# Localization Distillation for Object Detection

Zhaohui Zheng, Rongguang Ye, Qibin Hou, *Member, IEEE*, Dongwei Ren, *Member, IEEE*, Ping Wang,
Wangmeng Zuo, *Senior Member, IEEE*, and Ming-Ming Cheng, *Senior Member, IEEE*,

**Abstract**—Previous knowledge distillation (KD) methods for object detection mostly focus on feature imitation instead of mimicking the prediction logits due to its inefficiency in distilling the localization information. In this paper, we investigate whether logit mimicking always lags behind feature imitation. Towards this goal, we first present a novel localization distillation (LD) method which can efficiently transfer the localization knowledge from the teacher to the student. Second, we introduce the concept of valuable localization region that can aid to selectively distill the classification and localization knowledge for a certain region. Combining these two new components, for the first time, we show that logit mimicking can outperform feature imitation and the absence of localization distillation is a critical reason for why logit mimicking under-performs for years. The thorough studies exhibit the great potential of logit mimicking that can significantly alleviate the localization ambiguity, learn robust feature representation, and ease the training difficulty in the early stage. We also provide the theoretical connection between the proposed LD and the classification KD, that they share the equivalent optimization effect. Our distillation scheme is simple as well as effective and can be easily applied to both dense horizontal object detectors and rotated object detectors. Extensive experiments on the MS COCO, PASCAL VOC, and DOTA benchmarks demonstrate that our method can achieve considerable AP improvement without any sacrifice on the inference speed. Our source code and pretrained models are publicly available at https://github.com/HikariTJU/LD.

**Index Terms**—Object detection, localization distillation, knowledge distillation, rotated object detection.

✦

## 1 INTRODUCTION

As a model compression technology, knowledge distillation (KD) [1], [2] has been an efficient technique in learning compact models to mitigate the computational burden. It has been widely validated to be useful for boosting the performance of small-sized student networks by transferring the generalized knowledge captured by large-sized teacher networks [1], [2], [3], [4], [5], [6]. Speaking of KD in object detection, there are mainly three popular KD pipelines as shown in Fig. 1. Logit mimicking [1], also known as classification KD, is originally designed for image classification, where the KD process operates on the logits of the teacher-student pair. Feature imitation, motivated by the pioneer work FitNet [2], aims to enforce the consistency of the feature representations between the teacher-student pair. The last one, namely the pseudo bounding box regression, uses the predicted bounding boxes from the teacher as an addition supervision to the bounding box prediction branch of the student.

Among these methods, the original logit mimicking technique [1] for classification is often inefficient as it only transfers the classification knowledge while neglects the importance of localization knowledge distillation. Therefore, existent KD methods for object detection mostly focus on feature imitation, and demonstrate that distilling the feature representations is more advantageous than distilling the

- Z. Zheng, Q. Hou and M.M. Cheng are with TMCC, CS, Nankai University, Tianjin, China. E-mail: Zh_zheng@mail.nankai.edu.cn; {houqb,cmm}@nankai.edu.cn
- R. Ye and P. Wang are with the School of Mathematics, Tianjin University, China. E-mail: {ementon,wang_ping}@tju.edu.cn
- D. Ren and W. Zuo are with the School of Computer Science and Technology, Harbin Institute of Technology, China. E-mail: {rendongweihit,cswmzuo}@gmail.com
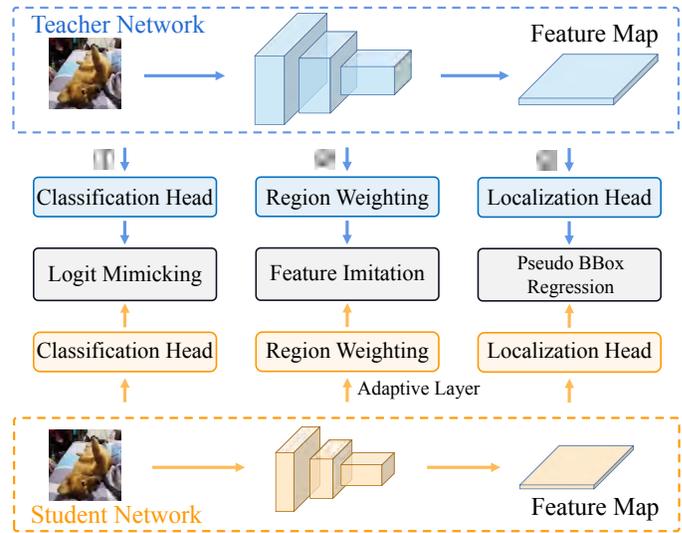- Q. Hou is the corresponding author.

Fig. 1. Existing KD pipelines for object detection. ① Logit Mimicking: classification KD in [1]. ② Feature Imitation: recent popular methods distill intermediate features based on various distillation regions, which usually need adaptive layers to align the size of the student's feature map. ③ Pseudo BBox Regression: treating teachers' predicted bounding boxes as additional regression targets [7], [8].

logits [9], [10], [11]. We summarize three crucial reasons for this phenomenon: First of all, the effectiveness of logit mimicking partially relies on the number of classes which may vary in different application scenarios [9]. Second, the logit mimicking can only be applied to the classification head, which cannot distill the localization information. At last, in the framework of multi-task learning, feature imitation can transfer the hybrid knowledge of classification and localization which can benefit the downstream classification and localization tasks.
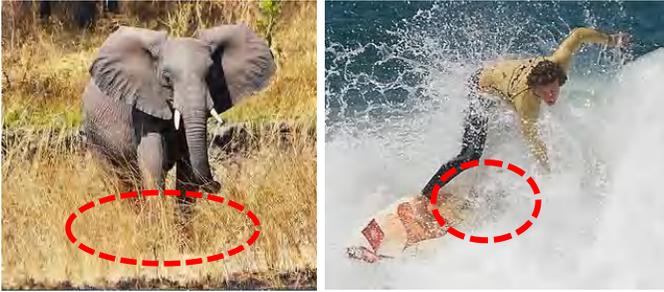
Fig. 2. Bottom edge for "*elephant*" and right edge for "*surfboard*" are ambiguous to locate.

