# Anchor-Intermediate Detector: Decoupling and Coupling Bounding Boxes for Accurate Object Detection

Yilong Lv, Min Li, Yujie He
Xi'an Institute of High Technology

{rfuyll@yeah.net, proflimin@163.com, ksy5201314@163.com}

Shaopeng Li
Tsinghua University
Department of Automation

sp-li16@mails.tsinghua.edu.cn

Zhuzhen He
National University of Defense Technology

hezhuzhen@163.com

Aitao Yang
Xi'an Institute of High Technology

824360083@qq.com

## Abstract

*Anchor-based detectors have been continuously developed for object detection. However, the individual anchor box makes it difficult to predict the boundary's offset accurately. Instead of taking each bounding box as a closed individual, we consider using multiple boxes together to get prediction boxes. To this end, this paper proposes the **Box Decouple-Couple(BDC) strategy** in the inference, which no longer discards the overlapping boxes, but decouples the corner points of these boxes. Then, according to each corner's score, we couple the corner points to select the most accurate corner pairs. To meet the BDC strategy, a simple but novel model is designed named the **Anchor-Intermediate Detector(AID)**, which contains two head networks, i.e., an anchor-based head and an anchor-free **Corner-aware head**. The corner-aware head is able to score the corners of each bounding box to facilitate the coupling between corner points. Extensive experiments on MS COCO show that the proposed anchor-intermediate detector respectively outperforms their baseline RetinaNet and GFL method by ∼2.4 and ∼1.2 AP on the MS COCO test-dev dataset without any bells and whistles. Code is available at: https://github.com/YilongLv/AID.*

## 1. Introduction

Object detection is a fundamental and challenging task in computer vision to classify and localize objects in images. Recently, as transformer has achieved good results in machine translation, it has begun to extend into the field of computer vision with success in tasks such as image classification [25, 23, 13, 31] and object detection [2, 6, 37, 8]. However, most current mainstream detectors are still based



(a) Raw Image
(b) Anchor-based
(c) Keypoint-base Anchor-free
(d) Our AID

Figure 1: The mechanism of the current mainstream detectors. Anchor-based detectors rely on predefined anchor boxes for localization. Anchor-free detectors recombine key points to achieve localization. Our AID takes advantage of both anchor-based and anchor-free methods. It decouples the corner points of the bounding boxes and then, according to the corner's score, the corner points are coupled to select the most accurate corner pair.

on convolutional neural networks.

Common anchor-based detectors, such as Faster RCNN [28], Cascade RCNN [1], YOLO [27], and RetinaNet [20], require pre-defined dense anchor boxes to cover the whole image. Although the anchor box are widely used, it is still a lack of accuracy in locating the object's boundary, because the anchor box is weakly perceptive to the boundary

Figure 2: An illustration of the difference between the proposed BDC strategy and existing NMS. For NMS, it retains only the first-ranked predicted box. In contrast, our BDC strategy retains the top-ranked multiple predicted boxes and re-pairs the coordinates according to their corner scores. $tl$ denotes the coordinates of the top-left corner point. $br$ denotes the coordinates of the bottom-right corner point. $S_n$ denotes the score of $n$ corner point. $Score$ denotes the classification score.

as shown in figure 1b.

In contrast, some anchor-free method avoid the difficulties of the weak boundary perception. In particular, the keypoint-based anchor-free detectors such as CornerNet [17], CenterNet [9], and CenterNet++ [10]. Instead of predicting the object's center and boundary' offsets, CornerNet pioneered the corner point prediction mechanism. Concretely, the model decouples each ground truth into lefttop and right-down corner points. In this paper, we name this step with a new term, **Box decouple**. Then, in the inference, CornerNet pairs all the top-left and bottom-right corner points, and this pipeline is shown in figure 1c, thus forming some bounding boxes. Similarly, we name this step **Box couple**. Thus the CornerNet don't consider the center-to-boundary perception performance, while improving the localization accuracy. Meanwhile, the net introduces another problem. i.e., the Box couple is very challenging. Suppose the size of the feature map is $\mathbb{R}^{w \times h}$. Then, the number of random corner pairs is $(w \times h)^2$. Too many possible pairs tend to lead to many false positive results, so the average precise drops. CenterNet and CenterNet++ have similar shortcomings. Therefore, how to significantly reduce the number of paired corner points is another challenge we need to address.

Based on the above analysis, it is found that the keypoint-based anchor-free detectors can circumvent the drawback of anchor-based detectors but also have their attendant difficulties. Therefore, we pondered whether a trade-off can be reached between anchor-based and anchor-free algorithms. Concretely, we can take advantage of the anchor-free detector to improve the shortcomings of the anchor-based detector as shown in figure 1d. Also, the Box couple dilemma in anchor-free detectors is alleviated by anchor-based detectors.

To this end, this paper proposes a novel architecture named the **Anchor-Intermediate Detector(AID)**, which is based on the mainstream detection method, including the anchor-based and anchor-free head. First, the anchor-based detection head maintains the conventional training pipeline. Then, we introduce an anchor-free corner-aware head for scoring the corner points of the bounding boxes, making it possible to enhance the boundary perception of the bounding boxes. In detail, during training, the corner-aware head generates two corner-aware heat maps for predicting the distribution of the object's top-left and bottom-right corner points. Similar to CornerNet, but we don't have to predict offsets and classification scores. The AID innovatively integrates two representative detection head frameworks, while the two heads are trained in parallel with each other.

