# OneNet: Towards End-to-End One-Stage Object Detection

Peize Sun[1], Yi Jiang[2], Enze Xie[1], Zehuan Yuan[2], Changhu Wang[2], Ping Luo[1]

[1]The University of Hong Kong    [2]ByteDance AI Lab

## Abstract

*End-to-end one-stage object detection trailed thus far. This paper discovers that the lack of classification cost between sample and ground-truth in label assignment is the main obstacle for one-stage detectors to remove Non-maximum Suppression(NMS) and reach end-to-end. Existing one-stage object detectors assign labels by only location cost, e.g. box IoU or point distance. Without classification cost, sole location cost leads to redundant boxes of high confidence scores in inference, making NMS necessary post-processing.*

*To design an end-to-end one-stage object detector, we propose Minimum Cost Assignment. The cost is the summation of classification cost and location cost between sample and ground-truth. For each object ground-truth, only one sample of minimum cost is assigned as the positive sample; others are all negative samples. To evaluate the effectiveness of our method, we design an extremely simple one-stage detector named OneNet. Our results show that when trained with Minimum Cost Assignment, OneNet avoids producing duplicated boxes and achieves to end-to-end detector. On COCO dataset, OneNet achieves 35.0 AP/80 FPS and 37.7 AP/50 FPS with image size of 512 pixels. We hope OneNet could serve as an effective baseline for end-to-end one-stage object detection. The code is available at: https://github.com/PeizeSun/OneNet.*

## 1. Introduction

Object detection is one of the fundamental tasks in the computer vision area and enables numerous downstream applications. It aims at localizing a set of objects and recognizing their categories in an image. One of the challenging topic for current object detectors is label assignment [11, 10, 28, 1, 26, 21, 19, 33, 41, 40, 2]. Specifically, how to define the positive samples for each object and negative samples for the background has always been an open problem.

For decades, the positive sample in object detection is box candidates whose intersection-over-union(IoU) with ground-truth boxes is larger than the threshold. In classical computer vision, the classifier is applied on sliding-
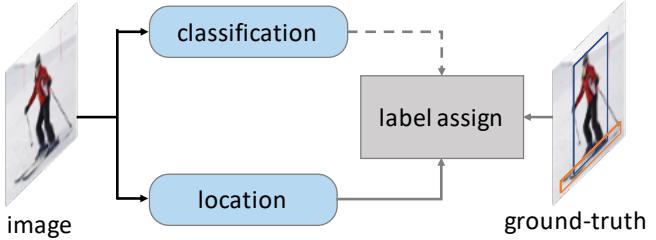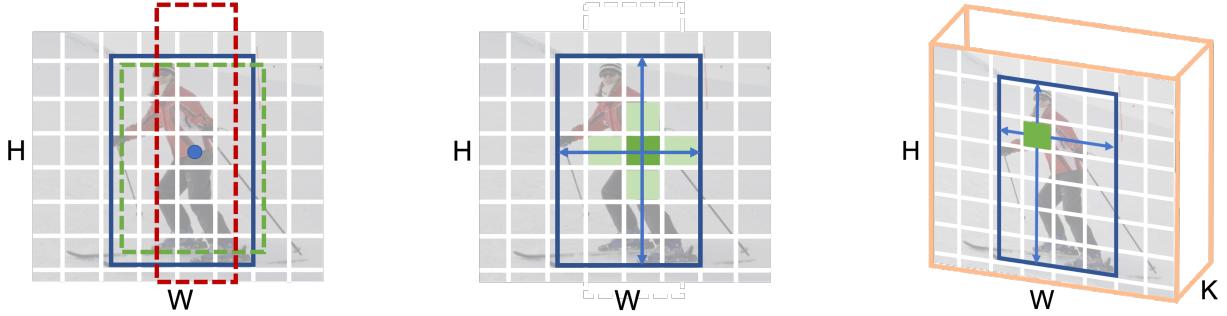


Figure 1: **The lack of classification information in label assignment of previous one-stage detectors**. Previous works assign labels by only location cost between sample and ground-truth, such as box IoU or point distance. The lack of classification cost leads to redundant boxes of high confidence scores, making NMS necessary post-processing.

windows enumerated on the image grid [4, 9, 36]. Modern detectors pre-define thousands of anchor boxes on the image grid and performs classification and regression on these candidates [11, 28, 19, 27]. We call these box-based label assign methods as "Box Assignment", shown in Figure 2a.

Despite box candidate methods dominate object detection for years, the detection performance is largely sensitive to sizes, aspect ratios, and the number of anchor boxes. To eliminate the hand-crafted design and complex computation of box candidates, anchor-free detectors [33, 41] are rising. These methods directly treat grid points in the feature map as object candidates and predict the offset from the grid point to the object box's boundaries. The label assignment in anchor-free detectors is simplified to point distance from the grid point to the object box's center. We call the point-based label assign paradigm as "Point Assignment", shown in Figure 2b.

Both Box Assignment and Point Assignment suffer from one common dilemma of many-to-one assignment [1, 39, 40]. For one ground-truth box, there is more than one positive sample. The detection performance is sensitive to hyper-parameters in the assignment process. Worsely, redundant and near-duplicate results are produced, thus making non-maximum suppression(NMS) necessary post-processing.

**(a) Box Assignment**: YOLO,RetinaNet. Matching cost is only location cost of box IoU. Boxes are labeled as positive(green) if IoU with ground-truth box is greater than high threshold, negative(red) if smaller than low threshold.

