

## Introduction

### Motivation

Object detection is a fundamental task in computer vision, consisting of both classification and localization tasks. Previous works mostly perform classification and localization with shared feature extractor like Convolution Neural Network. However, the tasks of classification and localization exhibit different sensitivities with regard to the same feature, hence the “**task spatial misalignment**” issue. This issue can result in a hedge issue between the performances of localizer and classifier.

### Contribution

To address these issues, we first propose a novel **Dynamic Coefficient Loss** to simultaneously consider and balance the performances of classification and localization tasks. To well address anchor label misjudgment issue in irregular-shaped object detection, we define a new **classification-aware IoU metric** to assign anchors intelligently. Finally, we further introduce the localization factor into NMS by proposing a **Classification-Localization balanced NMS**. Extensive experiments on MS COCO and PASCAL VOC demonstrate that our proposals can improve RetinaNet by around 1.5% AP with various backbones.

## Literature Review

To date, extensive efforts have been made to address the issues of object detection with deep learning technologies, and existing solutions can be divided into two categories, anchor-based and anchor-free models. Nevertheless, both anchor-based and anchor-free models employ CNN as a feature extractor, then feed extracted features into diverse localizer and classifier to respectively solve localization and classification problems. As a consequence, for an individual detector, its localizer and classifier share a set of identical features that are extracted by a same CNN. On the other hand, previous works have already pointed out that the sensitivities of classifier and localizer to same features are different. For instance, the existence of the phenomenon so called “task spatial misalignment” can be obviously observed in Figure 1. Specifically, from this figure, we discover that localization task is more interested in the marginal details of an object, while classification task focuses more on the specific robust features of an object which are mostly non-existent in marginal areas.

## Research Problems

In conclusion, the aforementioned phenomenon impacts the performance of object detection in three ways:

- **Problem 1:** There exists an obvious hedge within the performances of localizer and classifier, i.e., the performance of localizer is excellent while the performance of classifier is rather poor, or vice versa.
- **Problem 2:** For anchor-based methods, the labels of anchors are determined by the maximum overlapping between the anchors and the Ground Truth boxes. Regarding the task of irregular-shaped object detection, the aforementioned hedge within classifier and localizer may directly lead to the misjudgment of the anchors which are with poor localization performance but contain abundant classification information.
- **Problem 3:** Non-Maximum Suppression (NMS) algorithms are usually employed to suppress a portion of bounding boxes after predicting the bounding boxes with an object detector. However, most traditional NMS algorithms determine which part of bounding boxes should be reserved only by taking the single metric of classification score in account, and the phenomenon of “task spatial misalignment” indicates that this kind of operation is irrational.

## Study Methodology

1. **Dynamic coefficient loss:** The classification gradient should be relatively large in case that the localization error is relatively small. We define a dynamic coefficient, which is the reciprocal of localization error, as follows,

$$\lambda(x_i, \hat{x}_i) = \frac{1}{\text{Sigmoid}[(c_i - \hat{c}_i)^2] + \text{Tanh}\left(\left|\frac{w_i}{h_i} - \frac{\hat{w}_i}{\hat{h}_i}\right|\right) + \varepsilon} \quad (1)$$

We then normalize the coefficient  $\lambda$  by employing the classification loss function, i.e.,

$$\lambda'(x_i, \hat{x}_i) = \lambda(x_i, \hat{x}_i) \frac{\sum_{j=1}^N L_{cls}(p_j, \hat{p}_j)}{\sum_{j=1}^N \lambda(x_j, \hat{x}_j) L_{cls}(p_j, \hat{p}_j)} \quad (2)$$

The overall loss function of object detection can be denoted as,

$$\mathcal{L} = \lambda'(x_i, \hat{x}_i) \sum_{i=1}^N L_{cls}(p_i, \hat{p}_i) + \sum_{i=1}^M L_{cls}(p_i, \hat{p}_i) + \sum_{i=1}^N L_{reg}(x_i, \hat{x}_i) \quad (3)$$

2. **IoU-classification-aware anchor assignment:** For the  $i$ -th anchor, the IoU-classification-aware score can be calculated by,

$$s_i^{iou} = \alpha * \hat{p}_i + (1 - \alpha) * \frac{area_i^{intersect}}{area_i^{min}} \quad (4)$$

The new protocol for judging the label  $\Gamma_i$  of the  $i$ -th anchor can be written as,

$$\Gamma_i = \begin{cases} 1 & \text{in case } IoU_i \geq fg\_threshold \\ & \text{or } s_i^{iou} \geq score\_threshold \\ 0 & \text{in case } IoU_i \leq bg\_threshold \\ -1 & \text{otherwise} \end{cases} \quad (5)$$