In the inference, we propose a novel post-processing strategy, named **Box Decouple-Couple (BDC) strategy**. We use the proposed strategy to re-pair the corner points of each prediction box and its overlapping boxes to get more accurate localization results, as shown in figure 2 . In addition, we take into account that the classification score may not be consistent with the localization accuracy, resulting in the prediction results being not most accurately localized. Coincidentally, the corner score can be used as the localization score based on the corner-aware heat map. So, **Corner-Classification(CoCl) score** are presented, consisting of the classification score and corner confidence for ranking in BDC strategy.

**The main contributions of our work can be summarized as follows**

- By analyzing the disadvantages of the current anchor-based and anchor-free models, we propose the novel

AID, which can achieve a trade-off between the two detection frameworks to improve the detection accuracy of the model.

- We have redesigned the training and inference pipelines separately. In training, the anchor-based detection head and the corner-aware head share the backbone network and the feature pyramid network. Both are trained synchronously. During inference, we first propose a new corner-classification score for post-processing. Then, our proposed BDC strategy rethinks the value of each prediction box and its overlapping boxes, from which we refine predictions with better localization quality.

- The AID uses the corner-aware head as the anchor-free head, and the anchor-based head uses some main methods, including RetinaNet, GFL, etc., while based on some backbone networks, including ResNet-50, ResNet-101, ResNeXt-101, etc. Also, our method achieved state-of-the-art results on the MS COCO dataset.

## 2. Related Work

### 2.1. Anchor-based detector

**Two-stage object detection**. It first started with Faster RCNN, introducing anchor boxes into Fast RCNN [12] using sliding windows and candidate regions. The first stage generates a set of region proposals, and the second stage performs classification and localization fine-tuning on these region proposals. Based on this, many algorithms [14, 5, 4, 29, 15] are proposed to improve its performance, feature fusion, attention mechanism, multi-scale training, and training strategy. The two-stage algorithms generally have better detection accuracy but are relatively slow.

**Single-stage object detection.** For faster detection, single-stage methods have emerged. Instead of relying on RPNs, single-stage detectors directly localized and classified regions of interest in images. SSD [22] was the first approach to propose a single-stage detection strategy, which has attracted much attention due to its efficient training and inference. Since then, many redesigned architectures of single-stage object detectors [27, 20, 19, 35, 36] have been proposed. Based on them, many methods have been proposed to improve the performance of single-stage detectors, feature fusion, label assignment strategies, detection head network structures, loss functions, and localization refinement. Currently, there are more research results on single-stage than two-stage.

### 2.2. Anchor-free detector

**Keypoint-based detectors.** The keypoint-based approach focuses on extreme locations of instances, such as corner points and extreme points. CornerNet is one of the representative methods. The improved CornerNet-lite [18] introduces CornerNet-sweep and CornerNet-squeeze to improve its speed. CenterNet adds center points to the top-left and bottom-right corners to provide the ability to perceive internal information, thus improving precision and recall. ExtremeNet [38] even increases the number of key points to the topmost, bottommost, leftmost, rightmost, and center points. The training pipeline of the anchor-free approach [24, 32] is similar to CornerNet.

**Center-based detector.** These anchor-free methods have similarities with anchor-based methods in that both predict from the center to the boundary. The FCOS [30] method designs a new centrality branching detector head. The centerness scores of every location within the ground truth are also defined. And the label assignment strategy is set accordingly. Different from FCOS, FoveaBox [16] does not add any new branching network. It regards every position within the subregions of the ground truth as positive and performs label assignment. FSAF [39] connects an anchor-free branch with online feature selection to the RetinaNet. The newly added branch defines the central region of the object as positive and locates it by predicting the four distances to its boundaries. There are also center-based anchor-free detectors [33, 11, 34] that play an important role in object detection.

## 3. Proposed method

In this section, we describe the the AID in detail, including the anchor-based and corner-aware head, as shown in Figure 3. The anchor-based detection head we use is RetinaNet, but it is not limited to this one. We first introduce the model structure and the training pipeline. Then, in inference, the rules of BDC strategy will be presented in detail.

### 3.1. Anchor-Intermediate Detector

Figure 3 shows our proposed AID, including the corner-aware head and standard anchor-based head. Following the style of multi-task learning, the corner-aware head $f_{cor}$ accepts the features $\mathcal{F}_n$ from the FPN and then learns the corner-aware heat map $\mathcal{M}$. $\mathcal{F}_n$ is the the $n$ $(1 <= n <= N)$ -th level feature at the FPN. The lower the layer, the smaller the $n$.

Note that we first resize feature so that the multi-scale feature becomes a single scale, as the higher-resolution feature facilitates the differentiation of corner points. Then they were concatenated by channel, which is shown in Eq. 1. Thus, we could obtain single-scale fused features with the architecture shown in Figure 3 of the anchor-intermediate model.

$$\mathcal{F} = Cat(Resize(\mathcal{F}_n|n = 1, \ldots, N) \tag{1}$$

Figure 3: Pipeline of the AID. The left part shows the overall detection model which consists of backbone, FPN and detection head. The upper part is the conventional anchor-based head, which outputs classification scores ($H \times W \times C$) and boundary prediction ($H \times W \times 4$), respectively. The lower part is the corner-aware head, which outputs the top-left and bottom-right corner heat maps ($H \times W \times C$), respectively. In the box decouple, the top-ranked boxes are processed uniformly to separate their corner points. These corner points are scored according to the corner point heat map. Finally, the top-left and bottom-right corner points with good scores are coupled.



Figure 4: Network structure of the corner-aware head.

For better recognition of corner points, we separately predict the two extreme corners where corner pooling is used. The detection head finally predicts the corner feature heat map $\mathcal{M} = (\mathcal{M}_{tl}, \mathcal{M}_{br})$, including the top-left and bottom-right corner heat maps as $\mathcal{M}^{tl}$ and $\mathcal{M}^{br} \in \mathbb{R}^{w \times h \times c}$. Each set of heatmaps has $C$ channels, which represent the object's categories.