In this work, we examine the aforementioned common belief in object detection KD, and challenge whether feature imitation always stays ahead of logit mimicking? For this purpose, we firstly present a simple yet effective localization distillation (LD) method which is inspired by an interesting observation that the bounding box distributions generated by the teacher [12], [13] can serve as a strong supervision to the student detector. The bounding box distribution [12], [13] is originally designed to model the real distributions of bounding boxes, an efficient way to solve the localization ambiguity as shown in Fig. 2. With the discretized probability distribution representations, the localizer can reflect the localization ambiguity by the flatness and sharpness of the distribution, which is not held in the conventional Dirac delta representation of bounding boxes [14], [15], [16], [17]. This allows our LD to efficiently transfer richer localization knowledge from the teacher to the student than using pseudo bounding box regression (right part in Fig. 1).

Combining the proposed LD and the classification KD yields a unified KD method based on a pure logit mimicking framework for both the classification branch and the localization branch. As logit mimicking enables us to separately distill the classification knowledge and the localization one, we found that these two sub-tasks favor different distillation regions. Motivated by this, we introduce the concept of valuable localization region (VLR) and propose to conduct distillation in a selective region distillation manner. We will show the advantage of using VLR in our distillation framework in the experiment section.

Furthermore, we comprehensively discuss the technical details of LD and elaborate on the behavior of logit mimicking and feature imitation. Intriguingly, we observe that logit mimicking can outperform feature imitation for the first time, which indicates that the absence of localization distillation is actually the key reason why logit mimicking under-performs in object detection for years. Another observation is that we find the reason why logit mimicking works is not because of the consistency of the feature representations between the teacher-student pair. Conversely, the student learns significantly different feature representations from the teacher's in terms of the $l_n$ distance and linear correlation. We also observe that if the student is trained with feature imitation, it tends to produce a sharp AP score landscape in the feature subspace, and aggravates the training difficulty in the early training stage.

The above observations reflect the great potentials of logit mimicking over feature imitation: 1) being able to separately transfer different types of knowledge, 2) learning more robust feature representations, and 3) easing the training difficulty. Our method is simple and can be easily equipped with in both horizontal and rotated object detectors to improve their performance without introducing any inference overhead. Extensive experiments on MS COCO show that without bells and whistles, we can lift the AP score of the strong baseline GFocal [12] with ResNet-50-FPN backbone from 40.1 to 42.1, and $AP_{75}$ from 43.1 to 45.6. Our best model using ResNeXt-101-32x4d-DCN backbone can achieve a single-scale test of 50.5 AP, which surpasses all existing detectors under the same backbone, neck, and test settings. PyTorch [18] and Jittor [19] version of the source code and pretrained models are publicly available at https://github.com/HikariTJU/LD.

The main contributions of this paper are four-fold:

1) We present a novel localization distillation method that greatly improve the distillation efficiency of logit mimicking in object detection.
2) We provide exploratory experiments and analysis for the behavior of logit mimicking and feature imitation. To our best knowledge, this is the first work revealing the great potential of logit mimicking over feature imitation.
3) We present a selective region distillation based on the newly introduced valuable localization region to better distill the student detector.
4) We extend our LD to a rotated version so that it can be applied to arbitrary-oriented object detection.

This paper is a substantial extension of its previous conference version [20]. In particular, (a) We provide theoretical connection for the proposed LD and the classification KD that they share the equivalent optimization effects; (b) We conduct more detailed and insightful analysis for logit mimicking and feature imitation, including the different characteristics of the learned feature representations and logits, and the training difficulty of feature imitation; (c) We extend the original LD to a more generic version, namely rotated LD, which can distill arbitrary-oriented object detectors.

## 2 RELATED WORK

### 2.1 Knowledge Distillation

Knowledge distillation [1], [21], [22], [23], [24], [25], as a hot research topic, has been deeply studied recently. The fundamental idea is to use a well-performed large-sized teacher network to transfer the captured knowledge to the small-sized student network. Logit mimicking, *a.k.a.* classification KD, was first introduced by Hinton *et al.* [1], where the logit outputs of the student classifier are supervised by those of the teacher classifier. Later, FitNet [2] extends the teacher-student learning framework by mimicking the intermediate-level hints from the hidden layers of the teacher model. Knowledge distillation was first applied to object detection in [7], where the hint learning, classification KD, and pseudo bounding box regression were simultaneously used for multi-class object detection. However, an object detector requires not only precise classification ability, but also strong localization ability. The absence of localization knowledge

distillation limits the performance of the conventional KD method.

To tackle the above issue, many feature imitation methods have been developed, most of which focus on where to distill and loss function weighting. Among these, Li *et al.* [26] proposed to mimic the features within the region proposal for Faster R-CNN. Wang *et al.* [9] imitated the fine-grained features on close anchor box locations. Recently, Dai *et al.* [27] introduced the General Instance Selection Module to mimic deep features within the discriminative patches between teacher-student pairs. DeFeat [28] leverages different loss weights when conducting feature imitation on the object regions and the background region. There are also various feature imitation methods from the perspective of weighted imitation loss, including Gaussian mask weighted [8], feature richness weighted [29], and prediction-guided imitation loss [30]. Unlike the aforementioned methods, our work introduces localization distillation and demonstrate that logit mimicking can outperform feature imitation for KD in object detection.

## 2.2 Object Localization

Object localization is a fundamental issue in object detection [31], [32], [33], [34], [35], [36], [37], [38], [39], [40]. Bounding box regression is the most popular way so far for localization in object detection [14], [15], [16], [41], [42], where the Dirac delta distribution representation has been used for years. R-CNN series [16], [43], [44], [45] adopt multiple regression stages to refine the detection results, while YOLO series [14], [46], [47], [48], SSD series [15], [49], [50], and FCOS series [12], [17] adopt one-stage regression. In [51], [52], [53], [54], IoU-based loss functions are proposed to improve the localization quality of bounding boxes. Recently, bounding box representation has evolved from Dirac delta distribution [14], [15], [16] to Gaussian distribution [55], [56], and further to probability distribution [12], [13]. The probability distribution of bounding boxes is more comprehensive for describing the uncertainty of bounding boxes, and has been validated to be the most advanced bounding box representation so far.

## 2.3 Localization Quality Estimation

As the name suggests, Localization Quality Estimation (LQE) predicts a score that measures the localization quality of the bounding boxes predicted by the detector. LQE is usually used to cooperate with the classification task during training [57], *i.e.,* enhancing the consistency between classification and localization. It can also be applied in joint decision-making during post-processing [14], [17], [58], which considers both the classification score and LQE when performing NMS. Early research can be dated to YOLOv1 [14], where the predicted object confidence is used to penalize the classification score. Then, box/mask IoU [58], [59] and box/polar centerness [17], [60] are proposed to model the uncertainty of detections in object detection and instance segmentation, respectively. For bounding box representation, Softer-NMS [55] and Gaussian YOLOv3 [56] predict variances for each edge of the bounding boxes. LQE is a preliminary approach to model localization ambiguity.

