

LLA: Loss-aware Label Assignment for Dense Pedestrian Detection

Zheng Ge^{1,2}, Jianfeng Wang², Xin Huang¹, Songtao Liu², Osamu Yoshie¹
¹Waseda University, ²Megvii Technology

jokerzz@fuji.waseda.jp; wangjianfeng@megvii.com; koushin@toki.waseda.jp;
liusongtao@megvii.com; yoshie@waseda.jp

Abstract

Label assignment has been widely studied in general object detection because of its great impact on detectors' performance. However, none of these works focus on label assignment in dense pedestrian detection. In this paper, we propose a simple yet effective assigning strategy called Loss-aware Label Assignment (LLA) to boost the performance of pedestrian detectors in crowd scenarios. LLA first calculates classification (cls) and regression (reg) losses between each anchor and ground-truth (GT) pair. A joint loss is then defined as the weighted summation of cls and reg losses as the assigning indicator. Finally, anchors with top K minimum joint losses for a certain GT box are assigned as its positive anchors. Anchors that are not assigned to any GT box are considered negative. Loss-aware label assignment is based on an observation that anchors with lower joint loss usually contain richer semantic information and thus can better represent their corresponding GT boxes. Experiments on CrowdHuman and CityPersons show that such a simple label assigning strategy can boost MR by 9.53% and 5.47% on two famous one-stage detectors – RetinaNet and FCOS, respectively, demonstrating the effectiveness of LLA. The code is available at <https://github.com/Megvii-BaseDetection/LLA>.

1. Introduction

Pedestrian detection in crowd scenarios has attracted considerable attentions in the recent literature [27, 31, 23, 16] and applications (e.g., autonomous driving and video surveillance). It is widely used in many real-world scenarios where the density of people is high, i.e., airports, train stations, shopping malls etc. Compared to general object detection [15], target objects in crowd scenarios are more densely arranged, resulting in heavy occlusion between different objects. Although the performance of modern Convolutional Neural Network (CNN) based object detectors [14, 24, 20, 19] is growing rapidly, they usually suffer when applied to dense pedestrian detection (DPD). This

is mainly caused by following two issues: mis-classified occluded pedestrian and mis-placed detected results. Recent works mainly utilized additional information or regularization term to relieve these two problems. For example, Bibox [37] and R2-NMS [8] alleviate the first issue by introducing visible body annotations as extra supervisions. For the second issue, [27] imposes a novel regression penalty term on the misplacing predictions to tackle it. While these methods try to amend the poor predictions from the detectors, in this paper, we delve into the essential cause of these issues and find an important problem which has never been discussed in the literature of DPD – label assignment.

Anchor box is the basic processing unit in CNN based object detectors. The procedure in which anchors are assigned as positive or negative during training is called *label assignment*. For anchor-based detectors, the most common label assigning strategy usually utilizes Intersection-over-Union (IoU) between anchors and a ground-truth (GT) bounding boxes. For example, in RetinaNet [14], if an anchor's IoU with a certain GT box exceeds 0.5, this anchor is considered as *positive* and assigned to that GT box. Otherwise, anchors are assigned as *negative* or *ignore* according to its maximum IoU between GT boxes. For anchor-free detectors, we first denote the dense sample locations on feature maps as **anchor points**. In a typical anchor-free detector – FCOS [24], a group of anchor points are directly assigned as positive if they fall into a square region in the center of a GT box. Such hand-crafted assigning strategies tend to assign positive anchors near the geometric center of that GT box. While they are proved to work effectively in many famous detectors [14, 24, 20, 19], things are different when they come to pedestrian detection in crowd scenarios. If a person is heavily occluded, his/her geometric center may fall onto other's body, which will lead to inconsistency between the features of sampled points and their corresponding GT boxes. These twisty samples interfere with the training of detectors are certainly one of the main reasons for the mis-classifying and misplacing issues in DPD.

Several recent works [32, 9, 35] try to make the procedure of label assignment more adaptive for general object

detection. A typical pipeline of them follows: 1) Constructing a bag of positive candidate anchors for each GT. 2) Calculating a certain metric *e.g.* IoU [32], score function [9] or likelihood [35] for each GT’s candidate anchors. 3) Applying statistical tools or hard thresholds on the calculated metric to define positive and negative anchors. Although they can adaptively define what is positive and negative according to the network’s prediction, they share a common prior that a set of positive candidate anchors need to be constructed in advance, which will be based on hand-crafted rule – IoU to ensure that reasonable statistic values can be acquired. Such a constraint limits the positive regions near the geometric center of each GT box. As we stated above, in dense pedestrian detection, the anchor boxes/points near the center of a GT box could be improper and even harmful for being positives if the corresponding person is heavily occluded. These methods fail to take this case into consideration, making their assigning results sub-optimal in crowd scenarios.

To break the limit of current existing label assigning strategies, we propose an extremely simple but effective label assigning strategy called **Loss-aware Label Assignment (LLA)** for dense pedestrian detection. First, LLA calculates *cls* and *reg* losses between each anchor and GT pair. Then, the weighted summation of *cls* and *reg* losses is defined as the joint loss to estimate how well can one anchor learn one GT box. To help model converge better, we impose an “in box” term $C^{in\text{box}}$ into the loss term as a minimal constraint. Specifically, if the center of an anchor box fails to fall into any GT box, we will add a constant punishment term $C^{in\text{box}}$ ($C^{in\text{box}} > 0$) in its joint loss. Otherwise, $C^{in\text{box}}$ is set to 0. Finally, anchors with top K minimum joint losses for a certain GT box are assigned as its positive anchors. Anchors that are not assigned to any GT box are considered negative. Noted that in LLA, label assignment fully abandons the scale prior as proposed in FPN [13], and center prior as utilized in FreeAnchor/FCOS/ATSS. LLA defines positive anchors based on the model’s output, making LLA fully adaptive. LLA is proved to be **occlusion-aware** because heavily occluded regions tend to have a higher joint loss and thus are less likely to be assigned as positive. On CrowdHuman [23], LLA brings 9.53% and 5.47% improvements on MR when applied to RetinaNet and FCOS, respectively, demonstrating its effectiveness. Experiments on CityPersons [31] further reveals LLA’s capability on various pedestrian detection datasets.