The structure of the corner-aware head follow the CornerNet, as shown in Figure 4. The $1 \times 1$ convolution and activation $\delta$ are performed on the fused features $\mathcal{F}$ according to each direction. Then direction pooling $\mathcal{P} = \{\mathcal{P}_t, \mathcal{P}_b, \mathcal{P}_r, \mathcal{P}_l\}$ is performed separately to obtain the direction-aware feature maps $\{\mathcal{F}_t, \mathcal{F}_b, \mathcal{F}_r, \mathcal{F}_l\}$, and the

process is shown in Eq. 2.

$$\{\mathcal{F}_t, \mathcal{F}_l, \mathcal{F}_b, \mathcal{F}_r\} = \mathcal{P}(\delta(Conv(\mathcal{F}))), \tag{2}$$

The top-left and bottom-right features are element-wise summed. Then, we use the convolution operation on the features. The feature $\mathcal{F}$ are added to the corner features to prevent a vanishing gradient. Before outputting the heat map, we will perform sigmoid $\sigma$ activation on each element of the feature map so that all positions of the feature map will be between 0 and 1, which represent the confidence. The above is as follows:

$$\begin{aligned}
\mathcal{M}_{tl} &= \sigma(Conv(\mathcal{F}_t \oplus \mathcal{F}_l) + Conv\mathcal{F}) \\
\mathcal{M}_{br} &= \sigma(Conv(\mathcal{F}_b \oplus \mathcal{F}_r) + Conv\mathcal{F})
\end{aligned} \tag{3}$$

As for the training loss, suppose $\mathcal{M}_{cor} = \{\mathcal{M}_{tl}, \mathcal{M}_{br}\}$ is the corner score in the top-left and bottom-right corner-aware heat map. And let $\mathcal{Y}_{cor}$ be the labeled corner heat map expanded with an unnormalized Gaussian function. The loss function of corner-aware heat map is as follows:

$$L_{corner} = \begin{cases} (1 - \mathcal{M}_{cor})^\alpha log(\mathcal{M}_{cor}) & if\ \mathcal{Y}_{cor} = 1 \\ (1 - \mathcal{Y}_{cor})^\beta (\mathcal{M}_{cor})^\alpha log(1 - \mathcal{M}_{cor}) & otherwise \end{cases} \tag{4}$$

Next, the total loss is optimized as shown in Eq. 5.

$$L_{total} = L_{reg} + L_{cls} + \lambda L_{corner} \tag{5}$$

## 3.2. CoCl score for post-processing

In the inference, the overlapping boxes are suppressed based on the classification score $\mathcal{S}_{cls}$ only. However, the classification may not coincide with the localization, which leads to the suppression of the bounding boxes with high localization accuracy when the classification score is low.

We consider corner score as another important basis for post-processing, which is able to integrate the localization information. The corner's score comes from the top-left and bottom-right corner-aware heat map. Therefore, The new evaluation score is named CoCl score, as follows:

$$CoCl = \mathcal{S}_{cls} \times \mathcal{F}(\mathcal{M}_{tl}, \mathcal{M}_{br}) \tag{6}$$

Where $\mathcal{F}$ denotes the integration between $\mathcal{M}_{tl}$ and $\mathcal{M}_{br}$, see section 4.3.3 for details. The results has the same size as $\mathcal{S}_{cls} \in \mathbb{R}^{ijwh \times C}$.

## 3.3. Box Decouple-Couple strategy

**Box Decouple.** In NMS, the overlapping boxes are suppressed and then discarded. If a predicted box with highest classification scores fail to locate most accurately, then its overlapping boxes are likely located accurately because they have a large Intersection-over-Union (IoU) with this predicted box.

The conventional NMS treats each bounding box as an independent individual, but it does not predict the object's location well, especially in those with blurred boundaries. Therefore, we decouple all bounding boxes $\mathcal{B}$ as $\{\mathcal{S}_{tl}, \mathcal{S}_{br}\}$. The details are shown in Eq. 7, where $n$ denotes all predicted bounding boxes, $\mathcal{S}_{tl}$ denotes the top-left point set and $\mathcal{S}_{br}$ denotes the bottom-right point set.

$$\mathcal{B} = (\overbrace{x_1, y_1}^{S^{tl}}, \overbrace{x_2, y_2}^{S^{br}}), \tag{7}$$
$$\mathcal{S}_{tl} = \{s_i^{tl}\}_{i=1,\ldots,n}, \quad \mathcal{S}_{br} = \{s_i^{br}\}_{i=1,\ldots,n}.$$

For the predicted bounding box, suppose the box $\mathcal{P}$ corresponds to the overlapping box $\mathcal{O} = \{b \mid IoU(b, \mathcal{P}) > \tau, b \in \mathcal{B}\}$. We decouple the predicted box $\mathcal{P}$ and its overlapping boxes $\mathcal{O}$ into the top-left corner set $\{s_p^{tl}, s_{o1}^{tl}, \ldots, s_{oi}^{tl}\} \subseteq \mathcal{S}_{tl}$ and bottom-right corner sets $\{s_p^{br}, s_{o1}^{br}, \ldots, s_{oi}^{br}\} \subseteq \mathcal{S}_{br}$. Then, we map these corner points onto the corner-aware heat map $\mathcal{M}$ and the process is as follows:

$$
\begin{aligned}
(f_p^{br}, f_{o1}^{br}, \ldots, f_{oi}^{br}) &= f_{\mathcal{P}, \mathcal{O} \mapsto \mathcal{M}_{br}}(s_p^{br}, s_{o1}^{br}, \ldots, s_{oi}^{br}) \\
(f_p^{tl}, f_{o1}^{tl}, \ldots, f_{oi}^{tl}) &= f_{\mathcal{P}, \mathcal{O} \mapsto \mathcal{M}_{tl}}(s_p^{tl}, s_{o1}^{tl}, \ldots, s_{oi}^{tl})
\end{aligned} \tag{8}
$$

Therefore, we use box decoupling to transition the anchor-based inference to anchor-free inference, as shown in Figure 3 of box decouple. Here, these decoupled corner points are very similar to those in CornerNet. Next, we use the idea of anchor-free to process these corner points. It is worth emphasizing that our method dramatically reduces the number of corner point pair in our method, thus improving the average precise.