**(b) Point Assignment**: FCOS,CenterNet. Matching cost is only location cost of point distance. The point closest to center of ground-truth box is assigned to positive label. Nearby points have a reduced negative loss.

**(c) Minimum Cost Assignment**: OneNet. Matching cost is summation of classification cost and location cost. Positive label is assigned to only one sample of minimum cost, others are all negative ones.

Figure 2: **Comparisons of different label assignment methods.** $H$ and $W$ are height and width of feature map, respectively, $K$ is the number of object categories. Previous works on one-stage object detection assign labels by only location cost, such as (a) box IoU or (b) point distance between sample and ground-truth. In our method, however, (c) classification cost is additionally introduced. We discover that **classification cost is the key to the success of end-to-end**. Without classification cost, only location cost leads to redundant boxes of high confidence scores in inference, making NMS a necessary post-processing component.

Recently, one-to-one assignment achieves great success in sparse candidates and multiple-stage refinement detectors [2, 42, 31], where one ground-truth bounding box is assigned to only one positive sample, other samples are all negative ones. They directly output detection results without NMS. The detection accuracy performance of these detectors is promising. However, dense candidates and one-stage refinement detectors have the potential to be faster and simpler. *It inspires us to design a straightforward one-to-one label assignment strategy for end-to-end one-stage detectors.*

An intuitive idea is that we replace many-to-one label assignment in previous one-stage detectors to one-to-one, *e.g.*, only one positive sample is assigned to one ground-truth bounding box, others samples on the whole feature map are all negative ones. However, this leads to very poor performance, shown in Table 1. During inference, the detector still outputs multiple redundant boxes of high confidence scores for one object, and thus NMS is still needed.

We discover that the lack of classification cost between sample and ground-truth in label assignment is the main obstacle for a one-stage detector to remove NMS and reach end-to-end. By reviewing existing one-stage object detectors, we identity that they assign labels by only location cost [26, 19, 33, 41], such as box IoU (Figure 2a) or point distance (Figure 2b) Whereas, the well-established end-to-end detectors[2, 42, 31] assign labels by both location cost and classification cost. Our experimental and visual results

demonstrate that no classification cost leads to redundant boxes of high confidence scores in inference, making NMS a necessary post-processing component.

To design an end-to-end one-stage object detector, we propose a simple but effective label assignment strategy, named *Minimum Cost Assignment*, where the cost is the summation of classification cost and location cost. In detail, for each ground-truth box, only one sample of minimum cost among all samples is assigned as a positive sample; others are all assigned as negative samples, shown in Figure 2c. To demonstrate the proposed label assignment strategy's effectiveness, we design a detector called *OneNet*, named for its one-stage pipeline and one positive sample. Our experiments show that when trained with Minimum Cost Assignment, OneNet avoids producing duplicated boxes and achieves to an end-to-end detector. OneNet exhibits several appealing advantages:

- The whole network is fully-convolutional and end-to-end training. No RoI operation or Attention interaction is used in the pipeline.

- Label assignment is based on the minimum cost of classification and location, rather than hand-crafted heuristic rule or complex bipartite-matching.

- No any post-procedure, such as non-maximum suppression(NMS) or max-pooling, is involved during inference, making highly-efficient object detection for the industry.

On COCO dataset [20], OneNet achieves 35.0 AP/80 FPS using ResNet-50 [13], 37.7 AP/50 FPS using ResNet-101, with image size of 512 pixels. We claim that a bigger image size and a more complex backbone can further improve detection accuracy, however, inference speed is accordingly delayed, shown in Table 3 and Table 5. Our method is aimed to serve as an effective baseline for end-to-end one-stage object detection. Therefore, the reported performance is a trade-off between accuracy and speed.

Assigning labels by location cost is conceptually intuitive and popularize in object detection to date. However, we surprisingly discover that this widely-used method is the obstacle of the end-to-end detector. We explain that the main reason is the *misalignment between label assignment and network optimization objective*. Object detection is a multi-task of classification and location. The selected positive sample by only location cost can maximize contribution to location task but cannot ensure the optimal classification. The optimal solution is that one object has one prediction. However, if trained with only location cost, the classification branch is forced to output the approximate solution; that is, one object has multiple predictions. We hope our work could inspire re-thinking the label assignment in object detection and exploring the next generation of object detectors.

## 2. Related Work

**Object detection.** Object detection is one of the most fundamental and challenging topics in computer vision fields. Limited by classical feature extraction techniques [4, 36], the performance has plateaued for decades, and the application scenarios are limited. With the rapid development of deep learning [16, 30, 32, 13, 14], object detection achieves powerful performance [8, 20]. Modern object detection mainly consists of one-stage and two-stage detector.