2. Related Works

2.1. Occlusion Handling for Pedestrian Detection.

Recently, occlusion handling becomes a popular topic in the field of pedestrian detection. In heavy occlusion situations, the detector will get confused by the adjacent instance

which leads to an inaccurate boundaries regression. To solve this problem, Repulsion Loss [27] and OR-CNN [33] both impose additional penalty terms on the BBoxes that appear in the middle of two persons to force them to regress to the right person. Moreover, utilizing visible annotation as extra supervision is a common strategy to obtain more precision location information. Bi-box [37] adds a visible branch on Fast R-CNN to predict the full and visible body of a pedestrian at the same time. ATT [34] exploits the visible-region information as external guidance to handle various occlusion patterns in crowded situations. MGAN [18] forces the detector to focus on the visible regions of a pedestrian by adopting a novel attention branch to highlight the visible body region while suppressing the occluded part. What’s more, some researchers point out that in crowded scenario NMS may be trapped in a dilemma: a lower threshold of intersection over union (IoU) resulting in the miss of highly overlapped pedestrians while a higher IoU threshold naturally brings in more false positives. To solve this problem, Adaptive-NMS [16] proposes a subnet to predict the threshold for different anchors. R2-NMS [8] leverages the less occluded visible parts to remove the redundant boxes. PS-RCNN [5] utilizes two parallel R-CNN modules to detect slightly/none occluded and heavily occluded human instances in a divide-and-conquer manner. Different from all the existing works, our method handles the severe occlusion situation by a novel label assigning method, which neither requires additional annotation nor demands extra parameters.

2.2. Hand-crafted Label Assignment

Anchors are a set of pre-defined square boxes with different scale and aspect ratios which are densely assigned to each spatial location on feature maps. Traditional anchor-based object detectors [14, 20, 13, 1] assign labels for anchors based on hand-crafted IoU between anchors and GTs. Specifically, if an anchor’s IoU with a certain GT box exceeds a certain threshold (*e.g.* 0.5 for RetinaNet [14]), this anchor is defined as positive. The remaining anchors are either defined as *negative* or *ignore* according to its maximum IoU between GT boxes. Besides, extra scale constraint introduced by FPN [13] is imposed into these detectors [13, 1] to better handle the scale variations in objects. Instead of using pre-defined anchors, MetaAnchor [29] proposes an anchor function to provide anchors with dynamic shapes in both training and testing phase. Guided Anchoring [25] first learns a set of adjusted anchors in an anchor-free manner to better fit the shape of targets, and then design the subsequent anchor-based modules based on the learned anchors.

Anchor-free methods have drawn more and more attention recently due to its simple pipeline. FCOS [24] and FoveaBox [10] define anchors in the center region of targets as positive anchors. FSAF [41] further utilizes an on-

line Feature Selection Module to select the appropriate feature level for each GT. In keypoint-based anchor-free detectors [11, 39, 38, 30], only a single center point for each GT (anchor which is the closest to the center of that GT) is defined as positive, while other anchors are all negatives. Because the mechanism of keypoint-based detectors is much different from bounding box based detectors, they are out of the scope of this paper.

2.3. Dynamic Label Assignment

Recently, methods called dynamic label assignment are proposed to improve the training process of detectors. ATSS[32] introduces dynamic IoU thresholds by calculating the mean and standard deviation of IoU values on a set of pre-defined candidate anchors for each GT. FreeAnchor [35] designs a detection-customized likelihood that takes precision and recall into consideration to tackle the anchor-object matching problem. Specifically, anchors that have a higher likelihood are defined as positives. PAA [9] first designs a score function based on the classification and regression loss then applies one-dimensional GMM [21] for these calculated scores for each GT to choose the thresholds for separating positive and negative anchors. DeTR [3] and DeFCN [26] explore customized loss and quality terms as their indicators for one-to-one label assignment. Although one-to-one matching strategy is proved crucial to end-to-end detectors, it may not be the optimal choice for detectors followed by NMS.

Instead of solving the anchor-object matching problem directly, there are also a few works trying to re-weight the positives and negatives which can be categorized into generalized label assignment. PISA [2] first proposes two ranking strategies – IoU-HLR and Score-HLR to rank positive and negative proposals, respectively, to evaluate the importance of anchors, then forces the model focus on more important samples by giving them higher weights. Noisy Anchor [12] constructs a cleanliness score based on the detector’s outputs to re-weight anchors and soften classification labels. However, as we stated above, although they can adaptively define positive/negative anchors, they still limit the positive regions near the geometric center of each GT box, which could harm the detectors’ performance in DPD.

3. Method

3.1. Revisiting Label Assignment in General Object Detection

Let us take a brief look at how label assignment is conducted on two well-known one-stage detectors – RetinaNet and FCOS. Given an input image M , the ground-truth annotations are denoted as G , where a ground-truth box $g_i \in G$ is made up of a class label g_i^{cls} and a location g_i^{loc} . Due to the wide-spread of multi-scale feature pyramids network

(FPN), both scale and spatial constraints need to be considered during label assignment.

In RetinaNet, $a_j \in A$ stands for an **anchor box**. RetinaNet handles spatial and scale constraint simultaneously based on the Intersection-over-Union (IoU) matching rule. During training, a_j is assigned to GT g_i if $IoU(a_j, g_i^{loc}) > 0.5$, while a_j is defined as negative if $\forall g_i \in G, IoU(a_j, g_i^{loc}) < 0.4$. Anchors which are neither positive nor negative are ignored at that training step.

In FCOS, $a_j \in A$ stands for an **anchor point**. During training, a_j is assigned to GT g_i only if 1) a_j falls into the center area (within a fixed radius) of g_i . 2) a_j meets the pre-defined scale constraint introduced by FPN[13]. Anchor points which do not meet these two requirements are defined as negatives. Note that in FCOS, both spatial and scale constraints are explicitly imposed in the process of label assignment, making it less flexible.