**Box Couple.** Since the heat maps $\mathcal{M}$ are downsampled, each element in the corner-aware heat map does not correspond to the original image one by one, but is mapped to a region of image. In addition, the corner-aware heat map is fitted with a Gaussian model in the training, so the confidence for pixels at the same distance from the center is the same. Therefore, there are some locations in the heat map with the same confidence. To this end, we chose multiple decoupled corner points and use them to obtain a new prediction box. This process is specified in Eq. 9.

$$
\begin{aligned}
\mathcal{T}^{br} &= \operatorname*{arg\,max}_{s_p^{br}, s_{o1}^{br}, \ldots, s_{oi}^{br}} \overset{n}{top k}(f_p^{br}, f_{o1}^{br}, \ldots, f_{oi}^{br}) \\
\mathcal{T}^{tl} &= \operatorname*{arg\,max}_{s_p^{tl}, s_{o1}^{tl}, \ldots, s_{oi}^{tl}} \overset{n}{top k}(f_p^{tl}, f_{o1}^{tl}, \ldots, f_{oi}^{tl})
\end{aligned} \tag{9}
$$

Finally, the new corner points obtained are combined to form an updated bounding box $\mathcal{B}_{update}$. The details are shown in Eq. 10.

$$\mathcal{B}_{update} = (Mean(\{s_i^{tl}\}_{i=\mathcal{T}^{br}}), \ Mean(\{s_i^{tl}\}_{i=\mathcal{T}^{tl}})) \tag{10}$$

Following this method, the output boxes not only contain classification and localization's information but also improve the accuracy of localization. The detailed procedure of the BDC strategy is shown in Algorithm 1.

---

**Algorithm 1** Box Decouple-Couple strategy.

---

TRAIN: Corner heat map $(\mathcal{M}_{tl}, \mathcal{M}_{br})$
    **select** $\mathcal{S}_{ijc}$ for $i, j, c, \mathcal{M}$ in $W, H, C, [\mathcal{M}_{tl}, \mathcal{M}_{br}]$
    **if** $\mathcal{Y}_{ijc} == 1$:
        $(1 - \mathcal{S}_{ijc})^\alpha log(\mathcal{S}_{ijc}) \mapsto 0$
    **else:**
        $(\mathcal{S}_{ijc})^\alpha log(1 - \mathcal{S}_{ijc}) \mapsto 0$

PREDICT:
    $CoCl = \mathcal{S}_{cls} \times \mathcal{F}(\mathcal{M}_{tl}, \mathcal{M}_{br})$
    $\mathcal{B}^p, \mathcal{B}^o = $ **NMS**$(CoCl)$
    # $\mathcal{B}^p \in \mathbb{R}^{1 \times 4} \Rightarrow$ **Prediction Box**
    # $\mathcal{B}^o \in \mathbb{R}^{N \times 4} \Rightarrow$ **Overlapping Box**
    $\{\mathcal{S}^{tl}, \mathcal{S}^{br}\} = $ **Decouple**$(\mathcal{B}_n^p, \mathcal{B}_n^o)$
    $(f^{br}, f^{tl}) = f_{\mathcal{P}, \mathcal{O} \mapsto \mathcal{M}}(\mathcal{S}^{tl}, \mathcal{S}^{br})$
    **select the top-n** $[f^{br}, f^{tl}] \mapsto [f_{top-n}^{br}, f_{top-n}^{tl}]$
    $(\mathcal{S}_{top-n}^{tl}, \mathcal{S}_{top-n}^{br}) = f_{\mathcal{M} \mapsto \mathcal{P}, \mathcal{O}}(f_{top-n}^{br}, f_{top-n}^{tl})$
    $\mathcal{B}_{update} = $ **Couple**$(\mathcal{S}_{top-n}^{tl}, \mathcal{S}_{top-n}^{br})$
    **return** $\mathcal{B}_{update}$

---

# 4. Experiment

## 4.1. Dataset

We evaluate our method on the MS-COCO dataset [21] according to the commonly used settings. MS-COCO contains about 160K images of 80 classes. The dataset was partitioned into training 2017, val2017, and test 2017 subsets with 118K, 5K, and 41K images, respectively. The standard average precision (AP) metric reports the results at different IoU thresholds and target scales. We trained only on the 2017 train images in all our experiments without using any additional data. For the experiments of the ablation study, we evaluated the performance on a subset of val2017. Compared to state-of-the-art methods, we report the official results returned from the test server on the test-dev subset.

## 4.2. Implementation details

We use the generic "Backbone - FPN - Head" as our pipeline and the MMdetection toolbox [3] to implement our method. All models are trained on 4 TESLA A100 GPUs with four small batches per GPU. Unless otherwise specified, we used a stochastic gradient descent (SGD) optimizer with a weight decay of 0.0001 and momentum of 0.9. The initial learning rate was set to 0.01, and training was started using a linear warm-up strategy. We initialize our backbone network with the weights pre-trained on ImageNet [7].