## 2.4 Arbitrary-Oriented Object Detection

Driven by the success of object detection, rotated object detection has become a hot topic in computer vision recently [61]. The mainstream rotated object detectors, such as RRPN [62], generate rotated proposals based on Faster R-CNN [16], while Rotated-RetinaNet [63] directly predicts an additional rotated angle based on RetinaNet. To address the boundary discontinuity and square-like problems, SCRDet [37] and RSDet [64] propose IoU-smooth L1 loss and modulated loss respectively for attaining smoother boundary loss, and CSL [65] proposes to use angle classification instead of angle regression.

Different from the horizontal bounding box regression which can easily leverage the IoU-based losses (*e.g.,* GIoU [52], DIoU [53], and CIoU [54]) to enhance localization ability, the Skew IoU loss for rotated bounding box regression is quite difficult to implement due to the complexity of the backward propagation in existing deep learning libraries [18], [19], [66]. PIoU loss [67] approximates the Skew IoU by accumulating the pixels of the intersection and union of two rotated bounding boxes. GWD [38] and KLD [39] model the rotated bounding boxes via the 2D Gaussian Distribution representation and propose to use the Gaussian Wasserstein distance and KL divergence to simulate the Skew IoU loss, respectively. More recently, based on the 2D Gaussian distribution representation of rotated bounding boxes, Yang *et al.* [68] proposed the KFIoU loss by exploiting the Kalman filter formulation to mimic the Skew IoU in the trend level. To sum up, the rotated regression-based detectors are still dominating this task owing to their simplicity and strong performance.

## 3 APPROACH

To begin with, we revisit the knowledge distillation background, including logit mimicking and feature imitation. Next, we describe our simple yet effective localization distillation (LD) and explain how to apply LD for rotated object detection. Then, we analyze the property of the proposed LD loss, especially the theoretical connection to the classification KD. In addition, we also introduces the concept of valuable localization region for better distilling the localization knowledge in our framework. Finally, we describe the selective region distillation based on the newly introduced valuable localization region and give the optimization objective.

### 3.1 Preliminaries

In the KD pipeline of object detection, the input image is fed into two object detectors, *i.e.,* the student detector and the frozen teacher detector. The distillation process forces the outputs of the student to mimic those of the teacher. There are two mainstream paradigms of KD methods in object detection.

**Logit mimicking.** The logit mimicking (LM) is first developed for image classification [1], in which the student model can be improved by mimicking the soft output of the teacher classifier. Let $\boldsymbol{z}_S, \boldsymbol{z}_T \in \mathbb{R}^{W \times H \times C}$ be the logits predicted by the student and the teacher, respectively. $W$ and $H$ represent the output size of the logit maps. $C$ denotes the number of
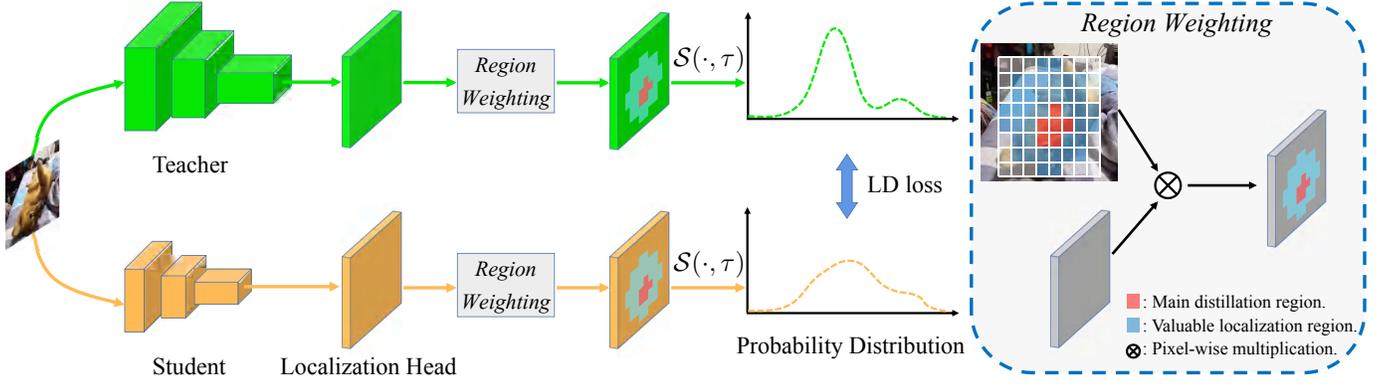
Fig. 3. Illustration of localization distillation (LD) for an edge $e \in \mathcal{B}$. Only the localization branch is visualized here. $\mathcal{S}(\cdot, \tau)$ is the generalized SoftMax function with temperature $\tau$. For a given detector, we first switch the bounding box representation to probability distribution. Then, we determine where to distill via region weighting on the main distillation region and the valuable localization region. Finally, we calculate the LD loss between two probability distributions predicted by the teacher and the student.

classes. These logits are then transformed into probability distributions $\boldsymbol{p}_\tau$ and $\boldsymbol{q}_\tau$ by using the generalized SoftMax function. We can train the network by minimizing the loss:

$$\mathcal{L} = \mathcal{L}_{\text{CE}} + \lambda \mathcal{L}_{\text{KD}} \qquad (1)$$
$$= \mathcal{H}(\boldsymbol{p}, \boldsymbol{g}) + \lambda \mathcal{H}(\boldsymbol{p}_\tau, \boldsymbol{q}_\tau), \qquad (2)$$

where $\boldsymbol{p}$ is the predicted probability vectors, $\boldsymbol{g} = \{0, 1\}^n$ is the one-hot ground-truth label, $\mathcal{H}$ is the cross-entropy loss, and $\lambda$ balances the two loss terms. For object detection, the distillation can be carried out on some pre-defined distillation region $\mathcal{R}$.

**Feature imitation.** Recently, it has been found that feature imitation (FI), which aims to transfer knowledge by imitating the deep features between teacher-student pairs, works better than the classification KD [2], [9]. Mathematically, the feature imitation procedure can be formulated as:

$$\mathcal{L}_{\text{FI}} = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} ||\tilde{M}_S(r) - M_T(r)||_2, \qquad (3)$$

where $\mathcal{R}$ is the imitation region, and $|\cdot|$ is the cardinality of the region. Note that an adaptive layer is needed to transform the size of student's feature map $M_S$ to be the same as the teacher's $M_T$, so that $\tilde{M}_S, M_T \in \mathbb{R}^{W \times H \times D}$.