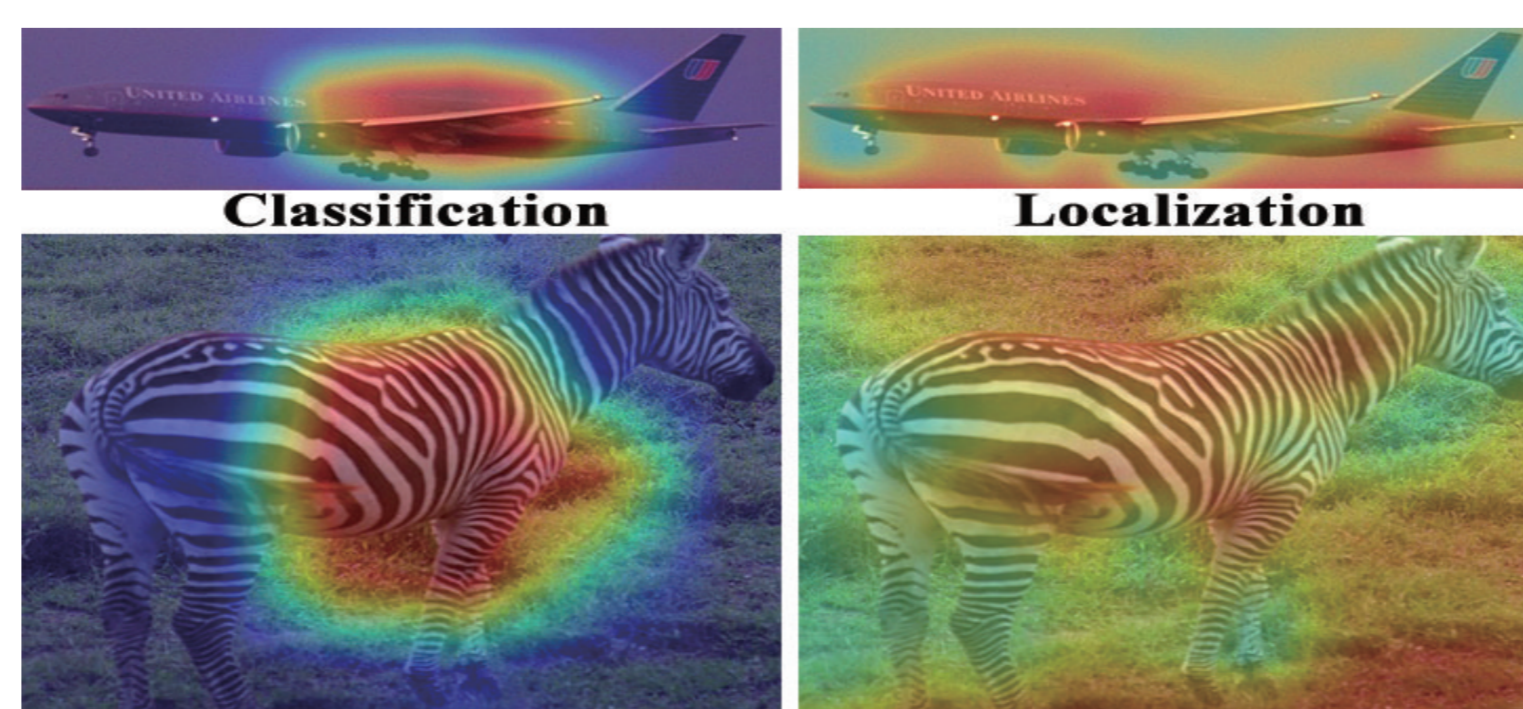
3. **Classification-localization-balanced NMS:** We calculate the weighted average center coordinates of all the boxes by using their classification scores as weights, and take the weighted average center as the approximate center  $\hat{c}^0$  of the ground truth  $a$ , i.e.,

$$\hat{c}^0 = \sum_{i=1}^n \frac{\exp(f_i^o)}{\sum_{j=1}^n \exp(f_j^o)} * c_i^o \quad (6)$$

Regarding a specific bounding box  $b_i^o (1 \leq i \leq n)$ , the new selection score of NMS algorithm can be written as,

$$s_{nms}_i^o = p_i^o + \text{dis}(c_i^o, \hat{c}^0) \quad (7)$$

Figure 1. Illustration of the phenomenon called “task spatial misalignment”. The images in the first column are the sensitive locations for classifier, and images in the second column correspond to the sensitive locations for localizer.



## Results and Discussion

We use them to modify RetinaNet and compare the modified network with other state-of-the-art detectors by conducting extensive experiments with different backbones. The results are shown in Table 1. In case of using RetinaNet with backbone ResNet-101 and ResNext-101, our methods achieve 1.5% and 1.4% AP improvements respectively, thus **verifying the superiority of our proposed methods**. Worth noting that we use DC Loss and CL-Balanced NMS to modify FCOS with Generalized Focal Loss and achieve an amazing high AP of 48.9%.

