandong Guo, Xuecheng Nie, and Jian Zhao. Single-stage is enough: Multi-person absolute 3d pose estimation. *CVPR*, 2022.
- [Li *et al.*, 2017] Jianshu Li, Jian Zhao, Yunchao Wei, Congyan Lang, Yidong Li, and Jiashi Feng. Towards real world human parsing: Multiple-human parsing in the wild. 2017.
- [Li *et al.*, 2021] Jie Li, Laiyan Ding, and Rui Huang. Imenet: Joint 3d semantic scene completion and 2d semantic segmentation through iterative mutual enhancement. *IJCAI*, 2021.
- [Lin *et al.*, 2016] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. *CVPR*, 2016.
- [Lin *et al.*, 2017] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. *IEEE TPAMI*, 2017.
- [Liu *et al.*, 2018] Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution. *NeurIPS*, 2018.
- [Long *et al.*, 2014] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. *IEEE TPAMI*, 2014.
- [McNally *et al.*, 2021] William McNally, Kanav Vats, Alexander Wong, and John McPhee. Rethinking keypoint representations: Modeling keypoints and poses as objects for multi-person human pose estimation. *CVPR*, 2021.
- [Milletari *et al.*, 2016] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. *ICCV*, 2016.
- [Neven *et al.*, 2019] Davy Neven, Bert De Brabandere, Marc Proesmans, and Luc Van Gool. Instance segmentation by jointly optimizing spatial embeddings and clustering bandwidth. *CVPR*, 2019.
- [Ni *et al.*, 2020] Zhen-Liang Ni, Gui-Bin Bian, Guan'an Wang, Xiao-Hu Zhou, Zeng-Guang Hou, Xiao-Liang Xie, Zhen Li, and Yu-Han Wang. Barnet: Bilinear attention network with adaptive receptive fields for surgical instrument segmentation. *IJCAI*, 2020.
- [Nie *et al.*, 2019] Xuecheng Nie, Jiashi Feng, Jianfeng Zhang, and Shuicheng Yan. Single-stage multi-person pose machines. *ICCV*, 2019.
- [Redmon and Farhadi, 2018] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv: Computer Vision and Pattern Recognition*, 2018.
- [Ruan *et al.*, 2019] Tao Ruan, Ting Liu, Zilong Huang, Yunchao Wei, Shikui Wei, and Yao Zhao. Devil in the details:

- Towards accurate single and multiple human parsing. *AAAI*, 2019.
- [Sofiiuk *et al.*, 2019] Konstantin Sofiiuk, Olga Barinova, and Anton Konushin. Adaptis: Adaptive instance selection network. *ICCV*, 2019.
- [Tian *et al.*, 2020] Zhi Tian, Chunhua Shen, and Hao Chen. Conditional convolutions for instance segmentation. *arXiv: Computer Vision and Pattern Recognition*, 2020.
- [Wang *et al.*, 2019] Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. Solo: Segmenting objects by locations. *ECCV*, 2019.
- [Wang *et al.*, 2020] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. *NeurIPS*, 2020.
- [Yang *et al.*, 2018] Lu Yang, Qing Song, Zhihui Wang, and Ming Jiang. Parsing r-cnn for instance-level human analysis. *CVPR*, 2018.
- [Yang *et al.*, 2019] Brandon Yang, Gabriel Bender, Quoc V. Le, and Jiquan Ngiam. Condconv: Conditionally parameterized convolutions for efficient inference. *NeurIPS*, 2019.
- [Yang *et al.*, 2020] Lu Yang, Qing Song, Zhihui Wang, Mengjie Hu, Chun Liu, Xin Xueshi, Jia Wenhe, and Songcen Xu. Renovating parsing r-cnn for accurate multiple human parsing. *ECCV*, 2020.
- [Ying *et al.*, 2021] Hui Ying, Zhaojin Huang, Shu Liu, Tianjia Shao, and Kun Zhou. Embedmask: Embedding coupling for instance segmentation. *IJCAI*, 2021.
- [Zhang *et al.*, 2021] Zihao Zhang, Lei Hu, Xiaoming Deng, and Shihong Xia. Sequential 3d human pose estimation using adaptive point cloud sampling strategy. *IJCAI*, 2021.
- [Zhang *et al.*, 2022a] Sanyi Zhang, Xiaochun Cao, Guo-Jun Qi, Zhanjie Song, and Jie Zhou. Aiparsing: Anchor-free instance-level human parsing. 2022.
- [Zhang *et al.*, 2022b] Zhengbo Zhang, Chunluan Zhou, and Zhigang Tu. Distilling inter-class distance for semantic segmentation. *IJCAI*, 2022.
- [Zhao *et al.*, 2017] Jian Zhao, Jianshu Li, Xuecheng Nie, Fang Zhao, Yunpeng Chen, Zhecan Wang, Jiashi Feng, and Shuicheng Yan. Self-supervised neural aggregation networks for human parsing. *CVPR*, 2017.
- [Zhao *et al.*, 2018] Jian Zhao, Jianshu Li, Yu Cheng, Terence Sim, Shuicheng Yan, and Jiashi Feng. Understanding humans in crowded scenes: Deep nested adversarial learning and a new benchmark for multi-human parsing. *ACM MM*, 2018.
- [Zhao *et al.*, 2020] Jian Zhao, Jianshu Li, Hengzhu Liu, Shuicheng Yan, and Jiashi Feng. Fine-grained multi-human parsing. *IJCV*, 2020.
- [Zhou *et al.*, 2021] Tianfei Zhou, Wenguan Wang, Si Liu, Yi Yang, and Luc Van Gool. Differentiable multi-granularity human representation learning for instance-aware human semantic parsing. *CVPR*, 2021.
- [Zhu *et al.*, 2018] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable convnets v2: More deformable, better results. *CVPR*, 2018.