**Bounding box representation.** For a given bounding box $\mathcal{B}$, the conventional representations have two forms, i.e., $\{\delta_x, \delta_y, \delta_w, \delta_h\}$ (encoding the coordinate mappings of the central point, the width and the height from the anchor box to the ground-truth box) [14], [15], [16] and $\{t, b, l, r\}$ (the distances from the sampled point to the top, bottom, left, and right edges) [17]. These two forms actually follow the Dirac delta distribution that only focuses on the ground-truth locations but cannot model the ambiguity of bounding boxes as shown in Fig. 2. This is also clearly demonstrated in some previous works [12], [55].

## 3.2 Localization Distillation

In this subsection, we present localization distillation (LD), a new way to enhance the distillation efficiency for object detection. Our LD is evolved from the view of probability distribution representation of bounding boxes (anchor free [12] and anchor-based [13]), which is originally designed for

generic object detection and carries abundant localization information. The working principle of our LD can be seen in Fig. 3. The procedure is the same to both anchor-based and anchor-free detectors.

Given an object detector, we follow [12], [13] to convert the bounding box representation from a quaternary representation to a probability distribution. Let $e \in \mathcal{B}$ be one of the regression variables of bounding box, whose regression range is $[e_{\min}, e_{\max}]$. The bounding box distribution quantizes the continuous regression range into a uniform discretized variable $\boldsymbol{e} = [e_0, e_1, \cdots, e_n] \in \mathbb{R}^{n+1}$ with $n$ sub-intervals, where $e_0 = e_{\min}$ and $e_n = e_{\max}$. The localization head predicts $n + 1$ logits $\boldsymbol{z} = \{z_0, z_1, \cdots, z_n\}$, corresponding to the endpoints of the subintervals $\{e_0, e_1, \cdots, e_n\}$. Each edge of the given bounding box can be represented as a probability distribution by using the SoftMax function. For the number of the subinterval $n$, we follow the settings of GFocal [12], and a recommended choice of $n$ is $8 \sim 16$. Different from [12], [13], we transform $\boldsymbol{z}_S$ and $\boldsymbol{z}_T$ into the probability distributions $\boldsymbol{p}_\tau$ and $\boldsymbol{q}_\tau$ using the generalized SoftMax function $\mathcal{S}(\cdot, \tau)$. Note that when $\tau = 1$, it is equivalent to the original SoftMax function. When $\tau \to 0$, it tends to be a Dirac delta distribution. When $\tau \to \infty$, it will be a uniform distribution. Empirically, $\tau > 1$ is set to soften the distribution, making the bounding box distribution carry more information. The localization distillation for measuring the similarity between the two probability vectors $\boldsymbol{p}_\tau, \boldsymbol{q}_\tau \in \mathbb{R}^n$ for one of the bounding box representation $\boldsymbol{e}$ is attained by:

$$\mathcal{L}_{\text{LD}}^{\boldsymbol{e}} = \mathcal{H}(\boldsymbol{p}_\tau, \boldsymbol{q}_\tau) \qquad (4)$$
$$= \mathcal{H}(\mathcal{S}(\boldsymbol{z}_S, \tau), \mathcal{S}(\boldsymbol{z}_T, \tau)). \qquad (5)$$

Then, LD for all the four edges of some bounding box $\mathcal{B}$ can be formulated as:

$$\mathcal{L}_{\text{LD}}(\mathcal{B}_S, \mathcal{B}_T) = \sum_{\boldsymbol{e} \in \mathcal{B}} \mathcal{L}_{\text{LD}}^{\boldsymbol{e}}, \qquad (6)$$

where $\mathcal{B}_S, \mathcal{B}_T$ are the predicted bounding boxes of the student and the teacher, respectively.

## 3.3 Rotated LD

Our LD can also be flexibly used to distill rotated bounding box detectors. Parametric regression is the most popular

manner in the classical dense regression-based rotated object detection [37], [38], [39], [69]. $\mathcal{B} = \{\delta_x, \delta_y, \delta_w, \delta_h, \delta_\theta\}$ is commonly used to represent a rotated bounding box, where $\delta_\theta$ denotes the encoded rotated angle. To conduct rotated localization distillation, we firstly generate the lower and upper bounds of the regression range $[e_{\min}, e_{\max}]$, where $e \in \mathcal{B}$.

Note that the rotated angle prediction $\delta_\theta$ usually has a different regression range from $\delta_x, \delta_y, \delta_w, \delta_h$. Thus, different lower and upper bounds of regression ranges are set for them. In practice, $[e_{min}, e_{max}] \subset [-5, 5]$ will be an acceptable choice. Then, we convert the rotated bounding box to rotated bounding box distributions, as Sec. 3.2 describes. Finally, the LD loss is calculated according to Eq. (6) for the rotated bounding box distributions.

### 3.4 Property of LD

We can see that our LD holds the formulation of the standard logit mimicking. The question one may ask is: Does LD also inherit the property of the classification KD, especially for the optimization process? Different from the classification task where a unique integer is treated as the ground-truth label, the ground-truth label of the localization task is a float point number $e^*$, whose value, for instance, is ranged in an interval $[e_i, e_{i+1}]$. In the following, we show an important property of LD, demonstrating that it can inherit the optimization effects held by the classification KD.

**Proposition 1.** *Let $s$ be the student's predicted probability vector, and $u_1, u_2$ are two constants with $u_1 + u_2 = 1$. Then, we have:*

1) *If $p, q$ are two classification probability vectors, LD effect on the linear combination $l = u_1 p + u_2 q$ is equal to the linear combination of KD effects on $p, q$;*
2) *If $l$ is a localization probability vector, LD effect on $l$ is equal to two KD effects on its decomposition $p$ and $q$.*

*The above two share the same expression,*

$$\partial LD_i^l = u_1 \partial KD_i^p + u_2 \partial KD_i^q, \tag{7}$$

*where $\partial KD_i^p$ denotes the derivatives of the KD loss of two probabilities $s, p$ w.r.t. a given logit $z_i$, and $\partial LD_i^p$ likewise for the LD loss.*

The proof can be found in the Appendix (A.1). Proposition 1 provides the theoretical connection between LD and the classification KD. It shows that the optimization effects of LD on a float point number localization problem is functionally equivalent to two KD effects on the integer position classification problems. Therefore, as a direct corollary of [70], LD holds the gradient rescaling to the distribution focal loss (DFL) [12] w.r.t. the relative prediction confidence at two near positions. For the details, we refer to the Appendix (A.2).