3.2. Rethinking Label Assignment in Dense Pedestrian Detection

Label assigning strategies in RetinaNet and FCOS are based on a strong assumption that the geometric center is the most appropriate spatial location to represent objects. In statistics, this assumption may hold when we deal with objects with a large variety of categories. However, compared to general object detection, dense pedestrian detection are different in the following two aspects:

- Poses between different individuals vary a lot.
- Human body can be heavily occluded by others.

Given these two specificities, heuristic label assigning strategies adopted in RetinaNet or FCOS are no more suitable for dense pedestrian detection. In Fig. 1, we visualize some typical scenarios in which inappropriate label assignment is introduced by RetinaNet and FCOS. As can be seen in Fig. 1, when the human body exhibits unusual postures (*i.e.* dancing, playing sports, etc), part of positive anchors of both RetinaNet and FCOS may fall onto background regions. In these cases, the classifier will get confused and then learns improper decision boundaries. When people are very close to each other, their positive anchors will be severely intertwined where the regressor can hardly determine which target should be approached for each anchor. These examples indicate that in dense pedestrian detection, hand-crafted label assigning strategies are no more qualified.

Some recent works [32, 9, 35] have taken advantage of the outputs of models to perform label assignment, referred to as *dynamic label assignment*. The success of *dynamic label assignment* is based on an observation that anchors with lower loss/higher likelihood can better represent the corresponding instances. However, all these strategies need to

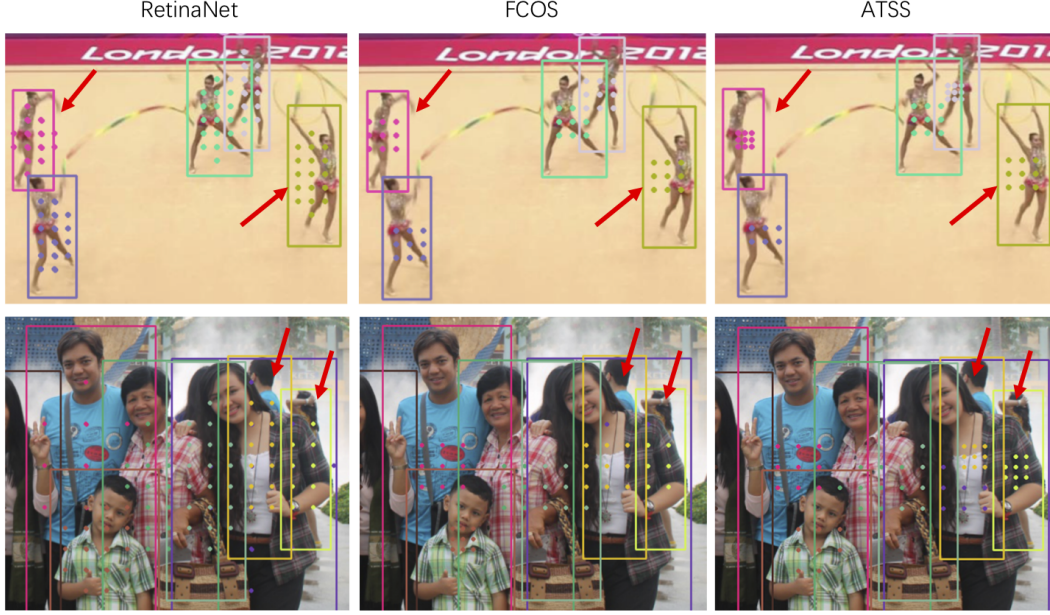


Figure 1. Inappropriate label assignment in RetinaNet, FCOS and ATSS [32]. Bounding boxes in this figure are GT annotations while dots are assigned positive anchors. **Red arrows highlight the typical inappropriate assignment.** Only FPN layers with the largest number of positive anchors are shown for better visualization.

first construct a positive candidate set of anchors for each object and then adaptively split positive and negative anchors according to the statistics from the candidate set. As the construction of the candidate set remains hand-crafted, the whole set is still near the geometric centers of their corresponding objects, which makes them *partially-dynamic*. As shown in Fig. 1, in ATSS, some of the positive anchors of a heavily occluded human instance, still fall into other instance’s body region, which indicates that partially-dynamic label assignments are sub-optimal solutions for dense pedestrian detection.

3.3. Loss-aware Label Assignment

Based on previous discussion, we propose **Loss-aware Label Assignment (LLA)**, a simple and fully-dynamic label assigning strategy for dense pedestrian detection which can be easily plugged into any anchor-based or anchor-free detectors. Without loss of generality, we extend the definition of $a_j \in A$ into an anchor box/point. Given an input image M , suppose there are J anchors and I GT annotations. In a single forward step, we can get score predictions $S(\theta, M) \in \mathbb{R}^{J \times N}$, where N is the number of classes and θ is the learnable parameters in detector, and get bounding box predictions $B(\theta, M) \in \mathbb{R}^{J \times 4}$. Unlike previous methods which only calculate losses between each anchor and its assigned GT, LLA calculates losses between all anchor-GT pairs, getting:

$$\begin{aligned} C^{cls} &= f^{cls}(G^{cls}, S(\theta, M)) \\ C^{reg} &= f^{reg}(G^{loc}, B(\theta, M)), \end{aligned} \quad (1)$$

where $C^{cls} \in \mathbb{R}^{I \times J}$ and $C^{reg} \in \mathbb{R}^{I \times J}$. G^{cls} and G^{loc} are ground-truth annotations for class and bounding box, respectively. f^{cls} denotes binary cross entropy or Focal Loss while f^{reg} can denote any of regression losses in SmoothL1 [6], IoU and GIoU [22] Loss. Then, a *Cost Matrix* can be formulated as:

$$C = C^{cls} + \lambda * C^{reg}, \quad (2)$$

where $C \in \mathbb{R}^{I \times J}$. Given the definition of C , C_{ij} represents the joint loss between a_j and g_i . The smaller C_{ij} is, the more possible anchor a_j is assigned to GT g_i . Hence, we select top K smallest values in each row of C and consider the corresponding anchor-GT pairs *matching*. However, in our experiments, we found that at very early training stage, LLA can hardly produce stable assignment results due to the under-fitting of the detector. To help model converge faster, we add a spatial prior – a_j can be assigned to g_i only if a_j (or the center of a_j) falls within the range of g_i^{loc} . Based on this prior, we introduce C^{inbox}

$$C_{ij}^{inbox} = \begin{cases} 0, & \text{if } a_j \text{ in } g_i^{loc} \\ +\infty, & \text{if } a_j \text{ not in } g_i^{loc}. \end{cases} \quad (3)$$