**One-Stage detector.** One-stage detector directly predicts the category and location of dense anchor boxes or points over different spatial positions and scales in a single-shot manner such as YOLO [26], SSD [21] and RetinaNet [19]. YOLO [26] divides the image into an S × S grid, and if the center of an object falls into a grid cell, the corresponding cell is responsible for detecting this object. SSD [21] spreads out anchor boxes on multi-scale feature map layers within a ConvNet to directly predict object category and anchor box offsets. RetinaNet [19] utilizes focal loss to ease the extremely unbalance of positive and negative samples based on the FPN [18]. Recently, anchor-free algorithms [15, 17, 33, 41, 6, 40] is proposed to make this pipeline much simpler by replacing hand-crafted anchor boxes with reference points. All of the above methods are built on dense points, and each point is directly classified and regressed. These points are assigned to ground-truth

object boxes in training time based on a pre-defined principle, *e.g.*, whether the anchor has a higher intersection-over-union (IoU) threshold with its corresponding ground truth, or whether the reference point falls in one of the object boxes. CornerNet [17] generates tons of keypoints by heatmap and group them by the Associative Embedding [24], CornerNet has impressive performance but will have too many false positives. CenterNet [41] directly uses the center point to regress the target object on a single scale. FCOS [33] assign the objects of different size and scales to multi-scale feature maps, with the power of FPN [18], FCOS can achieve state-of-the-art performance without bells and whistles. ATSS[40] reveals that the essential difference between anchor-based and anchor-free detection is how to define positive and negative training samples, leading to the performance gap between them. Despite obtaining remarkable performance, NMS post-processing is needed to remove redundant predictions during inference time for the one-stage dense detectors.

**Two-Stage detector.** The two-stage detector is another excellent pipeline and has dominated modern object detection for years [1, 3, 10, 12, 28]. This paradigm can be viewed as an extension of the one-stage detector. It firstly generates a high-quality set of foreground proposals from dense region anchors by region proposal networks and then refines each proposal's location and predicts its specific category. The region proposal algorithm plays an essential role in the first stage of these two-stage methods. Fast R-CNN [10] use Selective Search [34] to generate foreground proposals and refine the proposals in R-CNN [11] Head. Faster R-CNN [28] presents the region proposal network, which generates high-quality proposals by CNN-based Networks. Cascade RPN [37] improves the region proposal quality and detection performance by systematically addressing the limitation of the conventional RPN that heuristically defines the anchors and aligns the features to the anchors. Recently CPNDet [7] proposes a new corner proposal network instead of region proposal networks. CPNDet uses the [17] to generate better proposals, and it can be viewed as a two-stage anchor free object detection pipeline. In the refinement stage, the pre-defined sampling method for foreground and background is essential. Cascade R-CNN [1] iteratively uses multiple R-CNN heads with different label assign threshold to get high-quality detection boxes. Libra R-CNN [25] try to solve the unbalance problems in sample level, feature level, and objective level. Grid R-CNN [23] adopts a grid guided localization mechanism for accurate object detection instead of traditional bounding box regression. The two-stage methods show great detection accuracy in object detection. However, too many hyper-parameters and modules are introduced to design the region proposal stage and object recognition stage; thus, two-stage detectors are not easily used in the industry.

**End-to-end object detection.** The well-established end-to-end object detectors are based on sparse candidates and multiple-stage refinement. DETR [2] is proposed to directly output the predictions without any hand-crafted label assignment and post-processing procedure, achieving fantastic performance. DETR can be viewed as the first end-to-end object detection method; DETR utilizes a sparse set of object queries to interact with the global image feature. Benefit from the global attention mechanism [35] and the bipartite matching between predictions and ground truth objects, DETR can discard the NMS procedure while achieving remarkable performance. Deformable-DETR [42] is introduced to restrict each object query to a small set of crucial sampling points around the reference points, instead of all points in the feature map. Sparse R-CNN [31] starts from a fixed sparse set of learned object proposals and iteratively performs classification and localization to the object recognition head.

## 3. Label Assignment

The Minimum Cost Assignment is designed to address label assignment of end-to-end one-stage object detection. We introduce starting from previous label assignment methods of one-stage detectors, such as Box Assignment and Point Assignment.

### 3.1. Matching Cost

Previous methods assign samples by box IoU or point distance between samples and ground-truth. We summarize them as assigning labels by location cost. The location cost is defined as follows:

$$\mathcal{C}_{loc} = \lambda_{iou} \cdot \mathcal{C}_{iou} + \lambda_{L1} \cdot \mathcal{C}_{L1} \qquad (1)$$

where $\mathcal{C}_{L1}$ and $\mathcal{C}_{iou}$ are L1 loss and IoU loss between sample and ground truth box(center), respectively. $\lambda_{L1}$ and $\lambda_{iou}$ are coefficients. In Box Assignment, $\lambda_{L1} = 0$. In Point Assignment, $\lambda_{iou} = 0$.

Assigning labels by location cost are conceptually intuitive and offer excellent detection performance to date. However, object detection is a multi-task of location and classification. Only location cost causes sub-optimal classification performance. Specifically, it leads to redundant boxes of high confidence scores, making NMS necessary post-processing. More details are discussed in Section 5.

Towards end-to-end one-stage object detection, we introduce classification cost into matching cost. Cost is summation of classification cost and location cost between sample and ground-truth, defined as follows:

$$\mathcal{C} = \lambda_{cls} \cdot \mathcal{C}_{cls} + \mathcal{C}_{loc} \qquad (2)$$

where $\mathcal{C}_{cls}$ is classification loss of predicted classifications and ground truth category labels. $\mathcal{C}_{loc}$ is defined in Equa-

**Algorithm 1** Pseudocode of Minimum Cost Assignment in a PyTorch-like style.

```
# For simplicity,
# cross entropy loss as classification cost
# L1 loss as location cost

# Input:
# class_pred, box_pred: network output(HxWxK, HxWx4)
# gt_label, gt_box: ground-truth (N, Nx4)

# Output:
# src_ind: index of positive sample in HW sequence(N)
# tgt_ind: index of corresponding ground-truth (N)

# flattened class: HWxK
output_class = class_pred.view(-1, K)

# flattened box: HWx4
output_box = box_pred.view(-1, 4)

# classification cost: HWxN
cost_class = -torch.log(output_class[:, gt_label])

# location cost: HWxN
cost_loc = torch.cdist(out_box, gt_bbox, p=1)

# cost matrix: HWxN
cost_mat = cost_class + cost_loc

# index of positive sample: N
_, src_ind = torch.min(cost_mat, dim=0)

# index of ground-truth: N
tgt_ind = torch.arange(N)
```

tion 1. $\lambda_{cls}$ is coefficient. In Box Assignment and Point Assignment, $\lambda_{cls} = 0$.