### 3.5 Valuable Localization Region

Previous works mostly force the deep features of the student to mimic those of the teacher by minimizing the $l_2$ loss. However, a straightforward question arises: Should we use the whole imitation regions without discrimination to distill the hybrid knowledge? According to our observation, the

---

**Algorithm 1** Valuable Localization Region

**Require:** A set of anchor boxes $\boldsymbol{B}^a = \{\mathcal{B}_i^a\}$ and a set of ground truth boxes $\boldsymbol{B}^{gt} = \{\mathcal{B}_j^{gt}\}$, $1 \leqslant i \leqslant I$, $1 \leqslant j \leqslant J$. Positive threshold $\alpha_{pos}$ of label assignment.
**Ensure:** $\boldsymbol{V} = \{v_{ij}\}_{I \times J}, v_{ij} \in \{0, 1\}$ encodes final location of VLR, where 1 denotes VLR and 0 indicates ignore.
1: Compute DIoU matrix $\boldsymbol{X} = \{x_{ij}\}_{I \times J}$ with $x_{ij} = DIoU(\mathcal{B}_i^a, \mathcal{B}_j^{gt})$.
2: $\alpha_{vl} = \gamma \alpha_{pos}$.
3: Select locations with $\boldsymbol{V} = \{\alpha_{vl} \leqslant \boldsymbol{X} \leqslant \alpha_{pos}\}$.
4: **return** $\boldsymbol{V}$

---

answer is no. In this subsection, we describe the valuable localization region (VLR) to further improve the distillation efficiency, which we believe will be a promising way to train better student detectors.

Specifically, the distillation region is divided into two parts, the main distillation region and the valuable localization region. The main distillation region is intuitively determined by label assignment, *i.e.,* the positive locations of the detection head. The valuable localization region can be obtained by Algorithm 1. First, we calculate the DIoU [53] matrix $\boldsymbol{X}$ between all the anchor boxes $\boldsymbol{B}^a$ and the ground-truth boxes $\boldsymbol{B}^{gt}$. Then, we set the lower bound of DIoU to be $\alpha_{vl} = \gamma \alpha_{pos}$, where $\alpha_{pos}$ is the positive IoU threshold of label assignment. The VLR can be defined as $\boldsymbol{V} = \{\alpha_{vl} \leqslant \boldsymbol{X} \leqslant \alpha_{pos}\}$. Our method has only one hyperparameter $\gamma \leqslant 1$, which controls the range of the VLRs. When $\gamma = 0$, all the locations whose DIoUs between the preset anchor boxes and the GT boxes satisfy $0 \leqslant x_{ij} \leqslant \alpha_{pos}$ will be determined as VLRs. When $\gamma \to 1$, the VLR will gradually shrink to empty. Here we use DIoU [53] since it gives higher priority to the locations close to the center of the object.

Similar to label assignment, our method assigns attributes to each location across multi-level FPN. In this way, some of locations outside the GT boxes will also be considered. So, we can actually view the VLR as an outward extension of the main distillation region. Note that for anchor-free detectors, like FCOS, we can use the preset anchors on feature maps and do not change its regression form, so that the localization learning maintains to be the anchor-free type. While for anchor-based detectors which usually set multiple anchors per location, we unfold the anchor boxes to calculate the DIoU matrix, and then assign their attributes.

### 3.6 Selective Region Distillation

Given the above descriptions, the total loss of logit mimicking for training the student $\boldsymbol{S}$ can be represented as:

$$\begin{aligned}
\mathcal{L} = &\lambda_0 \mathcal{L}_{cls}(\mathcal{C}_S, \mathcal{C}^{gt}) + \lambda_1 \mathcal{L}_{reg}(\mathcal{B}_S, \mathcal{B}^{gt}) + \lambda_2 \mathcal{L}_{DFL}(\mathcal{B}_S, \mathcal{B}^{gt}) \\
&+ \lambda_3 \mathbb{I}_{Main} \mathcal{L}_{LD}(\mathcal{B}_S, \mathcal{B}_T) + \lambda_4 \mathbb{I}_{VL} \mathcal{L}_{LD}(\mathcal{B}_S, \mathcal{B}_T) \\
&+ \lambda_5 \mathbb{I}_{Main} \mathcal{L}_{KD}(\mathcal{C}_S, \mathcal{C}_T) + \lambda_6 \mathbb{I}_{VL} \mathcal{L}_{KD}(\mathcal{C}_S, \mathcal{C}_T),
\end{aligned} \tag{8}$$

where the first three terms are exactly the same to the classification and bounding box regression branches for any regression-based detector, *i.e.,* $\mathcal{L}_{cls}$ is the classification loss, $\mathcal{L}_{reg}$ is the bounding box regression loss and $\mathcal{L}_{DFL}$ is the distribution focal loss [12]. $\mathbb{I}_{Main}$ and $\mathbb{I}_{VL}$ are the distillation

TABLE 1
**Ablations**. We show ablation experiments for LD and VLR on MS COCO val2017.

| $\tau$ | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|
| – | 40.1 | 58.2 | 43.1 | 23.3 | 44.4 | 52.5 |
| 1 | 40.3 | 58.2 | 43.4 | 22.4 | 44.0 | 52.4 |
| 5 | 40.9 | 58.2 | 44.3 | 23.2 | **45.0** | 53.2 |
| 10 | **41.1** | **58.7** | **44.9** | **23.8** | 44.9 | **53.6** |
| 15 | 40.7 | 58.5 | 44.2 | 23.5 | 44.3 | 53.3 |
| 20 | 40.5 | 58.3 | 43.7 | **23.8** | 44.1 | 53.5 |

(a) **Temperature $\tau$ in LD**: The generalized Softmax function with large $\tau$ brings considerable gains. We set $\tau = 10$ by default. The teacher is ResNet-101 and the student is ResNet-50.

| $\varepsilon$ | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|
| – | 40.1 | 58.2 | 43.1 | 23.3 | 44.4 | 52.5 |
| 0.1 | 40.5 | 58.3 | 43.8 | 23.0 | 44.2 | 52.7 |
| 0.2 | 40.2 | 58.2 | 43.6 | 23.1 | 44.0 | 53.0 |
| 0.3 | 40.1 | 58.4 | 43.1 | 23.6 | 43.9 | 52.5 |
| 0.4 | 40.3 | 58.4 | 43.4 | 22.8 | 44.0 | 52.6 |
| LD | **41.1** | **58.7** | **44.9** | **23.8** | **44.9** | **53.6** |