## 4.3. Ablation study

This section will perform detailed ablation experiments on the proposed method. Besides exploring the effects of different baselines and backbone networks on the experimental results, the rest of the experiments are performed on the RetinaNet method based on ResNet-101 and trained for 12 epochs. The final validation is performed on the COCO-val2017 dataset.

### 4.3.1 Corner-aware head

In this section, we discuss the results after adding corner-aware head at different positions of the model, and the experimental results are shown in Table 1. From the table, the performance of our method improves by 0.9% AP on the classification branch and by 0.6% AP on the regression branch.

| Method | AP | AP50 | AP75 | APS | APM | APL |
|---|---|---|---|---|---|---|
| Baseline | 38.5 | 57.6 | 41 | 21.7 | 42.8 | 50.4 |
| + Cor head - FPN | 39.8(+1.3) | 58.7 | 42.4 | 21.8 | 43.5 | 52.8 |
| + Cor head - cls | 39.4(+0.9) | 58.4 | 42.4 | 21.4 | 43.5 | 53.2 |
| + Cor head - reg | 39.1(+0.6) | 58.6 | 42.1 | 21.4 | 42.9 | 52.6 |

Table 1: Comparison of performances when applying the our method to each position of the baseline model.

We perform the analysis. Corner-aware head belongs to the anchor-free head. In contrast, the classification and regression branches belong to the anchor-based head, and it seems more common sense for them to be trained independently. Therefore, we also separate the Corner-aware head network from the two branches and connect it to the back of the FPN. This way, the Corner-aware head and Cls-Reg branches will present a parallel structure. The performance of this structure is improved by 1.3% AP. The experimental results illustrate that separating the Corner-aware head from the classification and regression branches can improve the model's performance by making them independent.

### 4.3.2 Hyperparameters $\lambda$

Further, we discuss the effect of $\lambda$ in Eq. 5 on the experimental results, shown in Table 2. When the hyperparameters $\lambda$ is equal to 0, it indicates a baseline method without the corner-aware head, which has a detection performance of 38.5 AP. The detection performance of the AID reaches the highest 40.1 AP when the hyperparameters $\lambda$ is equal to 0.3. Meanwhile. The AP decreases as the $\lambda$ keeps increasing, which indicates that if the corner-aware head is trained with large weights, it will adversely affect the training of the traditional detection training loss, thus leading to a decrease in detection performance.

| Weight factor | AP | AP50 | AP75 | APS | APM | APL |
|---|---|---|---|---|---|---|
| 0(baseline) | 38.5 | 57.6 | 41 | 21.7 | 42.8 | 50.4 |
| 0.1 | 39.8(+1.3) | 58.7 | 42.4 | 21.8 | 43.5 | 52.8 |
| 0.3 | 40.1(+1.6) | 58.8 | 42.8 | 21.8 | 43.6 | 53.2 |
| 0.5 | 39.4(+0.9) | 58.1 | 42.2 | 22.2 | 42.5 | 52.2 |
| 0.8 | 39.0(+0.5) | 56.8 | 40.7 | 20.5 | 41.4 | 50.7 |
| 1 | 38.2(-0.3) | 56.8 | 40.7 | 20.5 | 41.4 | 50.7 |

Table 2: Peformance of the AID when changing the hyper-parameters $\lambda$ of the total loss. $\lambda$ weighting means weighting the loss of the corner-aware head.

### 4.3.3 CoCl score

During the inference, the detection confidence is calculated according to the classification score. We next discuss the contribution of several classical forms of $\mathcal{F}(\mathcal{M}_{tl}, \mathcal{M}_{br})$ to the detection accuracy, and the experimental results are shown in Table 4. First, the detection result obtained by multiplying the classification score and the corner score is 39.6, when the corner score take the maximum output of the top-left and bottom-right corner-aware heat map. Also, calculating the minimum output of the two corner-aware heat map as the corner score is a scheme with a detection accuracy of 39.6 AP. Finally, the detection performance AP is 39.8 when the average output of both heat map is calculated.

| Method | Backbone | AID | | PAFPN | ATSSAssigner | AP | AP50 | AP75 | APS | APM | APL |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | CoCl score | BDC | | | | | | | | |
| Retinanet | | | | | | 38.5 | 57.6 | 41 | 21.7 | 42.8 | 50.4 |
| AID | ResNet-101 | ✔ | | | | 39.8(+1.3) | 58.7 | 41.1 | 20.6 | 42 | 51.2 |
| | | | ✔ | | | 39.1(+0.6) | 58.3 | 41.5 | 21.6 | 43.1 | 51.8 |
| | | ✔ | ✔ | | | 40.0(+1.5) | 58.7 | 42.1 | 31.6 | 43.3 | 52.8 |
| | | ✔ | ✔ | ✔ | | 40.1(+1.6) | 58.8 | 43 | 22.1 | 43.9 | 53.1 |
| | | ✔ | ✔ | | ✔ | 40.6(+2.1) | 57.3 | 43.9 | 22.9 | 44.4 | 52.3 |
| | | ✔ | ✔ | ✔ | ✔ | 40.8(+2.3) | 57.4 | 44.3 | 23.5 | 44.4 | 52.4 |

Table 3: Individual contribution of the components in our method. The first row represents the results of the baseline trained with the focal loss. PAFPN is used as a more powerful feature pyramid network. ATSSAssigner is used to replace the conventional MaxIoUAssigner.