### 3.2. Minimum Cost Assignment

Minimum Cost Assignment is a straightforward method: for each ground-truth, only one sample of minimum cost among all samples is chosen as a positive sample; others are all negative ones. No hand-crafted heuristic rule or complex bipartite-matching[2, 42, 31] is involved. Algorithm 1 shows an illustrative example of Minimum Cost Assignment: cross-entropy loss as classification cost and L1 loss as location cost.

In dense detectors, classification loss is Focal Loss [19]. Following[2, 42, 31], location cost contains L1 loss and generalized IoU(GIoU) loss [29]. Finally, the cost is defined as follows:

$$\mathcal{C} = \lambda_{cls} \cdot \mathcal{C}_{cls} + \lambda_{L1} \cdot \mathcal{C}_{L1} + \lambda_{giou} \cdot \mathcal{C}_{giou} \qquad (3)$$

where $\mathcal{C}_{cls}$ is focal loss of predicted classifications and ground truth category labels, $\mathcal{C}_{L1}$ and $\mathcal{C}_{giou}$ are L1 loss and GIoU loss between normalized center coordinates and height and width of predicted boxes and ground truth box, respectively. $\lambda_{cls}$, $\lambda_{L1}$ and $\lambda_{giou}$ are coefficients of each component.

We note that [41] also chooses only one positive sample for one ground-truth, but nearby samples are labeled by Gaussian kernel. Our method is much more straightforward: Except for the chosen positive sample, all others are negative ones.
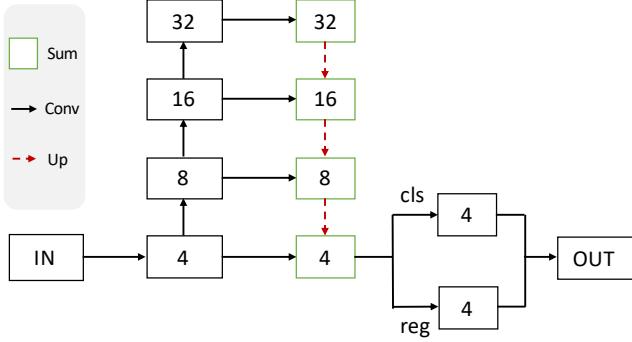
Figure 3: **Pipeline of OneNet.** For an input image of $H \times W \times 3$, the backbone generates feature map of $\frac{H}{4} \times \frac{W}{4} \times C$, the head produces classification prediction of $\frac{H}{4} \times \frac{W}{4} \times K$, where $K$ is number of categories, and regression prediction of $\frac{H}{4} \times \frac{W}{4} \times 4$, the final output is direct top-k(*e.g.*,100) scoring boxes.



Figure 4: **Multi-head Training(optional).** Cascading prediction heads are employed during training, where parameters of classification convolution and regression convolution are shared cross the heads, respectively. During inference, only first head is used, so the inference speed is not delayed.

## 4. OneNet

OneNet is an elementary fully-convolutional one-stage detector without any post-procedure such as NMS. The pipeline are shown in Figure 3.

**Backbone.** The backbone is a first bottom-up and then top-down structure. The bottom-up component is ResNet architecture [13] to produce multi-scale feature maps. The top-down architecture with lateral connections (FPN) [18] is developed to generate the final feature map for object recognition. The output feature is the shape of $H/4 \times W/4 \times C$, where $H$ and $W$ are the input image's height and width.

**Head.** The head performs classification and location on each grid point of the feature map $H/4 \times W/4$ by two parallel convolutional layers. The classification layer predicts the probability of object presence at each grid point for $K$ object categories. The location layer predicts the offset from each grid point to 4 boundaries of the ground-truth box.

**Training.** The label assignment is Minimum Cost Assignment. The training loss is similar to matching cost, composed of Focal Loss, L1 loss, and GIoU loss.

**Inference.** The final output is direct top-k(*e.g.*,100) scoring boxes. There is no any post-procedure, such as NMS or max-pooling operation [26, 19, 33, 41].

### 4.1. Multi-head Training

We propose an optional Multi-head Training strategy. It mainly consists of cascading prediction heads and weight-sharing mechanics.

**Cascading heads.** For the first stage, the input feature, denoted by $F_0$, is broadcast twice in channel dimension ($H/4 \times W/4 \times C \rightarrow H/4 \times W/4 \times 2C$). Then it is
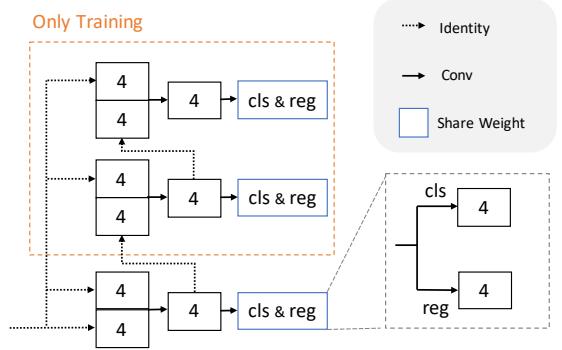
fed into a convolution layer to output feature $F_1$ of shape $H/4 \times W/4 \times C$. Based on $F_1$, classification prediction and regression prediction are produced. For the later stage $j$, the original feature $F_0$ of shape $H/4 \times W/4 \times C$ and the previous stage's feature $F_{j-1}$ of shape $H/4 \times W/4 \times C$ are concatenated in channel dimension to produce a composite feature of shape $H/4 \times W/4 \times 2C$. Then the feature map $F_j$ of shape $H/4 \times W/4 \times C$ is generated to perform classification and regression.