(b) **LD vs. Pseudo BBox Regression [7]**: The localization knowledge can be more efficiently transferred by our LD compared to the pseudo bbox regression. The teacher is ResNet-101 and the student is ResNet-50.

| $\gamma$ | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|
| – | 40.1 | 58.2 | 43.1 | 23.3 | 44.4 | 52.5 |
| 1 | 41.1 | 58.7 | 44.9 | 23.8 | 44.9 | 53.6 |
| 0.75 | 41.2 | 58.8 | 44.9 | 23.6 | 45.4 | 53.5 |
| 0.5 | 41.7 | 59.4 | 45.3 | 24.2 | 45.6 | 54.2 |
| 0.25 | **41.8** | **59.5** | **45.4** | 24.2 | 45.8 | **54.9** |
| 0 | 41.7 | **59.5** | **45.4** | **24.5** | **45.9** | 54.0 |

(c) **Role of $\gamma$ in VLR**: Conducting LD on valuable localization region has a positive effect on performance. We set $\gamma = 0.25$ by default. The teacher is ResNet-101 and the student is ResNet-50.

masks for the main distillation region and the valuable localization region, respectively. $\mathcal{L}_{KD}$ is the KD loss [1], $\mathcal{C}_S$ as well as $\mathcal{C}_T$ denote the classification head output logits of the student and the teacher, respectively, and $\mathcal{C}^{gt}$ is the ground-truth class label.

All the distillation losses will be weighted by the same weight factors according to their types. More clearly, the weight factor of the LD loss follows that of the bbox regression term and the weight factor of the KD loss follows that of the classification term. Also, it is worth mentioning that the DFL loss term can be disabled since LD loss has sufficient guidance ability. In addition, we can enable or disable the four types of distillation losses so as to distill the student in different regions selectively.

# 4 EXPERIMENT

In this section, we conduct comprehensive ablation studies and analysis to demonstrate the superiority of the proposed LD and distillation scheme on the challenging large-scale MS COCO [71] benchmark, PASCAL VOC [72], and aerial image DOTA dataset [73].

## 4.1 Experiment Setup

**MS COCO.** The train2017 (118K images) is utilized for training and val2017 (5K images) is used for validation. We also obtain the evaluation results on MS COCO test-dev 2019 (20K images) by submitting to the COCO server. The experiments are conducted under the mmDetection [74] framework. Unless otherwise stated, we use ResNet [75] with FPN [76] as our backbone and neck networks, and the FCOS-style [17] anchor-free head for classification and localization. The training schedule for ablation experiments is set to single-scale $1\times$ mode (12 epochs). For other training and testing hyper-parameters, we follow exactly the GFocal [12] protocol, including QFL loss for classification and GIoU loss for bbox regression, etc. We use the standard COCO-style measurement, i.e., average precision (AP), for evaluation. All the baseline models are retrained by adopting the same settings so as to fairly compare them with our LD.

**PASCAL VOC.** We also provide experimental results on another popular object detection benchmark, i.e., PASCAL VOC [72]. We use the VOC 07+12 training protocol, i.e., the union of VOC 2007 trainval set and VOC 2012 trainval set

(16551 images) for training, and VOC 2007 test set (4952 images) for evaluation. The initial learning rate is 0.01 and the total training epochs are set to 4. The learning rate decreases by a factor of 10 after the 3rd epoch. For comprehensively evaluating the localization performance, the average precision (AP) along with 5 mAP across different IoU thresholds are reported, i.e., $AP_{50}$, $AP_{60}$, $AP_{70}$, $AP_{80}$ and $AP_{90}$.

**DOTA.** As for the evaluation of rotated LD, we report the detection results on the classic aerial image dataset DOTA [73]. We follow the standard mmRotate [61] training and testing protocol. The train set and validation set consist of 1403 images and 468 images, respectively, which are randomly selected in literature. These huge images are cropped into smaller subimages with shape $600 \times 600$, which is in line with the cropping protocol in official implementation. In practice, we obtain about 15,700 training and 5,300 validation patches. Unless otherwise stated, all the hyper-parameters follow the default settings of mmRotate for a fair comparison. We report results in terms of AP and 5 mAPs under different IoU thresholds, which is consistent with PASCAL VOC. Due to the memory limitation, the teachers are ResNet-34 FPN with $2\times$ training schedule (24 epochs), and the students are ResNet-18 FPN with $1\times$ training schedule (12 epochs).

## 4.2 Ablation Study

**Temperature $\tau$ in LD.** Our LD introduces a hyper-parameter, i.e., the temperature $\tau$. Tab. 1(a) reports the results of LD with various temperatures, where the teacher model is ResNet-101 with AP 44.7 and the student model is ResNet-50. Here, only the main distillation region is adopted. Compared to the first row in Tab. 1(a), different temperatures consistently lead to better results. In this paper, we simply set the temperature in LD as $\tau = 10$, which is fixed in all the other experiments.

**LD vs. Pseudo BBox Regression.** The teacher bounded regression (TBR) loss [7] is a preliminary attempt to enhance the student on the localization head, i.e., the pseudo bbox regression in Fig. 1, which is represented as:

$$\mathcal{L}_{TBR} = \lambda\mathcal{L}_{reg}(\mathcal{B}^s, \mathcal{B}^{gt}), \text{ if } \ell_2(\mathcal{B}^s, \mathcal{B}^{gt}) + \varepsilon > \ell_2(\mathcal{B}^t, \mathcal{B}^{gt}),$$

(9)

where $\mathcal{B}^s$ and $\mathcal{B}^t$ denote the predicted boxes of student and teacher respectively, $\mathcal{B}^{gt}$ denotes the ground truth boxes,

TABLE 2
Evaluation of **selective region distillation** for KD and our LD. The teacher-student pair is ResNet-101→ResNet-50 for COCO, and ResNet-101→ResNet-18 for VOC 07+12.

| LD | | KD | | MS COCO val2017 | | | VOC 07+12 | | |
|---|---|---|---|---|---|---|---|---|---|
| Main | VLR | Main | VLR | AP | $AP_{50}$ | $AP_{75}$ | AP | $AP_{50}$ | $AP_{75}$ |
| | | | | 40.1 | 58.2 | 43.1 | 51.8 | 75.8 | 56.3 |
| ✓ | | | | 41.1 | 58.7 | 44.9 | 53.0 | 75.9 | 57.6 |
| ✓ | ✓ | | | 41.8 | 59.5 | 45.4 | 53.4 | 76.3 | **58.3** |
| ✓ | ✓ | ✓ | | **42.1** | **60.3** | **45.6** | 53.1 | 76.8 | 57.6 |
| ✓ | ✓ | ✓ | ✓ | 42.0 | 60.0 | 45.4 | **53.7** | **77.3** | 58.2 |