# Supplementary Material for Single-stage Multi-human Parsing via Point Sets and Center-based Offsets

Paper ID 341

1 This supplementary material provides the detailed infor-  
2 mation that the paper cannot show due to the length limit,  
3 including: 1) The performance comparison of our single stage  
4 method on the Densepose COCO dataset [?]. 2) Visualization  
5 of multi-human parsing results compared with groundtruth on  
6 MHPv2.0 [?]. 3) Some detailed information during training  
7 and inference.



Figure 1: The visualization of the results on DensePose COCO.

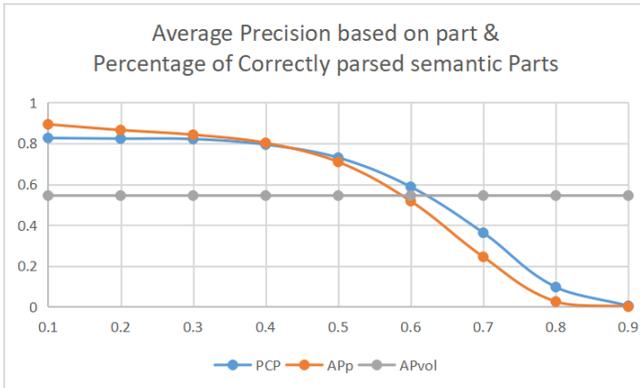


Figure 2: The values of  $AP^p$  and PCP at different thresholds.

Table 1: The results performed on Densepose-COCO dataset.

Methods	$AP_{50}^p$	$AP_{vol}^p$	PCP <sub>50</sub>
<i>Top-Down</i>			
P-RCNN (Yang et al. CVPR'18.)	43.5	53.1	51.8
M-CE2P (Ruan et al. AAAI'19.)	43.7	52.9	51.2
RP-RCNN (Yang et al. ECCV'20.)	48.5	54.4	51.1
<i>Bottom-Up</i>			
PGN (Gong et al. ECCV'18.)	23.4	35.9	32.5
NAN (Zhao et al. IJCV'20.)	37.6	48.3	43.9
DSPF (Zhou et al. CVPR'21.)	49.7	<b>54.7</b>	52.8
<i>Single-stage</i>			
SMP (Ours)	<b>70.9</b>	54.4	<b>73.0</b>

## 1 The Performance on Densepose-COCO dataset

We apply our method on the DensePose-COCO dataset.  $AP_{50}^p$ ,  $AP_{vol}^p$  and PCP<sub>50</sub> are also selected as metrics and the results are compared with other state of the art (SOTA) top-down and bottom-up methods. The backbone selected in the experiment is ResNet101 with deformable convolution [?]. Since the instance scales of the COCO dataset are small [?], we change the size of the grid in model to twice the original one.