In our implementation, the term $+\infty$ is replaced by a

large positive value (e.g. 10^2). Then, the *Restricted Cost Matrix* can be formulated as:

$$C_r = C^{cls} + \lambda * C^{reg} + C^{inbox}. \quad (4)$$

After selecting the top K smallest values in C , finally, the assignment matrix $\pi_{ij} \in \{0, 1\}$ can be obtained:

$$\pi_{ij} = \begin{cases} 1, & \text{if } a_j \text{ matches } g_i \\ 0, & \text{if } a_j \text{ does not match } g_i. \end{cases} \quad (5)$$

Note that **if an anchor is assigned to multiple GTs simultaneously, we assign this anchor to GT with the smallest cost.** After label assigning, the detector is updated in the same as in RetinaNet and FCOS.

Compared to label assigning strategies with complex hand-crafted rules, LLA only leverages the minimal ‘‘in box’’ prior. Which FPN layer should each GT assigned to is automatically determined by LLA according to the feedback of the model’s outputs. Compared to partially-dynamic strategies that restrict the positive candidates near the center of each GT, LLA assigns anchors in a fully dynamic manner which helps LLA better handle the severe occlusion situations.

4. Experiments

In this section, we carry out heavy experiments on CrowdHuman [23] to illustrate the effectiveness of LLA. In CrowdHuman, there are 15000, 4370, and 5000 images in the training set, validation set, and testing set respectively. It provides three categories of bounding boxes annotations for each human instance: head bounding-box, human visible-body bounding-box and human full-body bounding-box. Because annotations of full-body are more crowded and challenging than visible-body, thus in this work, all of our experiments are conducted on full-body annotations. Further experimental results on CityPersons are also provided to prove LLA’s general applicability on other datasets.

4.1. Training Details

For CrowdHuman, the network structure in our experiment follows [14]. We resize the input images so that the short edge is 800 pixels while the long edge is smaller than 1400 pixels. We train our model on 8 GPUs with 16 images per mini-batch. SGD with the momentum of 0.9 is adopted as our optimizer. The initial learning rate is 0.01 for both RetinaNet and FCOS and is decayed by a factor of 10 after 8th epoch and 11th epoch. The training process finishes at the end of the 12th epoch. If not specified, the default backbone in our experiments is ResNet-50 [7]. We adopt IoU Loss for f^{reg} in Eq. 1. For back-propagation, the regression loss is replaced by GIoU Loss because we find such a usage of IoU/GIoU Loss can yield the best performances. For evaluation, we follow the standard Caltech [4]

Table 1. Results of LLA on two detectors – RetinaNet and FCOS. #A=1 denotes that only one anchor on each spatial location is used. Lower is better for MR.

Backbone	Method	MR	AP	Recall
ResNet-50	RetinaNet [14]	59.13	81.04	88.25
	RetinaNet*	54.00	80.54	87.33
	RetinaNet*(#A=1)	60.92	80.21	86.70
	LLA.RetinaNet*(#A=1)	49.60	84.69	89.90
ResNet-50	FCOS [24]	54.95	86.36	94.05
	LLA.FCOS	49.48	87.25	93.43
ResNet-101	RetinaNet	57.72	81.20	87.77
	LLA.RetinaNet*(#A=1)	47.29	86.25	91.27

evaluation metric – MR, which stands for the Log-Average Missing Rate over false positives per image (FPPI) ranging in $[10^{-2}, 100]$.

4.2. Ablation Studies

Effect of LLA. We present experimental results on RetinaNet and FCOS. λ is set to 1 and 1.3 while the NMS threshold is set to 0.5 and 0.6 for RetinaNet and FCOS, respectively. As seen in Table 1, RetinaNet achieves 59.13% MR and 81.04% AP on CrowdHuman. We further implement a modified version of RetinaNet, termed as RetinaNet*, in which we replace SmoothL1 Loss with recently proposed GIoU Loss to better regress instances’ boundaries. To better reflect the power of LLA, we reduce the number of anchors in RetinaNet* from 9 to 1. Our modified RetinaNet* achieves worse MR but better AP and Recall. Finally, we introduce LLA into RetinaNet*. RetinaNet* with LLA improves MR and AP by a large margin – 9.53% on MR and 3.65% on AP, respectively. Similar performance gain can be observed on a better backbone – ResNet-101. For FCOS, as suggested in ATSS [32], we adopt a better label assigning strategy – center sampling as our baseline. For LLA.FCOS, we remove the *Centerness* branch and replace center sampling with LLA. In this case, MR and AP still get boosted by 5.47% and 0.89%.

Visualizations of positive anchors assigned by LLA’s are shown in Fig. 2. As seen in the first column, instead of evenly distributed in each GT box, positive anchors for these dancers fall onto the foreground more compactly. In the second and third columns, positive anchors of heavily occluded human instances fall onto the visible regions (e.g. heads, shoulders, etc.) which are far away from the geometric center of GT boxes, demonstrating the effectiveness of LLA.