TABLE 3
Quantitative results of LD for lightweight detectors. The teacher is ResNet-101. The results are reported on MS COCO val2017.

| Student | LD | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|---|
| ResNet-18 | | 35.8 | 53.1 | 38.2 | 18.9 | 38.9 | 47.9 |
| | ✓ | 37.5 | 54.7 | 40.4 | 20.2 | 41.2 | 49.4 |
| ResNet-34 | | 38.9 | 56.6 | 42.2 | 21.5 | 42.8 | 51.4 |
| | ✓ | 41.0 | 58.6 | 44.6 | 23.2 | 45.0 | 54.2 |
| ResNet-50 | | 40.1 | 58.2 | 43.1 | 23.3 | 44.4 | 52.5 |
| | ✓ | 42.1 | 60.3 | 45.6 | 24.5 | 46.2 | 54.8 |

TABLE 4
Quantitative results of LD on various popular dense object detectors. The teacher is ResNet-101 and the student is ResNet-50. The results are reported on MS COCO val2017.

| Student | LD | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|---|
| RetinaNet [63] | | 36.9 | 54.3 | 39.8 | 21.2 | 40.8 | 48.4 |
| | ✓ | 39.0 | 56.4 | 42.4 | 23.1 | 43.2 | 51.1 |
| FCOS [17] | | 38.6 | 57.2 | 41.5 | 22.4 | 42.2 | 49.8 |
| | ✓ | 40.6 | 58.4 | 44.1 | 24.3 | 44.1 | 52.3 |
| ATSS [77] | | 39.2 | 57.3 | 42.4 | 22.7 | 43.1 | 51.5 |
| | ✓ | 41.6 | 59.3 | 45.3 | 25.2 | 45.2 | 53.3 |

$\varepsilon$ is a predefined margin, and $\mathcal{L}_{reg}$ represents the GIoU loss [52]. Here, only the main distillation region is adopted. From Tab. 1(b), we can see that the TBR loss does yield performance gains (+0.4 AP and +0.7 $AP_{75}$) when using proper threshold $\varepsilon = 0.1$ in Eq. (9). However, it uses the coarse bbox representation, which does not contain any localization uncertainty information of the detector, leading to sub-optimal results. On the contrary, our LD directly produces 41.1 AP and 44.9 $AP_{75}$, since it utilizes the probability distribution of bounding boxes which contains rich localization knowledge.

**Various $\gamma$ in VLR.** The newly introduced VLR has the parameter $\gamma$ which controls the range of VLR. As shown in Tab. 1(c), AP is stable when $\gamma$ ranges from 0 to 0.5. The variation in AP in this range is around 0.1. As $\gamma$ increases, the VLR gradually shrinks to empty. The performance also gradually drops to 41.1, *i.e.,* conducting LD on the main distillation region only. The sensitivity analysis experiments on the parameter $\gamma$ indicate that conducting LD on the VLR has a positive effect on performance. In the rest experiments, we set $\gamma$ to 0.25 for simplicity.

**Selective Region Distillation.** There are several interesting observations regarding the roles of KD and LD and their preferred regions. We report the relevant ablation study results in Tab. 2, where "Main" means that the logit mimicking is conducted on the main distillation region, *i.e.,* the positive locations of label assignment, and "VLR" denotes the valuable localization region. For MS COCO, it can be seen that conducting "Main LD", "VLR LD", and "Main KD" all benefits the student's performance. This indicates that the main distillation regions contain the valuable knowledge for both classification and localization and the classification KD benefits less compared to LD. Then, we impose the classification KD on a larger range, *i.e.,* the VLR. However, we observe that further incorporating "VLR KD" yields no improvement (the last two rows of Tab. 2). This is the main reason why we adopt the proposed selective region distillation for COCO.

Next, we check the roles of KD and LD on PASCAL VOC. Tab. 2 shows that it is beneficial to transfer the localization knowledge to both the main distillation region and the VLR. However, due to the different knowledge distribution patterns, it shows a similar degradation of the classification KD. Comparing the 3rd row and the 4th row of Tab. 2, "Main KD" leads to a performance drop, while "VLR KD" produces a positive effect to the student. This indicates that the selective region distillation can take the advantages of both KD and LD on their respective favorable regions.

**LD for Lightweight Detectors.** Tab. 3 reports the results of our distillation scheme ("Main LD + VLR LD + Main KD" on COCO), where a series of lightweight students are distilled, including ResNet-18, ResNet-34, and ResNet-50. For all given students, our LD can stably improve the detection performance without any bells and whistles. From these results, we can see that our LD improves the students ResNet-18, ResNet-34, ResNet-50 by +1.7, +2.1, +2.0 in AP, and +2.2, +2.4, +2.5 in $AP_{75}$, respectively.

**Application to Other Dense Object Detectors.** Our LD can be flexibly applied to other dense object detectors, including either anchor-based or anchor-free types. We employ LD with the divide-and-conquer distillation scheme to several recently popular detectors, such as RetinaNet [63] (anchor-based), FCOS [17] (anchor-free) and ATSS [77] (anchor-based). According to the results in Tab. 4, we can see that our LD can consistently improve the baselines by around 2 AP scores.

**Arbitrary-Oriented Object Detectors.** As a direct extension of our LD, the rotated bounding box requires an additional probability distribution, i.e., the rotated angle distribution. We make the necessary and minimum modification to two arbitrary-oriented object detectors, 1) the foundation of dense regression-based rotated detector—Rotated-RetinaNet [63] and 2) the recently popular 2D Gaussian distribution modeling detector—GWD [38]. We follow the mmRotate [61] training and testing protocols. We use ResNet-34 as the teacher and ResNet-18 as the student for GPU memory saving. The results are reported on the validation set of DOTA-v1.0 [73].