| id | F | AP | AP50 | AP75 |
|---|---|---|---|---|
| 1 | $\mathcal{S} \times e^{avg(\mathcal{M}_{tl}, \mathcal{M}_{br})}$ | 39.8 | 58.7 | 42.4 |
| 3 | $\mathcal{S} \times e^{max(\mathcal{M}_{tl}, \mathcal{M}_{br})}$ | 39.6 | 58.6 | 42.2 |
| 3 | $\mathcal{S} \times e^{min(\mathcal{M}_{tl}, \mathcal{M}_{br})}$ | 39.6 | 58.5 | 42.2 |
| 4 | $\mathcal{S}^0 \times avg(\mathcal{M}_{tl}, \mathcal{M}_{br})^1$ | 39.0 | 58.3 | 41.5 |
| 5 | $\mathcal{S}^{0.3} \times avg(\mathcal{M}_{tl}, \mathcal{M}_{br})^{0.7}$ | 39.9 | 58.5 | 42.7 |
| 6 | $\mathcal{S}^{0.5} \times avg(\mathcal{M}_{tl}, \mathcal{M}_{br})^{0.5}$ | 39.1 | 57.1 | 41.8 |
| 7 | $\mathcal{S}^{0.8} \times avg(\mathcal{M}_{tl}, \mathcal{M}_{br})^{0.2}$ | 35.7 | 52.0 | 38.2 |
| 8 | $\mathcal{S}^1 \times avg(\mathcal{M}_{tl}, \mathcal{M}_{br})^0$ | 28.5 | 44.3 | 29.2 |

Table 4: Comparison of performances of different CoCl score functions.

In addition, We also perform a weighted average of the classification score and the corner score, using a balance parameter $\alpha$ to adjust the ratio between them, as shown in Eq. 11. When $\alpha$ equals to 1, the detection confidence degenerates to the classification score. Similarly, the detection confidence degenerates to the corner score when we $\alpha$ equals to 0. In the inference, we studied the different parameters. The experimental results are shown in Table 4. From the experimental results, the method in this paper can achieve 0.3 when $\alpha$ is equal to 39.9.

$$CoCl = \mathcal{S}_{cls}^{\alpha} \times avg(\mathcal{M}_{tl}, \mathcal{M}_{br})^{\alpha} \tag{11}$$

#### 4.3.4 Box Decouple-Couple strategy

In the inference, for the prediction boxes and the overlapping boxes, we next investigate the contribution of the strategies of different box coupling to the experimental results, which are shown in Table 5. The simple strategy is to select the highest detection score, which has a detection accuracy of 39.7 AP, 1.2 AP higher than the baseline method. We continued our study by averaging all the top-left and bottom-right corner points. Its detection result is 36.0 AP, which is lower than the performance of the baseline method. We analyze that because some corner points with small detection scores are not suitable for predicting object. There-

fore, we consider the sum of the mean and deviation of the detection scores as the threshold value $\tau_{score}$. The positions of the corner points with detection scores larger than $\tau_{score}$ are averaged, and the detection performance under this strategy is 40.0 AP, which is 1.5 AP higher than the baseline method. Also, we average the corner points with the top n(=10) scores, resulting in 39.8 AP, which is 1.3 AP higher than the baseline method.

| id | F | AP | AP50 | AP75 | APS | APM | APL |
|---|---|---|---|---|---|---|---|
| 1 | Top-n(=10) | 39.8(+1.3) | 58.6 | 41.1 | 20.6 | 42 | 51.2 |
| 2 | Max | 39.7(+1.2) | 58.7 | 42.1 | 21.6 | 43.3 | 52.8 |
| 3 | All | 36.0(-2.5) | 58.6 | 39 | 20.6 | 40.6 | 46.2 |
| 4 | $\tau_{score}(= 0.5)$ | 40.0(+1.5) | 58.7 | 42.7 | 21.7 | 43.7 | 53.4 |

Table 5: Comparison of performances of different box couple method. **Max** means to select the corner point with the maximum score. **Top-n** means select the first n maximum corner point. **All** means select all corner points. $\tau_{score}(= 0.5)$ means to select the corner point coordinates greater than the threshold (=0.5).

#### 4.3.5 Stronger components

We conducted experiments using the components of the BDC strategy, different feature pyramid networks and label assignment strategies to improve the detection performance of the AID further, which are shown in Table 3. When PAFPN is used as the feature pyramid network, our method achieves 40.1 AP, which improves 1.6 AP compared to the baseline method. In addition, we use ATSSAssigner alone instead of the conventional maximum IoU strategy. Our method improves 40.6 AP compared to the baseline method. When the PAFPN and ATSSAssigner are used together, the proposed method achieves 40.8 AP.

### 4.4. Comparison with State-of-the-Arts

For experiments comparing with the state-of-the-art dense detector on COCO test-dev2017, we train 1x and 2x

(i.e., 12 and 24 epochs) models.

To demonstrate the effectiveness of the proposed method, we have conducted a series of experiments based on two baselines and various advanced backbone networks. The results are presented in Table 6, from which it can be seen that our model achieves excellent performance. As shown in Table 6, using ResNet-101 and ResNeXt-101-64×4d, our method based on the RetinaNet model achieves 40.1∼41.7 AP, consistently outperforming current baseline. Meanwhile, using ResNet-101 and ResNeXt-101-64×4d, our method based on the GFL model achieves 45.7∼49.4 AP, which outperforms the baseline GFL methods by approximately 0.3∼1.2 AP.

## 5. Conclusion

We focus on improving the performance of the anchor-based detector. The novel AID is proposed and contains three innovations: Firstly, The corner-aware head is proposed to quantify the localization of each corner point. Then, The BDC strategy is proposed. The overlapping boxes are fully utilized to decouple the corner points, and then the corner points are re-paired. In addition, the CoCl score is proposed, which contains both classification and corner scores and can comprehensively evaluate the quality of prediction boxes. Experimental results demonstrate the effectiveness of these improvements. The performance of our method exceeds that of many state-of-the-art models. However, there are improvements in our work: The corner-aware head increases the training burden of the model, and we will consider using knowledge distillation to achieve a lighter model. And, the BDC strategy has many manually predefined hyperparameters, and adaptive hyperparameters will be designed in the future.