Analysis of Each Component in Restricted Cost Matrix. In Restricted Cost Matrix C_r , each term has its unique value. We start from the minimal requirement C^{cls} because C^{cls} helps identify the region of foreground instances. However, as seen in Table 2, without other two terms in C_r , the model fails to converge. This is mainly

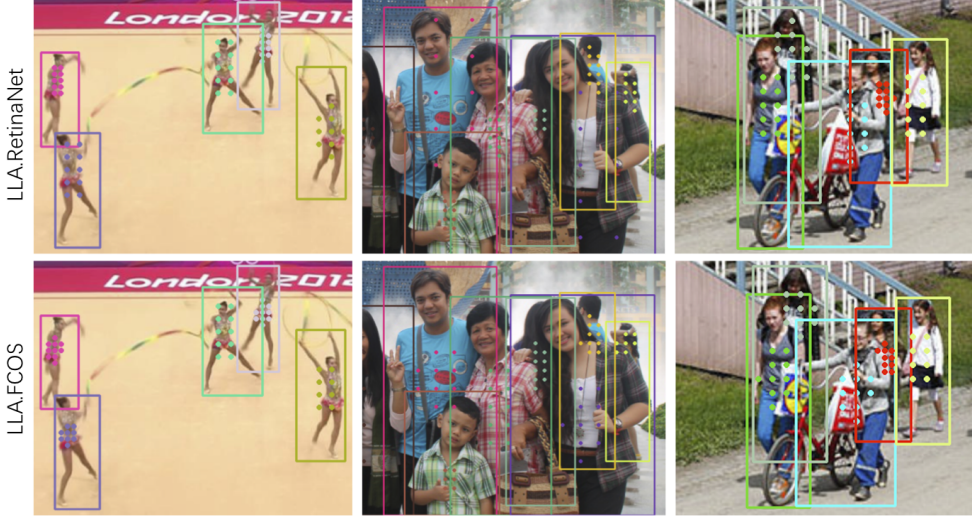


Figure 2. Label assigning results of LLA. Only FPN layers with the largest number of positive anchors are shown for better visualization.

Table 2. Contributions of each component in C_r on RetinaNet*. “-” denotes fail to converge.

C^{cls}	C^{reg}	C^{inbox}	MR	AP	Recall
✓			-	-	-
✓	✓		50.55	84.15	89.03
✓	✓	✓	49.60	84.69	89.90

Table 3. Performance of different usages of IoU/GIoU Loss. Values in bold stand for best results for MR and AP. “-” denotes failing to converge.

Before Assignment	After Assignment	MR	AP
IoU	IoU	49.78	84.48
IoU	GIoU	49.60	84.69
GIoU	IoU	-	-
GIoU	GIoU	52.52	80.54

due to that C^{cls} can not help model distinguish different instances in same category by incorporating spatial information. Without term C^{reg} in C_r , an anchor a in instance A can also be in topk list of instance B if A and B are in the same category. Such a mis-assignment would lead the optimization process to wrong directions. After introducing C^{reg} into C_r , the detector achieves 50.55% MR which already surpasses baseline by a large margin. Based on that, C^{inbox} can further improve MR by 0.95% to help stabilize training process at early training stage.

Different Usage of IoU/GIoU Loss. Intersection over Union is used twice in our work: 1). Calculating regression loss between each anchor-gt pair before assignment, defined as the C^{reg} term in Cost Matrix C. 2). Calculating regression loss for each assigned anchor-gt pair for back-propagation. As stated in Sec. 4.1, we use IoU Loss be-

fore assignment and GIoU Loss after assignment, because we found such a usage can yield the best detection performance. Here, we present the detection performances of other different settings in Table 3. Noted that our proposed LLA focuses on label assignment in object detection. The exploration of different IoU variants (e.g. DIoU [36], CIoU [36]) is beyond our scope. Hence, we did not conduct further experiments.

Effect of K . Hyper-parameter K can be viewed as the number of positive anchors we want for each GT. Intuitively, too large K will introduce many low-quality candidates while too small K will lead to an insufficient number of candidates and then hurt the detector’s accuracy. Thus we conduct heavy experiments to study the best value of K . We vary K from 1 to 16. As shown in Table 4, $K = 10$ achieves the best MR. We also observe that K is quite insensitive within a broad range which is a desired property for generalization on different datasets.

4.3. Further Analysis.

Ambiguous Assignment Ratio. Given a label assigning strategy, an anchor can be possibly assigned to several GTs, leading to the ambiguous assignment. For example, in RetinaNet, an anchor box may have an IoU value greater than 0.5 with multiple GTs. In this case, this anchor will be assigned to the GT with maximum IoU value. We call such kind of anchors which need further post-processing on assignment results as ambiguous anchors and define Ambiguous Assignment Ratio (AAR) as:

$$AAR = \frac{\#Ambiguous\ Anchors}{\#Positive\ Anchors} \quad (6)$$

A lower AAR means the adopted label assigning strategy assigns positive anchors in a more deterministic way, which

Table 4. Performance of LLA.RetinaNet under different values of K . Values in bold are the best results for MR and AP.

Top K	1	2	3	4	5	6	7	8
MR	53.78	52.19	51.73	51.25	50.38	51.38	50.46	50.06
AP	82.72	83.68	84.14	84.66	85.64	84.52	84.42	84.68
Top K	9	10	11	12	13	14	15	16
MR	49.87	49.60	49.86	50.03	50.30	50.16	50.38	50.54
AP	84.88	84.69	84.30	84.32	84.30	84.26	84.00	83.81

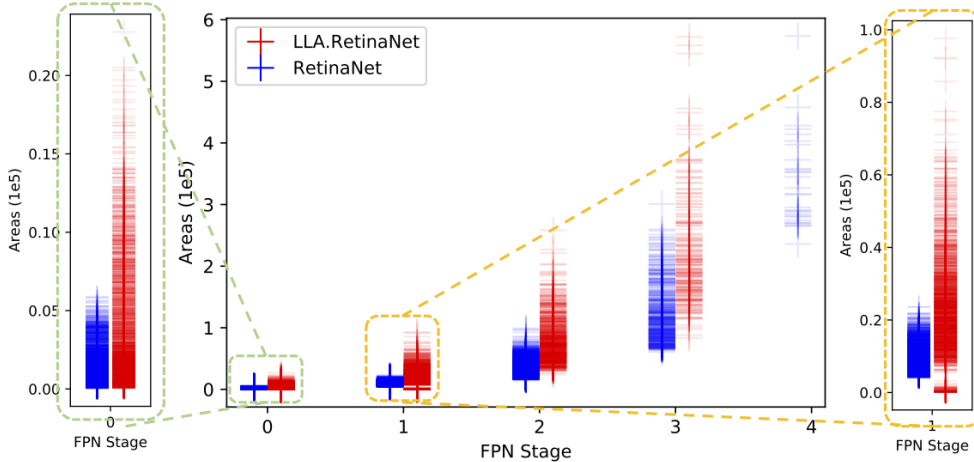


Figure 3. Visualization of objects’ areas and the FPN layers they are assigned to. Each light cross represents a GT box. We consider the layer which is assigned to the highest number of positive samples as the layer that GT is assigned to. “stage 0” denotes the highest resolution layer (*i.e.* “P3” in FPN [13]).