**Weight-sharing.** The classification convolution and regression convolution in each head share the same weight, respectively.

Simply applying Multi-head Training brings no improvement on detection accuracy and the inference speed is also delayed. We introduce two modifications to enable the final performance benefit from Multi-head Training: Large learning rate and Single-head Inference.

**Large learning rate.** Large learning rate has the potential to bring accuracy improvement. However, directly increasing the learning rate of single prediction head unstablize the training and degenerates the detection accuracy. When equipped with cascading heads and weight-sharing, the training learning rate could be increased and leads to improved accuracy.

**Single-head Inference.** During inference, only the first stage is used to output final results and other stages are discarded. This inference only leads to negligible accuracy drop compared with multiple-head inference. As a result, no additional computation cost is introduced in inference and accuracy performance is improved. The Multi-head Training is an optional training strategy for OneNet. Our experiments show that three stages are able to improve about 1 AP on COCO dataset, shown in Table 2.
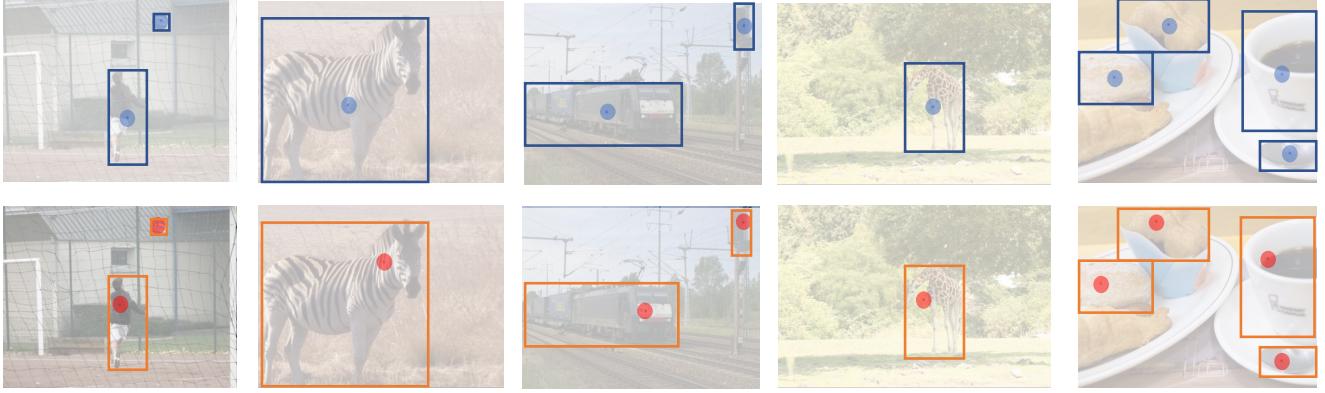
5

Figure 5: **Visualization of the positive sample.** 1st row is only location cost. 2nd row is summation of classification cost and location cost. The positive grid points are highlighted by surrounding circles for better visualization. The drawn boxes are ground-truth boxes. The positive samples assigned by only location cost are grid points closest to the ground-truth box center. Adding classification cost, positive samples are grid points in more discriminative areas, *e.g.* head of zebra.

## 5. Experiments

**Dataset.** Our experiments are conducted on the challenging MS COCO benchmark [20]. Using the standard metrics for object detection. All models are trained on the COCO `train2017` split (∼118k images) and evaluated with `val2017` (5k images).

**Implementation details.** The optimizer is AdamW [22] with weight decay 0.0001. The batch size is 64 images with 8 NVIDIA V100 GPUs. The initial learning rate is set to $5 \times 10^{-5}$ on single head training and $1 \times 10^{-4}$ on multiple head training. The default training schedule is 270k iterations, and the learning rate is divided by 10 at iteration 210k and 250k, respectively. Default backbone is ResNet-50 [13]. The backbone is initialized with the pre-trained weights on ImageNet [5]. Following [2, 42, 31], $\lambda_{cls} = 2$, $\lambda_{L1} = 5$, $\lambda_{giou} = 2$. Data augmentation includes random horizontal, scale jitter of resizing the input images such that the shortest side is at least 416 and at most 512 pixels. The default image size in inference is the shortest side of 512 pixels, and the longest size is at most 853 pixels. In inference, top-100 scoring boxes are selected as the final output.

### 5.1. Visualization of Positive Sample

We visualize the positive samples of CenterNet [41] and OneNet, as shown in Figure 5. The main difference between these two methods is that CenterNet follows location cost for label assignment, while OneNet follows the minimum cost assignment with classification cost and location cost.

For CenterNet, the positive sample lies in the grid point closest to the center of the object ground-truth box. This assignment is beneficial for box regression but is not a practical choice for positive and negative classification. In the first image, the person's twisted pose makes the positive

sample the point on the human body's edge. It is not the most discriminative region for object recognition.