## 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] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 213–229. Springer, 2020.

[3] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, et al. Mmdetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.

[4] Bowen Cheng, Yunchao Wei, Honghui Shi, Rogerio Feris, Jinjun Xiong, and Thomas Huang. Revisiting rcnn: On awakening the classification power of faster rcnn. In *Proceedings of the European conference on computer vision (ECCV)*, pages 453–468, 2018.

[5] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. *Advances in neural information processing systems*, 29, 2016.

[6] Xiyang Dai, Yinpeng Chen, Jianwei Yang, Pengchuan Zhang, Lu Yuan, and Lei Zhang. Dynamic detr: End-to-end object detection with dynamic attention. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2988–2997, 2021.

[7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

[9] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6569–6578, 2019.

[10] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet++ for object detection. *arXiv preprint arXiv:2204.08394*, 2022.

[11] Chengjian Feng, Yujie Zhong, Yu Gao, Matthew R Scott, and Weilin Huang. Tood: Task-aligned one-stage object detection. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3490–3499. IEEE Computer Society, 2021.

[12] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.

[13] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing Xu, and Yunhe Wang. Transformer in transformer. *Advances in Neural Information Processing Systems*, 34:15908–15919, 2021.

[14] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.

[15] Yihui He, Chenchen Zhu, Jianren Wang, Marios Savvides, and Xiangyu Zhang. Bounding box regression with uncertainty for accurate object detection. In *Proceedings of the ieee/cvf conference on computer vision and pattern recognition*, pages 2888–2897, 2019.

[16] Tao Kong, Fuchun Sun, Huaping Liu, Yuning Jiang, Lei Li, and Jianbo Shi. Foveabox: Beyound anchor-based object detection. *IEEE Transactions on Image Processing*, 29:7389–7398, 2020.

[17] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European conference on computer vision (ECCV)*, pages 734–750, 2018.

[18] Hei Law, Yun Teng, Olga Russakovsky, and Jia Deng. Cornernet-lite: Efficient keypoint based object detection. *arXiv preprint arXiv:1904.08900*, 2019.

[19] Xiang Li, Wenhai Wang, Lijun Wu, Shuo Chen, Xiaolin Hu, Jun Li, Jinhui Tang, and Jian Yang. Generalized focal loss:

| Method | Backbone | Schedule | AP | AP50 | AP75 | APS | APM | APL |
|---|---|---|---|---|---|---|---|---|
| **Anchor-based two-stage:** | | | | | | | | |
| Faster RCNN [28] | ResNet-101 | 2x | 39.7 | 60.7 | 43.2 | 22.5 | 42.9 | 49.9 |
| Cascade R-CNN [1] | ResNet-101 | 2x | 42.8 | 62.1 | 46.3 | 23.7 | 45.5 | 55.2 |
| Mask R-CNN [14] | ResNet-101 | 2x | 39.8 | 62.3 | 43.4 | 22.1 | 43.2 | 51.2 |
| R-FCN [5] | ResNet-101 | 2x | 41.4 | 63.4 | 45.2 | 24.5 | 44.9 | 51.8 |
| TridentNet [26] | ResNet-101 | 2x | 42.7 | 63.6 | 46.5 | 23.9 | 46.6 | 56.6 |
| **Anchor-free Keypoint-based:** | | | | | | | | |
| CornerNet [18] | hourglass-104 | 200e | 40.5 | 56.6 | 43.1 | 18.9 | 43 | 53.4 |
| CenterNet [9] | hourglass-104 | 190e | 44.9 | 62.4 | 48.1 | 25.6 | 47.4 | 57.4 |
| ExtremeNet [38] | hourglass-104 | 200e | 40.2 | 55.5 | 43.2 | 20.4 | 43.2 | 53.1 |
| Grid R-CNN [24] | ResNet-101 | 2x | 43.2 | 63.0 | 46.6 | 25.1 | 46.5 | 55.2 |
| RepPoints [32] | ResNet-101-dcn | 2x | 45.0 | 66.1 | 49.0 | 26.6 | 48.6 | 57.5 |
| **Anchor-free Center-based:** | | | | | | | | |
| FCOS [30] | ResNet-101 | 2x | 43.0 | 61.7 | 46.3 | 26.0 | 46.8 | 55.0 |
| FoveaBox [16] | ResNeXt-101 | 2x | 42.1 | 61.9 | 45.2 | 24.9 | 46.8 | 55.6 |
| FSAF [39] | ResNext-101-64x4d | 2x | 42.9 | 63.8 | 46.3 | 26.6 | 46.2 | 52.7 |
| **Anchor-based one-stage:** | | | | | | | | |
| SSD512 [22] | VGG16 | 2x | 28.8 | 48.5 | 30.3 | 10.9 | 31.8 | 43.5 |
| RefineDet [35] | ResNet-101 | 2x | 41.8 | 62.9 | 45.7 | 25.6 | 45.1 | 54.1 |
| FreeAnchor [36] | ResNet-101 | 2x | 43.1 | 62.2 | 46.4 | 24.5 | 46.1 | 54.8 |
| **Baseline:** | | | | | | | | |
| RetinaNet [20] | ResNet-50 | 1x | 36.9 | 56.2 | 39.3 | 20.5 | 39.9 | 46.3 |
| RetinaNet | ResNet-50 | 2x | 37.7 | 57.2 | 40.2 | 20.4 | 40.2 | 48.1 |
| RetinaNet | ResNet-101 | 1x | 39 | 58.6 | 41.7 | 21.9 | 42.2 | 49.3 |
| RetinaNet | ResNet-101 | 2x | 39.6 | 59.1 | 42.4 | 21.3 | 42.6 | 51.1 |
| RetinaNet | ResNeXt-101-32x4d | 1x | 40.4 | 60.3 | 43.1 | 23 | 43.7 | 50.6 |
| RetinaNet | ResNeXt-101-64x4d | 1x | 41.4 | 61.5 | 44.4 | 24.2 | 44.8 | 52.1 |
| GFL [19] | ResNet-50 | 1x | 40.5 | 58.9 | 43.8 | 22.8 | 43.6 | 50.6 |
| GFL | ResNet-50 | 2x | 43.2 | 61.7 | 46.9 | 26.5 | 46.7 | 51.8 |
| GFL | ResNet-101 | 2x | 45.2 | 63.9 | 49.3 | 27.5 | 48.8 | 55.3 |
| GFL | ResNet-101-dcn | 2x | 47.3 | 66 | 51.5 | 28.3 | 50.9 | 59.3 |
| GFL | ResNeXt-101-dcn | 2x | 48.2 | 67.2 | 52.5 | 29.2 | 51.6 | 60.2 |
| **Ours:** | | | | | | | | |
| RetinaNet+AID | ResNet-50 | 1x | 38(+1.1) | 56.9 | 40.6 | 20.8 | 40.4 | 47.8 |
| RetinaNet+AID | ResNet-50 | 2x | 39.2(+1.5) | 58.3 | 41.8 | 21.1 | 41.4 | 49.9 |
| RetinaNet+AID | ResNet-101 | 1x | 40.1(+1.1) | 59.4 | 42.9 | 22.2 | 42.9 | 51 |
| RetinaNet+AID | ResNet-101 | 2x | 41.2(+1.6) | 60.5 | 44 | 22.3 | 43.8 | 53 |
| RetinaNet+AID | ResNeXt-101-32x4d | 1x | 42.8(+2.4) | 62.6 | 46 | 24.6 | 45.8 | 54.3 |
| RetinaNet+AID | ResNeXt-101-64x4d | 1x | 41.7(+0.3) | 61.3 | 44.6 | 23.6 | 44.4 | 52.88 |
| GFL+AID | ResNet-50 | 1x | 41.1(+0.6) | 58.9 | 44.3 | 23.2 | 43.9 | 51 |
| GFL+AID | ResNet-50 | 2x | 43.9(+0.7) | 62.1 | 47.7 | 26.7 | 47.1 | 52.8 |
| GFL+AID | ResNet-101 | 2x | 45.7(+0.5) | 64.1 | 49.7 | 27.6 | 49.3 | 55.7 |
| GFL+AID | ResNet-101-dcn | 2x | 48.4(+1.1) | 66.9 | 52.6 | 29 | 51.9 | 60.8 |
| GFL+AID | ResNeXt-101-dcn | 2x | 49.4(+1.2) | 67.9 | 53.6 | 29.6 | 53 | 62 |

Table 6: Performance comparison with the state-of-the-art methods on the MS-COCO test-dev dataset in single-model and single-scale results. Our method respectively outperforms their baseline RetinaNet and GFL method by ~2.4 and ~1.2 AP without any bells and whistles.

Learning qualified and distributed bounding boxes for dense object detection. *Advances in Neural Information Processing Systems*, 33:21002–21012, 2020.

[20] 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.

[21] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.

[22] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.

[23] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.

[24] Xin Lu, Buyu Li, Yuxin Yue, Quanquan Li, and Junjie Yan. Grid r-cnn. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7363–7372, 2019.

[25] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In *International conference on machine learning*, pages 4055–4064. PMLR, 2018.

[26] David Paz, Hengyuan Zhang, and Henrik I Christensen. Tridentnet: A conditional generative model for dynamic trajectory generation. In *Intelligent Autonomous Systems 16: Proceedings of the 16th International Conference IAS-16*, pages 403–416. Springer, 2022.

[27] 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.

[28] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015.

[29] 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. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14454–14463, 2021.

[30] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9627–9636, 2019.

[31] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 568–578, 2021.

[32] Ze Yang, Shaohui Liu, Han Hu, Liwei Wang, and Stephen Lin. Reppoints: Point set representation for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9657–9666, 2019.

[33] Mohsen Zand, Ali Etemad, and Michael Greenspan. Objectbox: From centers to boxes for anchor-free object detection. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part X*, pages 390–406. Springer, 2022.

[34] 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.

[35] Shifeng Zhang, Longyin Wen, Xiao Bian, Zhen Lei, and Stan Z Li. Single-shot refinement neural network for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4203–4212, 2018.

[36] Xiaosong Zhang, Fang Wan, Chang Liu, Rongrong Ji, and Qixiang Ye. Freeanchor: Learning to match anchors for visual object detection. *Advances in neural information processing systems*, 32, 2019.

[37] Minghang Zheng, Peng Gao, Renrui Zhang, Kunchang Li, Xiaogang Wang, Hongsheng Li, and Hao Dong. End-to-end object detection with adaptive clustering transformer. *arXiv preprint arXiv:2011.09315*, 2020.

[38] Xingyi Zhou, Jiacheng Zhuo, and Philipp Krahenbuhl. Bottom-up object detection by grouping extreme and center points. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 850–859, 2019.

[39] Chenchen Zhu, Yihui He, and Marios Savvides. Feature selective anchor-free module for single-shot object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 840–849, 2019.