Table 5. AAR evaluated on RetinaNet and FCOS when w/ and w/o LLA. As can be seen, LLA can effectively reduce AAR.

w/ LLA	AAR(%)	
	RetinaNet	FCOS
No	7.4	13.2
Yes	6.2	4.6

is a desired property in the crowd scenario. We calculate the AAR for RetinaNet and FCOS when w/ and w/o LLA on CrowdHuman. Results shown in Table 5 illustrate that our proposed LLA can effectively decrease the AAR and reduce the ambiguity introduced in label assignment.

FPN Level Allocation. Which FPN level should each GT assigned to is crucial in label assignment. RPN and FCOS assign labels based on explicit geometric constraint (*i.e.* *area ranges* and *regression ranges*), while RetinaNet imposes implicit scale constraint by estimating IoUs between a set of pre-defined anchors and GTs. To compare the scale constraint learned by LLA and RetinaNet, we sample 5,000 GT annotations from CrowdHuman training set, for each GT, we plot its area and corresponding FPN layer which has the largest number of positive anchors of it. The results in Fig. 3 show that compared to RetinaNet, LLA tends to assign objects to a more fine-grained feature layer with higher resolution, in addition, no GT box is assigned to “P7”. Such

a phenomenon is reasonable because in the crowd scenario, dense anchors are more desired to precisely assign positive anchors for those heavily occluded human instances. Thus LLA can also be termed as **occlusion-aware** which is an appealing property in many real-world applications.

Evolution of Positive Anchors. We visualize the evolution of positive anchors during a training process in Fig. 4. We use LLA.RetinaNet and only visualize the FPN layer with the largest number of positive anchors. At early training stage, due to the under-fitting of the detector, some of the assigned positive anchors may fall onto the background or on other GT’s foreground. As the training process continues, positive anchors tend to migrate onto the foreground region of their corresponding GTs, demonstrating the effectiveness of LLA. Noted that although detectors like SSD [17], RetinaNet [14] and FCOS [24] tend to assign anchors close to the objects’ geometric center as their positive anchors, the success of PAA, AutoAssign [40] and our proposed LLA reveals that geometric center is not the best prior. However, we do not argue that semantic center is always a better prior location for positive anchors than geometric center. Instead, we can see in Fig. 4 that many anchors near objects’ geometric centers are also defined positive. Utilizing geometric center or semantic center or both of them is totally adaptive and only based on model’s predictions itself. That’s why LLA can achieve SOTA performance.

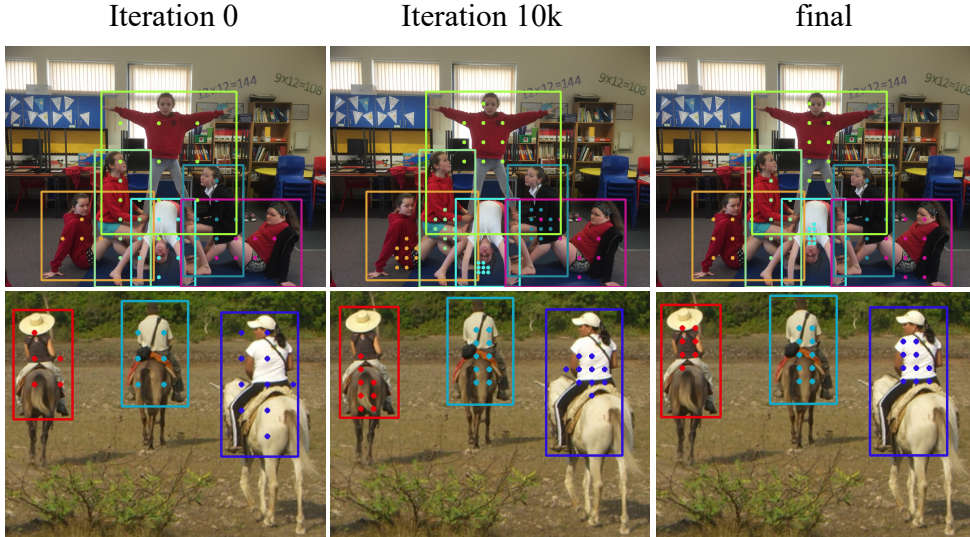


Figure 4. Evolution of positive anchors during a training process.

Table 6. Performance comparison with state-of-the-art label assigning strategies on CrowdHuman. PISA and Noisy Anchor are re-implemented based on RetinaNet. For PAA, we use the official code released by its author.

Method	Remarks	MR	AP	Recall
<i>Anchor-Based</i>				
RetinaNet [14]	-	59.13	81.04	88.25
Noisy Anchor [12]	Based on RetinaNet	53.03	84.01	89.12
PISA [2]	Based on RetinaNet	52.17	84.43	89.68
FreeAnchor [35]	-	50.95	82.61	86.90
PAA [9]	w/ IoU Branch [9],GN [28]	50.77	84.24	89.60
LLA.RetinaNet(Ours)	-	49.60	84.69	89.90
<i>Anchor-Free</i>				
FCOS [24]	w/o Centerness	71.47	83.34	93.76
	w/ Centerness	54.95	86.36	94.05
ATSS [32]	w/o Centerness	56.25	86.47	92.88
	w/ Centerness	49.51	87.41	94.19
LLA.FCOS(Ours)	w/o IoU Branch	49.48	87.25	93.43
	w/ IoU Branch	47.90	88.04	93.95