**Dataset.** DensePose-COCO [?] has 27,659 images (26,151/1,508 for train/test splits) gathered from COCO [?]. It annotates 14 human parts and 17 keypoints.

**Metrics.** As shown in paper,  $AP^p$  [?] calculates mean part instance level pixel intersection of union (IoU) of different semantic part categories within a human instance to determine if the human instance is a true positive. We choose the  $AP_{50}^p$  and  $AP_{vol}^p$  as the evaluation metrics. The former defines the instance whose IoU is larger than the threshold of 0.5 as the positive, and the latter is the average of  $AP^p$  in the IoU threshold range from 0.1 to 0.9 in increments of 0.1. In addition, we also report the official metric, Percentage of Correctly parsed semantic Parts (PCP) [?].

**Results.** As shown in Tab. 1, our method achieves comparable performance with the SOTA methods on the DensePose-COCO. In  $AP_{50}^p$ , our model achieves 70.9%, which is 21.2% higher than the previous optimal model DSPF [?]. In PCP<sub>50</sub>, SMP reaches 73.0%, which is also 20.2% higher. For  $AP_{vol}^p$ , our method reaches 54.4%, which is only 0.3% lower than DSPF [?], but higher than other SOTA methods [?]. From the



Figure 3: The visual results on MHPv2.0 including separated instances.

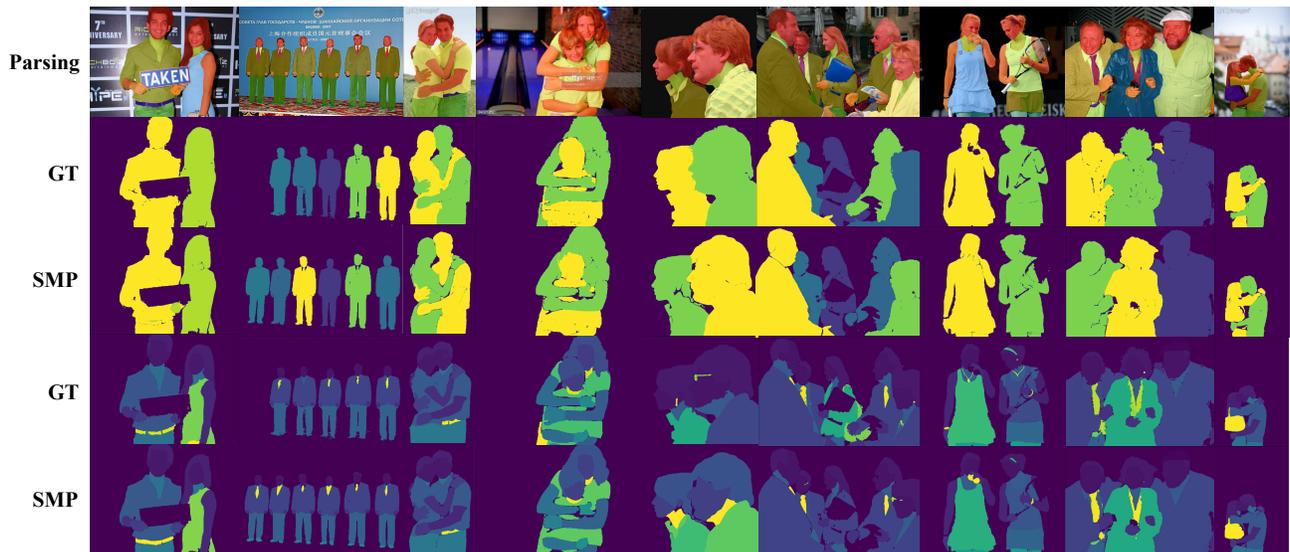


Figure 4: The visual results on MHPv2.0 compared with groundtruth.