From OneNet, we see that the positive sample lies in the more discriminative region of the object, *e.g.*, the inside of the human body, the head of the zebra. These choices are much more useful for classification. Meanwhile, it does not hurt the box regression since the positive samples are inside the object ground-truth box. We explain that the Minimum Cost Assignment's positive samples consider the location cost and the classification cost, so the final choice is a better choice overall.

### 5.2. Ablation Study

**Label assignment.** We carry out a series of experiments to study the effect of label assignment, shown in Table 1. The location cost of predicted box and classification loss is defined as Section 3. The pre-defined location cost is the distance between the fixed position of a grid point in the feature map and the ground-truth box center position. Without classification loss, pre-defined location cost is exactly the cost used in the label assignment method in CenterNet. From Table 1, the classification cost is the key to remove NMS. Without classification cost, NMS improves the performance by a remarkable margin. Whereas adding classification cost eliminates the necessity of NMS.

We also visualize prediction results of the model in Table 1 in Figure 6. We find that the overall confidence scores from different models vary greatly, so the visualization score threshold is adjusted in different models. From Figure 6, redundant boxes of high confidence scores are produced by the model without classification cost, which makes NMS post-processing a necessary component. Instead, adding classification costs reduces duplicated boxes and leads to the success of the end-to-end detector.

(a) **Cost is pre-defined location cost**. 1st row shows boxes of score higher than 0.4, 2nd shows 0.2.

(b) **Cost is summation of pre-defined location cost and classification cost**. Boxes of score higher than 0.4 are shown.

(c) **Cost is predicted location cost**. 1st row shows boxes of score higher than 0.4, 2nd shows 0.05.

(d) **Cost is summation of predicted location cost and classification cost**. Boxes of score higher than 0.4 are shown.

Figure 6: **Visualization of prediction results**. The visualization score threshold is adjusted in different models. Without classification cost, redundant boxes of (relative) high confidence scores are produced, making NMS necessary post-processing. Instead, adding classification cost reduces duplicated boxes and leads to the success of the end-to-end detector.

| Location cost |  | Classification cost | AP | AP$_{50}$ | AP$_{75}$ | AP$_s$ | AP$_m$ | AP$_l$ | FPS | AP(+NMS) |
| pre-defined | predicted |  |  |  |  |  |  |  |  |  |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| ✓ |  |  | 21.5 | 30.5 | 23.6 | 14.4 | 26.6 | 26.0 | 67 | 31.0 (+9.5) |
| ✓ |  | ✓ | 32.5 | 51.7 | 34.6 | 13.8 | 35.8 | 47.7 | 67 | 32.4 (-0.1) |
|  | ✓ |  | 2.1 | 3.9 | 2.0 | 3.8 | 2.0 | 2.6 | 67 | 8.7 (+6.6) |
|  | ✓ | ✓ | 35.7 | 54.3 | 38.4 | 17.9 | 39.3 | 48.6 | 67 | 35.5 (-0.2) |

Table 1: **Ablation on Label Assignment.** The pre-defined location cost is the distance between the fixed position of a grid point in the feature map and the ground-truth box center position. The location cost of predicted box and classification loss is defined as Section 3. The classification cost is the key to remove NMS.

| Learning rate | Train. | Infer. | AP | AP$_{50}$ | AP$_{75}$ | AP$_s$ | AP$_m$ | AP$_l$ | FPS |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 5e-5 | 1 | 1 | 35.7 | 54.3 | 38.4 | 17.9 | 39.3 | 48.6 | 67 |
| 1e-4 | 1 | 1 | 31.7 | 51.8 | 33.7 | 15.9 | 36.3 | 42.7 | 67 |
| 5e-5 | 3 | 3 | 35.4 | 55.0 | 38.0 | 17.0 | 38.8 | 49.6 | 12 |
| 1e-4 | 3 | 3 | 36.5 | 55.5 | 39.4 | 18.7 | 39.8 | 50.2 | 12 |
| 1e-4 | 3 | 1 | 36.3 | 55.4 | 39.3 | 18.4 | 39.6 | 50.1 | 67 |

Table 2: **Ablation on Multi-head Training.** Multi-head training enables a larger learning rate and obtains higher accuracy. Meanwhile, multi-head training and single-head inference achieve better accuracy and similar inference speed with baseline.

| size | AP | AP$_{50}$ | AP$_{75}$ | AP$_s$ | AP$_m$ | AP$_l$ | FPS |
| --- | --- | --- | --- | --- | --- | --- | --- |
| (512, 640) | 35.0 | 53.8 | 37.4 | 15.9 | 38.4 | 50.1 | 80 |
| (512, 853) | 35.7 | 54.3 | 38.4 | 17.9 | 39.3 | 48.6 | 67 |

Table 3: **Effect of image size.** Larger image size brings higher accuracy, but also slower inference speed.

| epoch | AP | AP$_{50}$ | AP$_{75}$ | AP$_s$ | AP$_m$ | AP$_l$ | FPS | Time |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 36 | 30.9 | 49.4 | 32.9 | 13.9 | 34.1 | 42.0 | 67 | 10h |
| 72 | 34.1 | 52.5 | 37.1 | 16.1 | 37.7 | 46.6 | 67 | 19h |
| 144 | 35.7 | 54.3 | 38.4 | 17.9 | 39.3 | 48.6 | 67 | 36h |
| 216 | 36.5 | 55.0 | 39.6 | 18.0 | 39.7 | 51.0 | 67 | 60h |

Table 4: **Effect of training schedule.** The increasing rate of detection accuracy to the training epoch starts to decrease at 144 epochs. Therefore, we choose 144 epoch as the baseline configuration.