4.4. Comparison with State-of-the-art.

Note that LLA can be compatible with recent advances in DPD, thus in this section, we only compare LLA with other state-of-the-art label assigning strategies. NMS thresholds 0.5 and 0.6 are adopted for anchor-based and anchor-free methods, respectively. As seen in Table 6, For the anchor-based method, our LLA built upon RetinaNet surpasses all other methods without any bells and whistles (*e.g.* Group-Norm [28] and IoU Branch). For anchor-free methods, LLA is built upon FCOS. It needs to be mentioned that both FCOS and ATSS restrict positive anchors in the central region of an object, and thus they can benefit from the Centerness branch to eliminate false positives. Especially under MR – a metric which is extremely sensitive to false

positives, using Centerness can remarkably reduce MR for FCOS and ATSS. However, LLA does not acknowledge and use *center prior* in crowd scenarios, under this circumstance, adopting the Centerness branch will instead hurt the detector’s performance. For fair comparison, we adopt the IoU branch proposed in PAA as a replacement of Centerness. As shown in Table 6, LLA without IoU branch leads FCOS and ATSS *w/o* Centerness by more than 5% MR, also shows better results than ATSS *w/* Centerness. After further utilizing IoU branch, LLA surpasses ATSS by a clear margin.

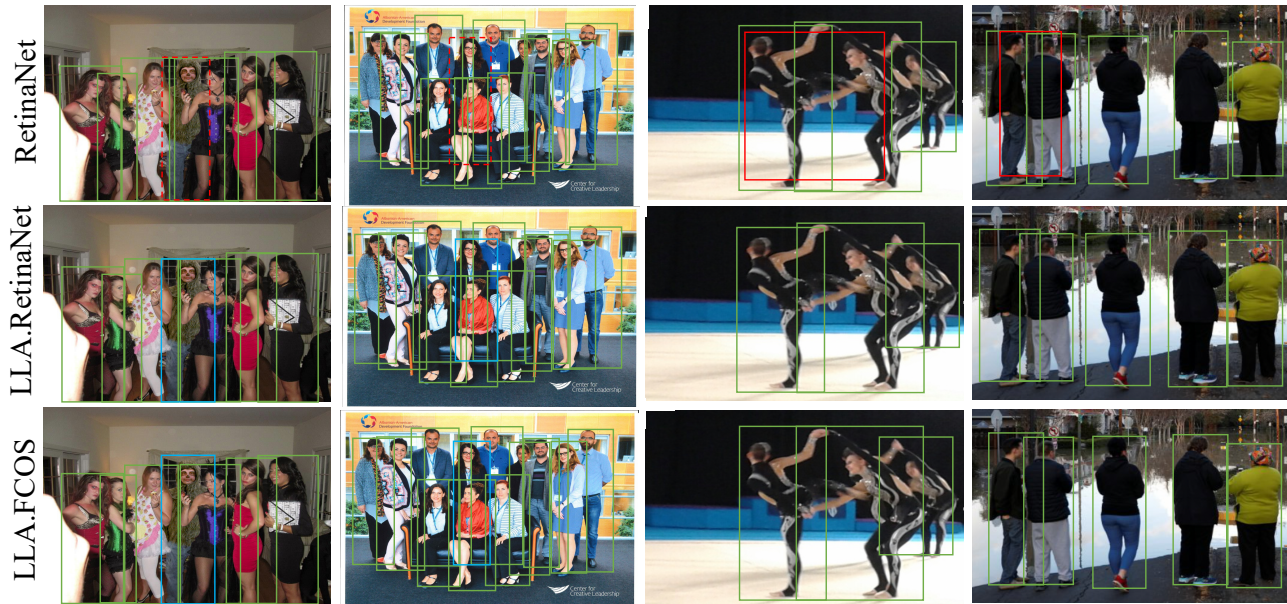


Figure 5. Prediction results of RetinaNet before and after adopting LLA, as well as LLA.FCOS. Red dotted and solid lines denote miss-detections and false positives in original RetinaNet, respectively. Blue solid lines represent that the miss-detections in RetinaNet are successfully detected in LLA.RetinaNet and LLA.FCOS.

4.5. Visualizing Prediction Results on CrowdHuman.

We compare the prediction results with and without LLA in Fig. 5. The first two columns exhibit that LLA can successfully detect the mis-detected instances by original RetinaNet. The last two columns show that LLA can effectively reduce false positives. We believe such two merits in LLA mainly benefit from the better placement of positive anchors, which greatly reduces the ambiguity during the training stage.

Table 7. Performance comparison of different label assigning strategies on CityPersons [31].

Method	MR	
	Reasonable	Heavy
<i>Anchor-Based</i>		
RetinaNet	16.83	46.92
FreeAnchor	15.01	46.72
LLA.RetinaNet(Ours)	14.34	44.70
<i>Anchor-Free</i>		
FCOS	15.27	46.82
ATSS	13.74	43.86
LLA.FCOS(Ours)	12.08	43.72

4.6. Experiments on CityPersons.

CityPersons [31] is another dataset for pedestrian detection which consists of 2975 images for training, 500 and

1575 images for validation and testing. Following the evaluation protocol in CityPersons, objects whose heights are less 50 are ignored. Besides, the validation set is further divided into two subsets according to visibility – Reasonable and Heavy Occlusion. MR on these two subsets is reported in this work. We follow the same training details as in CrowdHuman. IoU branch is adopted in Anchor Free LLA. As shown in Table 7, our anchor-based and anchor-free LLA reduce MR on the Heavy Occlusion subset by 2.22% and 3.10% respectively, exceeding all other existing label assigning strategies.

5. Conclusion.

In this work, we propose **Loss-aware Label Assignment** (LLA), an extremely simple but effective label assigning strategy for pedestrian detection in crowd scenarios. It defines positive/negative anchors based on the values of joint losses and defines anchors with smaller loss as positives. LLA does not utilize any human prior such as spatial (center/IoU constraint in ATSS/RetinaNet) and scale prior (scale constraint in FPN), making LLA fully adaptive. Experimental results on CrowdHuman and CityPersons demonstrate LLA’s superiority over other label assigning strategies as well as its generalization ability on various pedestrian detection datasets.