37 results, we can see that our method has excellent performance  
 38 on the Average Precision with 0.5 threshold, and the average  
 39 performance is basically the same as the current optimal  
 40 method. This indicates that our method has a higher lower  
 41 limit than other methods, but it performs undesirably in the  
 42 upper limit.

43 As shown in Fig. 2, when the threshold is smaller than 0.5,  
 44 the  $AP^p$  drops very slowly, while the threshold is larger than  
 45 0.5, the  $AP^p$  drops rapidly. It means that our method can get a  
 46 good result for most inputs. However, it is difficult to generate  
 47 multi-human parsing whose IoU are above 0.8.

## 48 2 Visualization of Multi-human Parsing 49 Results

50 In this appendix, we will show the visualization results of  
 51 the model in detail, including the display by instance and the  
 52 comparison with the ground truth. As shown in Fig. 3, We  
 53 show the prediction results of different instances in the same  
 54 image. In the first line, we display the results of combined  
 55 multi-human parsing. In the second line, we provide the human  
 56 parsing separate by instances. SMP does need neither  
 57 bounding boxes [?] to locate nor keypoints [?] or edge [?] to  
 58 match. The output results of SMP are directly human parsing  
 59 of individual instances, as shown in the second line of Fig. 3.

60 In Fig. 4, we provide the visualization of groundtruth and  
 61 prediction of human instance level segmentation and those of  
 62 part instance segmentation. Our method has achieved good re-  
 63 sults in the segmentation of two different granularity instances.  
 64 Even in the face of complex occlusion or long-distance body  
 65 parts, SMP can still properly match the human body and part  
 66 instances.

## 67 3 The Detailed Setting in Training and 68 Inference

69 Some hyperparametric settings involved in training and rea-  
 70 soning will be introduced in this appendix.

71 **Training.** During training, our initial learning rate is set as  
 72 0.02, and the batch size per card is 4. The total batch size on 8  
 73 NVIDIA Tesla V100 is 32.

74 **Inference.** During inference, we utilize  $\theta_c$  as the threshold  
 75 for localizing the barycenters of the human body in center  
 76 heatmap  $H_C$  with a value of 0.1. The threshold for localizing  
 77 part instance  $\theta_p$  has a value of 0.3. To avoid similar instances  
 78 responding at the same location, we applied matrix NMS [?]  
 79 to filter the prediction of instances with large overlap with a  
 80 threshold  $\theta_u$  whose value is 0.1.

81 And in the last of the appendix, we also provide the config  
 82 file of our method in mmdetection style. It is very easy to  
 83 understand, and is divided to model settings, dataset settings,  
 84 optimizer settings and learning policy.

Listing 1: config.py

```

1 # model settings
2 model = dict(
3     type='SingleStageInsParsingDetector',
4     backbone= dict(
5         type='ResNet',
6         depth=101,
7         num_stages=4,
8         out_indices=(0, 1, 2, 3),
9         frozen_stages=1,
10        init_cfg= dict(
11            type='Pretrained', checkpoint='
12            torchvision://resnet101'),
13            style='pytorch'),
14        neck= dict(
15            type='FPN',
16            in_channels=[256, 512, 1024, 2048],
17            out_channels=256,
18            start_level=0,
19            num_outs=5),
20        mask_head= dict(
21            type='SMP_head',
22            num_classes=59,
23            in_channels=256,
24            stacked_convs=4,
25            seg_feat_channels=512,
26            strides=[8, 8, 16, 32, 32],
27            scale_ranges=((1, 96), (48, 192), (96,
28                384), (192, 768), (384, 2048)),
29            sigma=0.2,
30            num_grids=[40, 36, 24, 16, 12],
31            ins_out_channels=256,
32            enable_center=True,
33            enable_offset=True,
34            enable_moi = True,
35            mask_feature_head= dict(
36                in_channels=256,
37                feat_channels=128,
38                start_level=0,
39                end_level=3,
40                out_channels=256,
41                mask_stride=4,
42                norm_cfg= dict( type='GN', num_groups
43                    =32, requires_grad=True)),
44            loss_ins= dict(
45                type='DiceLoss',
46                use_sigmoid=True,
47                loss_weight=3.0),
48            loss_cate= dict(
49                type='FocalLoss',
50                use_sigmoid=True,
51                gamma=2.0,
52                alpha=0.25,
53                loss_weight=1.0),
54            loss_center= dict(
55                type='FocalLoss',
56                use_sigmoid=True,
57                gamma=2.0,
58                alpha=0.25,
59                loss_weight=1.0),
60            loss_offset= dict(
61                type='SmoothL1Loss',
62                beta=1. / 9,
63                reduction='mean',
64
  