**Multiple-head Training.** Simply increasing the learning rate of OneNet degenerates the detection accuracy, shown in the top sections of Table 2. The multi-head training strategy is aimed to enable the model to benefit from a larger learning rate, shown in the bottom sections of Table 2. Multi-head increases the final detection accuracy but also slower inference speed. Surprisingly, multi-head training and single-head inference achieve better accuracy and similar inference speed with baseline. As a result, no additional computational cost is introduced in inference, and the detection accuracy is effectively improved.

**Image size.** The effect of image size on detection accuracy and inference are shown in Table 3. The first number in column "size" is the shortest side, the second is the maximum pixel of the longest side. Larger image size brings higher accuracy but also slower inference speed.

**Training schedule.** Longer training schedule yields higher accuracy but consumes more time. In Table 4, we show the effect of the training schedule. The increased rate of detection accuracy to the training epoch starts to decrease at 144 epoch. Therefore, we choose 144 epoch as the baseline configuration to obtain a better trade-off of detection accuracy and training time.

### 5.3. Comparison to CenterNet

We compare OneNet to one of the most popular one-stage object detectors, CenterNet [41]. In general, OneNet can be seen as replacing the original label assignment strategy in CenterNet with Minimum Cost Assignment. From Table 5, OneNet achieves comparable performance in both detection accuracy and inference speed with CenterNet. This shows that OneNet successfully removes NMS post-procedure and achieves to end-to-end one-stage object detector.

### 5.4. Label assignment in sparse detectors

Since classification cost is the key to end-to-end dense detectors' success, we further conduct experiments to study

| Method | Backbone | Size | AP | $AP_{50}$ | $AP_{75}$ | $AP_s$ | $AP_m$ | $AP_l$ | FPS |
|---|---|---|---|---|---|---|---|---|---|
| CenterNet [41, 38] | ResNet-18 | (512, 512) | 29.8 | 46.6 | 31.6 | 10.9 | 33.1 | 44.9 | 113 |
| | ResNet-50 | (512, 512) | 35.0 | 53.5 | 37.2 | 15.0 | 40.4 | 52.1 | 71 |
| | ResNet-101 | (512, 512) | 36.8 | 55.4 | 38.9 | 16.1 | 42.4 | 55.5 | 50 |
| OneNet | ResNet-18 | (512, 853) | 29.5 | 45.7 | 31.3 | 14.4 | 31.5 | 40.1 | 109 |
| | ResNet-50 | (512, 640) | 35.0 | 53.8 | 37.4 | 15.9 | 38.3 | 50.1 | 80 |
| | ResNet-50 | (512, 853) | 35.7 | 54.3 | 38.4 | 17.9 | 39.3 | 48.6 | 67 |
| | ResNet-101 | (512, 853) | 37.7 | 57.0 | 40.6 | 40.6 | 41.1 | 52.0 | 50 |

Table 5: **Comparison of CenterNet and OneNet.** OneNet achieves comparable performance in both detection accuracy and inference speed with CenterNet.

| Method | Location cost | Classification cost | AP | $AP_{50}$ | $AP_{75}$ | $AP_s$ | $AP_m$ | $AP_l$ | AP(+NMS) |
|---|---|---|---|---|---|---|---|---|---|
| Deformable DETR [42] | ✓ | | 12.0 | 19.5 | 12.4 | 13.1 | 14.9 | 8.8 | 23.8 (+11.8) |
| | ✓ | ✓ | 44.4 | 63.6 | 48.7 | 27.1 | 47.6 | 59.6 | 44.3 (-0.1) |
| Sparse R-CNN [31] | ✓ | | 20.1 | 29.7 | 21.2 | 18.0 | 22.1 | 22.7 | 33.1 (+13.0) |
| | ✓ | ✓ | 45.0 | 64.1 | 49.0 | 27.8 | 47.6 | 59.7 | 44.9 (-0.1) |

Table 6: **Effect of classification cost on sparse detectors.** Without classification cost, sparse detectors significantly drop the detection accuracy and heavily rely on NMS post-procedure.

the effect of classification cost on sparse detectors, including Deformable DETR [42] and Sparse R-CNN [31]. The implementation details follow the original papers, and the coefficient of classification cost is set to 0.

From Table 6, the classification cost also plays a vital role in removing NMS in sparse detectors. Without classification cost, NMS improves the performance by a remarkable margin. Whereas adding classification cost eliminates the necessity of NMS.

The experimental results of sparse detectors are similar to that of dense detectors. Therefore, we identify cost classification is the core of end-to-end object detector for both sparse and dense detectors.

## 5.5. Discussion

We provide a preliminary discussion about *why the proposed Minimum Cost Assignment strategy can successfully eliminate duplicated results and avoid NMS?*

We point out that there exists **misalignment** between label assignment and network optimization objective in previous methods. Generally speaking, object detectors' optimization objective is the sum of the classification cost and localization cost. However, in one-to-one Box Assign and Point Assign, the positive sample is selected for only minimizing location cost, *e.g.* box IoU or point distance. In such a case, the selected positive sample can minimize the regression cost but cannot minimize the classification cost, especially for some objects with irregular shapes or poses. In other words, there may exist another potential positive

sample that can contribute to the minimum total cost but is forcibly treated as negative samples. As a result, the network is trained into a sub-optimal status: both the "real" positive sample and selected positive sample are classified as high scores. During inference, it is reasonable to predict redundant positive results.