References

- [1] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the*

- IEEE Conference on Computer Vision and Pattern Recognition*, pages 6154–6162, 2018.
- [2] Yuhang Cao, Kai Chen, Chen Change Loy, and Dahua Lin. Prime sample attention in object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11583–11591, 2020.
 - [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. 2020.
 - [4] Piotr Dollár, Christian Wojek, Bernt Schiele, and Pietro Perona. Pedestrian detection: An evaluation of the state of the art. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(4):743–761, 2011.
 - [5] Zheng Ge, Zequn Jie, Xin Huang, Rong Xu, and Osamu Yoshie. Ps-rnn: Detecting secondary human instances in a crowd via primary object suppression. In *2020 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2020.
 - [6] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.
 - [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
 - [8] Xin Huang, Zheng Ge, Zequn Jie, and Osamu Yoshie. Nms by representative region: Towards crowded pedestrian detection by proposal pairing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10750–10759, 2020.
 - [9] Kang Kim and Hee Seok Lee. Probabilistic anchor assignment with iou prediction for object detection. 2020.
 - [10] Tao Kong, Fuchun Sun, Huaping Liu, Yuning Jiang, Lei Li, and Jianbo Shi. Foveabox: Beyond anchor-based object detection. *IEEE Transactions on Image Processing*, 29:7389–7398, 2020.
 - [11] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European Conference on Computer Vision*, pages 734–750, 2018.
 - [12] Hengduo Li, Zuxuan Wu, Chen Zhu, Caiming Xiong, Richard Socher, and Larry S Davis. Learning from noisy anchors for one-stage object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10588–10597, 2020.
 - [13] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2117–2125, 2017.
 - [14] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2980–2988, 2017.
 - [15] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of the European Conference on Computer Vision*, pages 740–755. Springer, 2014.
 - [16] Songtao Liu, Di Huang, and Yunhong Wang. Adaptive nms: Refining pedestrian detection in a crowd. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6459–6468, 2019.
 - [17] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *Proceedings of the European Conference on Computer Vision*, pages 21–37. Springer, 2016.
 - [18] Yanwei Pang, Jin Xie, Muhammad Haris Khan, Rao Muhammad Anwer, Fahad Shahbaz Khan, and Ling Shao. Mask-guided attention network for occluded pedestrian detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4967–4975, 2019.
 - [19] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 779–788, 2016.
 - [20] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems*, pages 91–99, 2015.
 - [21] Douglas A Reynolds. Gaussian mixture models. *Encyclopedia of biometrics*, 741, 2009.
 - [22] Hamid Rezaatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 658–666, 2019.
 - [23] Shuai Shao, Zijian Zhao, Boxun Li, Tete Xiao, Gang Yu, Xiangyu Zhang, and Jian Sun. Crowdhuman: A benchmark for detecting human in a crowd. *arXiv preprint arXiv:1805.00123*, 2018.
 - [24] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9627–9636, 2019.
 - [25] Jiaqi Wang, Kai Chen, Shuo Yang, Chen Change Loy, and Dahua Lin. Region proposal by guided anchoring. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2965–2974, 2019.
 - [26] Jianfeng Wang, Lin Song, Zeming Li, Hongbin Sun, Jian Sun, and Nanning Zheng. End-to-end object detection with fully convolutional network. *arXiv preprint arXiv:2012.03544*, 2020.
 - [27] Xinlong Wang, Tete Xiao, Yuning Jiang, Shuai Shao, Jian Sun, and Chunhua Shen. Repulsion loss: Detecting pedestrians in a crowd. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7774–7783, 2018.
 - [28] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision*, pages 3–19, 2018.
 - [29] Tong Yang, Xiangyu Zhang, Zeming Li, Wenqiang Zhang, and Jian Sun. Metaanchor: Learning to detect objects with customized anchors. In *Advances in Neural Information Processing Systems*, pages 320–330, 2018.

- [30] Ze Yang, Shaohui Liu, Han Hu, Liwei Wang, and Stephen Lin. Reppoints: Point set representation for object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9657–9666, 2019.
- [31] Shanshan Zhang, Rodrigo Benenson, and Bernt Schiele. Citypersons: A diverse dataset for pedestrian detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3213–3221, 2017.
- [32] Shifeng Zhang, Cheng Chi, Yongqiang Yao, Zhen Lei, and Stan Z Li. Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9759–9768, 2020.
- [33] Shifeng Zhang, Longyin Wen, Xiao Bian, Zhen Lei, and Stan Z Li. Occlusion-aware r-cnn: detecting pedestrians in a crowd. In *Proceedings of the European Conference on Computer Vision*, pages 637–653, 2018.
- [34] Shanshan Zhang, Jian Yang, and Bernt Schiele. Occluded pedestrian detection through guided attention in cnns. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 6995–7003, 2018.
- [35] Xiaosong Zhang, Fang Wan, Chang Liu, Rongrong Ji, and Qixiang Ye. Freeanchor: Learning to match anchors for visual object detection. In *Advances in Neural Information Processing Systems*, pages 147–155, 2019.
- [36] 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 *AAAI*, pages 12993–13000, 2020.
- [37] Chunluan Zhou and Junsong Yuan. Bi-box regression for pedestrian detection and occlusion estimation. In *Proceedings of the European Conference on Computer Vision*, pages 135–151, 2018.
- [38] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019.
- [39] Xingyi Zhou, Jiacheng Zhuo, and Philipp Krahenbuhl. Bottom-up object detection by grouping extreme and center points. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 850–859, 2019.
- [40] Benjin Zhu, Jianfeng Wang, Zhengkai Jiang, Fuhang Zong, Songtao Liu, Zeming Li, and Jian Sun. Autoassign: Differentiable label assignment for dense object detection. *arXiv preprint arXiv:2007.03496*, 2020.
- [41] Chenchen Zhu, Yihui He, and Marios Savvides. Feature selective anchor-free module for single-shot object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 840–849, 2019.