Table 1. performance of alternative detectors with different backbones on MS-COCO test-dev.

| Method           | Backbone                     | AP          | AP <sup>50</sup> | AP <sup>75</sup> | AP <sup>s</sup> | AP <sup>m</sup> | AP <sup>l</sup> |
|------------------|------------------------------|-------------|------------------|------------------|-----------------|-----------------|-----------------|
| YOLOv3           | DarkNet-53                   | 33.0        | 57.9             | 34.4             | 18.3            | 35.4            | 41.9            |
| Faster R-CNN     | ResNet-101                   | 36.2        | 59.1             | 39.0             | 18.2            | 39.0            | 48.2            |
| Mask R-CNN       | ResNet-101                   | 38.2        | 60.3             | 41.7             | 20.1            | 41.1            | 50.2            |
| Deformable R-FCN | Aligned-Inception-ResNet     | 37.5        | 58.0             | 40.8             | 19.4            | 40.1            | 52.5            |
| RetinaNet        | ResNet-101                   | 39.1        | 59.1             | 42.3             | 21.8            | 42.7            | 50.2            |
| IoU-Net          | ResNet-101                   | 40.0        | 59.0             | -                | -               | -               | -               |
| <b>Ours</b>      | <b>ResNet-101</b>            | <b>40.6</b> | <b>60.2</b>      | <b>43.5</b>      | <b>23.9</b>     | <b>43.9</b>     | <b>51.1</b>     |
| GHM              | ResNeXt-101                  | 41.6        | 62.8             | 44.2             | 22.3            | 45.1            | 55.3            |
| Faster R-CNN     | ResNeXt-101                  | 40.3        | 62.7             | 44.0             | 24.4            | 43.7            | 49.8            |
| Mask R-CNN       | ResNext-101                  | 41.4        | 63.4             | 45.2             | 24.5            | 44.9            | 51.8            |
| FCOS             | ResNeXt-101                  | 42.1        | 62.1             | 45.2             | 25.6            | 44.9            | 52.0            |
| CornerNet        | Hourglass-104                | 40.5        | 56.5             | 43.1             | 19.4            | 42.7            | 53.9            |
| RetinaNet        | ResNeXt-101                  | 40.8        | 61.1             | 44.1             | 24.1            | 44.2            | 51.2            |
| <b>Ours</b>      | <b>ResNeXt-101</b>           | <b>42.2</b> | <b>62.5</b>      | <b>45.1</b>      | <b>25.5</b>     | <b>44.9</b>     | <b>52.3</b>     |
| ATSS             | ResNeXt-101-64x4d-DCN        | 47.7        | 66.5             | 51.9             | 29.7            | 50.8            | 59.4            |
| GFL              | ResNeXt-101-32x4d-DCN        | 48.2        | 67.4             | 52.6             | 29.2            | 51.7            | 60.2            |
| <b>Ours</b>      | <b>ResNeXt-101-32x4d-DCN</b> | <b>48.9</b> | <b>67.8</b>      | <b>53.2</b>      | <b>30.0</b>     | <b>52.2</b>     | <b>61.7</b>     |

To demonstrate the **effectiveness of the three individual components** of our methods, we use these components to modify RetinaNet. The performance of these variations is evaluated with both COCO val2017 and VOC 2007 test. The results are respectively demonstrated in Table 2.

Table 2. Component analysis on COCO val and PASCAL VOC 2007 test. The baseline is ResNet-50 RetinaNet.

| IAC    | DC   | CL-B | VOC 2007 test | MS COCO 2017                         |
|--------|------|------|---------------|--------------------------------------|
| -score | Loss | NMS  | mAP           | AP AP <sup>50</sup> AP <sup>75</sup> |
|        |      |      | 78.2          | 35.9 56.1 38.6                       |
| ✓      |      |      | 78.6          | 36.4 56.5 39.1                       |
| ✓      | ✓    |      | 79.2          | 37.1 56.3 39.7                       |
| ✓      | ✓    | ✓    | 79.7          | 37.5 56.9 40.0                       |

## Conclusion

In this paper, we propose a systematical method to enhance the performance of object detection by energizing existing solutions from three different aspects to neutralize the negative impacts of task spatial misalignment. Specifically, a novel DC Loss is proposed to address the “task spatial misalignment” in object detection. Further, IoU-classification-aware score is devised to involve classification scores during assigning labels of anchors. Finally, CL-Balanced NMS is designed to address the misalignment between classification and localization via adding the localization score of candidate bounding boxes into conventional NMS algorithm.

## Acknowledgement

This work was partially supported by the National Natural Science Foundation of China (No.62072427, No.12227901), the Project of Stable Support for Youth Team in Basic Research Field, CAS (No.YSBR-005), and the Academic Leaders Cultivation Program, USTC.