On the contrary, one-to-one Minimum Cost Assignment ensures the selected positive sample has the lowest total cost (classification and regression), and no other samples have a lower cost. In this way, the network's training is more effective and is finally optimized into a desirable status, where each object has one prediction.

## 6. Conclusion

Assigning labels by location cost are conceptually intuitive and offers excellent performance in object detection to date. However, we surprisingly discover that this widely-used method is the obstacle of the end-to-end detector. On the contrary, when adding classification cost, a straightforward label assignment method of the minimum cost could effectively eliminate NMS. We hope our work could inspire re-thinking the label assignment in object detection and exploring the next generation of object detectors.

# References

[1] Zhaowei Cai and Nuno Vasconcelos. Cascade R-CNN: Delving into high quality object detection. In *CVPR*, 2018. 1, 3

[2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-End object detection with transformers. In *ECCV*, 2020. 1, 2, 4, 6

[3] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-FCN: Object detection via region-based fully convolutional networks. In *NeurIPS*, 2016. 3

[4] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *CVPR*, 2005. 1, 3

[5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009. 6

[6] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. CenterNet: Keypoint triplets for object detection. In *ICCV*, 2019. 3

[7] Kaiwen Duan, Lingxi Xie, Honggang Qi, Song Bai, Qingming Huang, and Qi Tian. Corner proposal network for anchor-free, two-stage object detection. In *ECCV*, 2020. 3

[8] Mark Everingham, Luc. Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (VOC) challenge. *IJCV*, 88(2):303–338, 2010. 3

[9] Pedro Felzenszwalb, Ross Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part based models. *T-PAMI*, 32(9):1627–1645, 2010. 1

[10] Ross Girshick. Fast R-CNN. In *ICCV*, 2015. 1, 3

[11] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014. 1, 3

[12] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask R-CNN. In *ICCV*, 2017. 3

[13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3, 5, 6

[14] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, 2017. 3

[15] Lichao Huang, Yi Yang, Yafeng Deng, and Yinan Yu. DenseBox: Unifying landmark localization with end to end object detection. *arXiv preprint arXiv:1509.04874*, 2015. 3

[16] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolutional neural networks. In *NeurIPS*, 2012. 3

[17] Hei Law and Jia Deng. CornerNet: Detecting objects as paired keypoints. In *ECCV*, 2018. 3

[18] Tsung-Yi Lin, Piotr Dollar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 3, 5

[19] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. Focal loss for dense object detection. In *ICCV*, 2017. 1, 2, 3, 4, 5

[20] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014. 3, 6

[21] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C. Berg. SSD: Single shot multibox detector. In *ECCV*, 2016. 1, 3

[22] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018. 6

[23] Xin Lu, Buyu Li, Yuxin Yue, Quanquan Li, and Junjie Yan. Grid R-CNN. In *CVPR*, 2019. 3

[24] Alejandro Newell, Zhiao Huang, and Jia Deng. Associative embedding: End-to-end learning for joint detection and grouping. In *Advances in neural information processing systems*, pages 2277–2287, 2017. 3

[25] Jiangmiao Pang, Kai Chen, Jianping Shi, Huajun Feng, Wanli Ouyang, and Dahua Lin. Libra R-CNN: Towards balanced learning for object detection. In *CVPR*, 2019. 3

[26] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 1, 2, 3, 5

[27] Joseph Redmon and Ali Farhadi. YOLO9000: Better, faster, stronger. In *CVPR*, 2017. 1

[28] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 1, 3

[29] Hamid Rezatofighi, 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 *CVPR*, 2019. 4

[30] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 3

[31] Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chenfeng Xu, Wei Zhan, Masayoshi Tomizuka, Lei Li, Zehuan Yuan, Changhu Wang, et al. Sparse r-cnn: End-to-end object detection with learnable proposals. *arXiv preprint arXiv:2011.12450*, 2020. 2, 4, 6, 9

[32] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 3

[33] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. FCOS: Fully convolutional one-stage object detection. In *ICCV*, 2019. 1, 2, 3, 5

[34] Jasper RR Uijlings, Koen EA Van De Sande, Theo Gevers, and Arnold WM Smeulders. Selective search for object recognition. *IJCV*, 104(2):154–171, 2013. 3

[35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 4

[36] Paul Viola and Michael Jones. Rapid object detection using a boosted cascade of simple features. In *Proceedings of the 2001 IEEE computer society conference on computer vision*

and pattern recognition. CVPR 2001, volume 1, pages I–I. IEEE, 2001. 1, 3

[37] Thang Vu, Hyunjun Jang, Trung X Pham, and Chang D Yoo. Cascade RPN: Delving into high-quality region proposal network with adaptive convolution. In *NeurIPS*, 2019. 3

[38] Feng Wang. Centernet-better. https://github.com/FateScript/CenterNet-better, 2020. 9

[39] Hongkai Zhang, Hong Chang, Bingpeng Ma, Naiyan Wang, and Xilin Chen. Dynamic R-CNN: Towards high quality object detection via dynamic training. In *ECCV*, 2020. 1

[40] 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 *CVPR*, 2020. 1, 3

[41] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019. 1, 2, 3, 4, 5, 6, 8, 9

[42] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 2, 4, 6, 9