```

```

149 60         loss_weight=10.0)),
150 61     train_cfg = dict(),
151 62     test_cfg = dict(
152 63         nms_pre=500,
153 64         ctr_score=0.1,
154 65         cate_score_thr=0.3,
155 66         mask_thr=0.5,
156 67         cate_update_thr=0.1,
157 68         update_thr=0.1,
158 69         kernel='gaussian', # gaussian/linear
159 70         sigma=2.0,
160 71         max_per_img=100,
161 72         debug_heatmap=False)
162 73     )
163 74
164 75
165 76 # dataset settings
166 77 dataset_type = 'MHP'
167 78 data_root = 'LV-MHP-v2/'
168 79 img_norm_cfg = dict(
169 80     mean=[123.675, 116.28, 103.53],
170 81     std=[58.395, 57.12, 57.375],
171 82     to_rgb=True)
172 83 train_pipeline = [
173 84     dict( type='LoadImageFromFile_mhp' ),
174 85     dict( type='LoadAnnotations_mhp' ,
175 86         with_mask=True, with_seg=True ,
176 87         with_mhp_parsing=True,
177 88         with_keypoints=True),
178 89     dict( type='Resize_Parsing' ,
179 90         img_scale=[(1333, 800), (1333, 768),
180 91                   (1333, 736),
181 92                   (1333, 704), (1333, 672),
182 93                   (1333, 640)],
183 94         multiscale_mode='value' ,
184 95         keep_ratio=True),
185 96     dict( type='RandomFlip_Parsing' ,
186 97         flip_ratio=0.5, dataset='mhp' ),
187 98     dict( type='Normalize' , **img_norm_cfg),
188 99     dict( type='Pad_Parsing' , size_divisor=32),
189 100    dict( type='DefaultFormatBundle_Parsing' ),
190 101    dict( type='Collect' , keys=['img' ,
191 102                                'gt_bboxes' ,
192 103                                'gt_labels' ,
193 104                                'gt_masks' ,
194 105                                'gt_parsing' ,
195 106                                'gt_semantic_seg' ,
196 107                                'gt_keypoints' ],
197 108        meta_keys=('filename' ,
198 109                  'ori_shape' ,
199 110                  'img_shape' ,
200 111                  'pad_shape' ,
201 112                  'scale_factor' ,
202 113                  'flip' ,
203 114                  'img_norm_cfg' ),)
204 115 ]
205 116 test_pipeline = [
206 117     dict( type='LoadImageFromFile_mhp' ),
207 118     dict(
208 119         type='MultiScaleFlipAug' ,
209 120         img_scale=(1333, 800),
210 121         flip=False,
211 122         transforms=[
212 123             dict( type='Resize_Parsing' ,
213 124                 keep_ratio=True),

```

```

214         dict( type='RandomFlip_Parsing' ,
215             dataset='mhp' ),
216         dict( type='Normalize' , **
217             img_norm_cfg),
218         dict( type='Pad_Parsing' ,
219             size_divisor=32),
220         dict( type='ImageToTensor' , keys=['
221             img' ]),
222         dict( type='Collect' , keys=['img' ]),
223     ]
224 ]
225 data = dict(
226     samples_per_gpu=4,
227     workers_per_gpu=4)
228 # optimizer
229 optimizer = dict( type='SGD' ,
230                 lr=0.02,
231                 momentum=0.9,
232                 weight_decay=0.0001)
233 optimizer_config = dict(
234     grad_clip=dict(
235         max_norm=35,
236         norm_type=2))
237 # learning policy
238 lr_config = dict(
239     policy='step' ,
240     warmup='linear' ,
241     warmup_iters=500,
242     warmup_ratio=0.02,
243     step=[8, 11])
244 runner = dict(
245     type='EpochBasedRunner' , max_epochs=12)

```