The results have been shown in Tab. 5, which demonstrates that our LD can also be successfully applied to rotated object detectors and attain considerable improvement in aerial image detection. Particularly, we obtain impres-

TABLE 5
Quantitative results of rotated LD on the popular arbitrary-oriented object detectors. The teacher is ResNet-34 and the student is ResNet-18. The results are reported on the validation set of DOTA-v1.0.

| Student | AP | $AP_{50}$ | $AP_{60}$ | $AP_{70}$ | $AP_{80}$ | $AP_{90}$ |
|---|---|---|---|---|---|---|
| R-RetinaNet [63] | 33.7 | 58.0 | 54.5 | 42.3 | 22.9 | 4.7 |
| LD (ours) | 39.1 | 63.8 | 61.1 | 48.8 | 28.7 | 8.8 |
| GWD [38] | 37.1 | 63.1 | 60.1 | 46.7 | 24.7 | 6.2 |
| LD (ours) | 40.2 | 66.4 | 63.6 | 50.3 | 28.2 | 8.5 |

sive improvements for the mAP under more rigorous IoU thresholds, e.g., $AP_{70}$, $AP_{80}$, $AP_{90}$. This shows the excellent compatibility of our LD, which can be applied to not only horizontal bounding boxes but also the rotated ones. In addition, it is worth mentioning that our LD does not rely on the representations of bounding boxes and the optimization way of modeling (IoU-based loss for horizontal bounding box prediction [52], [53] and 2D Gaussian modeling for rotated bounding box prediction [38]).

## 4.3 Logit Mimicking *v.s.* Feature Imitation.

Thus far, we have validated the effectiveness of our LD and the selective region distillation in distilling different types of object detectors. The proposed LD along with the classification KD provides a unified logit mimicking framework. It naturally raises several interesting questions:

- In terms of detection performance, how does logit mimicking perform compared to feature imitation? Does feature imitation stay ahead of logit mimicking?
- What are the characteristics of these two different distillation techniques? Are the deep feature representations and logits learned different?

In this subsection, we shall provide answers to the above questions.

**Quantitative Comparison on Numerical Results.** We first compare our proposed LD with several state-of-the-art feature imitation methods. We adopt the selective region distillation, i.e., performing KD and LD on the main distillation region, and performing LD on the VLR. Since modern detectors are usually equipped with FPN [76], following previous works [9], [27], [28], we re-implement their methods and impose all the feature imitations on multi-level FPN for a fair comparison. Here, "FitNets" [2] distills the whole feature maps. "DeFeat" [28] means the loss weights of feature imitation outside the GT boxes are larger than those inside the GT boxes. "Fine-Grained" [9] distills the deep features on the close anchor box locations. "GI Imitation" [27] selects the distillation regions according to the discriminative predictions of the student and the teacher. "Inside GT Box" means we select the ground-truth boxes with the same stride on the FPN layers as the feature imitation regions. "Main Region" means we imitate the features within the main distillation region.

From Tab. 6, we can see that distillation within the whole feature maps attains +0.6 AP gains. By setting a larger loss weight for the locations outside the GT boxes (DeFeat [28]), the performance is slightly better than that using the same

TABLE 6
**Logit Mimicking *vs.* Feature Imitation.** "Ours" means we use the selective region distillation, i.e., "Main LD + VLR LD + Main KD". "*" denotes we remove the "Main KD". The teacher is ResNet-101 and the student is ResNet-50 [75]. The results are reported on MS COCO val2017.

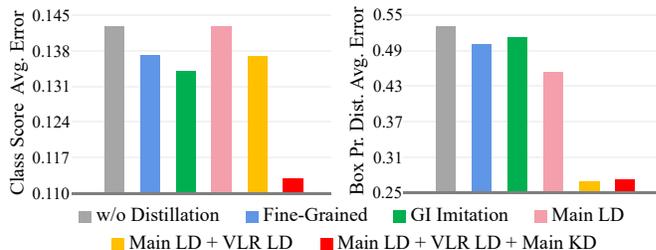| Method | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|
| Baseline (GFocal [12]) | 40.1 | 58.2 | 43.1 | 23.3 | 44.4 | 52.5 |
| FitNets [2] | 40.7 | 58.6 | 44.0 | 23.7 | 44.4 | 53.2 |
| Inside GT Box | 40.7 | 58.6 | 44.2 | 23.1 | 44.5 | 53.5 |
| Main Region | 41.1 | 58.7 | 44.4 | 24.1 | 44.6 | 53.6 |
| Fine-Grained [9] | 41.1 | 58.8 | 44.8 | 23.3 | 45.4 | 53.1 |
| DeFeat [28] | 40.8 | 58.6 | 44.2 | 24.3 | 44.6 | 53.7 |
| GI Imitation [27] | 41.5 | 59.6 | 45.2 | 24.3 | 45.7 | 53.6 |
| Ours | 42.1 | 60.3 | 45.6 | 24.5 | 46.2 | 54.8 |
| Ours + FitNets | 42.1 | 59.9 | 45.7 | 25.0 | 46.3 | 54.4 |
| Ours + Inside GT Box | 42.2 | 60.0 | 45.9 | 24.3 | 46.3 | 55.0 |
| Ours + Main Region | 42.1 | 60.0 | 45.7 | 24.6 | 46.3 | 54.7 |
| Ours + Fine-Grained | 42.4 | 60.3 | 45.9 | 24.7 | 46.5 | 55.4 |
| Ours* + Fine-Grained | 42.1 | 59.7 | 45.6 | 24.8 | 46.1 | 54.8 |
| Ours + DeFeat | 42.2 | 60.0 | 45.8 | 24.7 | 46.1 | 54.4 |
| Ours + GI Imitation | 42.4 | 60.3 | 46.2 | 25.0 | 46.6 | 54.5 |



Fig. 4. Visual comparisons of SOTA feature imitation and our LD. We show the average L1 error of classification scores and box probability distributions between teacher and student at the P4, P5, P6 and P7 FPN levels. The teacher is ResNet-101 and the student is ResNet-50. The results are evaluated on MS COCO val2017.

loss weight for all locations. Fine-Grained [9] focusing on the locations near GT boxes, produces 41.1 AP, which is comparable to the results of feature imitation using the Main Region. GI imitation [27] searches the discriminative patches for feature imitation and gains 41.5 AP. Due to the large gap in predictions between student and teacher, the imitation regions may appear anywhere.

Despite the notable improvements of these feature imitation methods, they do not explicitly consider the knowledge distribution patterns. On the contrary, our method can transfer the knowledge via a selective region distillation, which directly produces 42.1 AP. It is worth noting that our method operates on logits instead of deep features, indicating that our LD is a critical component for logit mimicking to outperform the feature imitation. Moreover, our method is orthogonal to the aforementioned feature imitation methods. Tab. 6 shows that with these feature imitation methods, our performance can be further improved. Particularly, with GI imitation, we improve the strong GFocal baseline by +2.3 AP and +3.1 $AP_{75}$.

**Teacher-Student Error Comparison.** We first check the average teacher-student errors of the classification scores
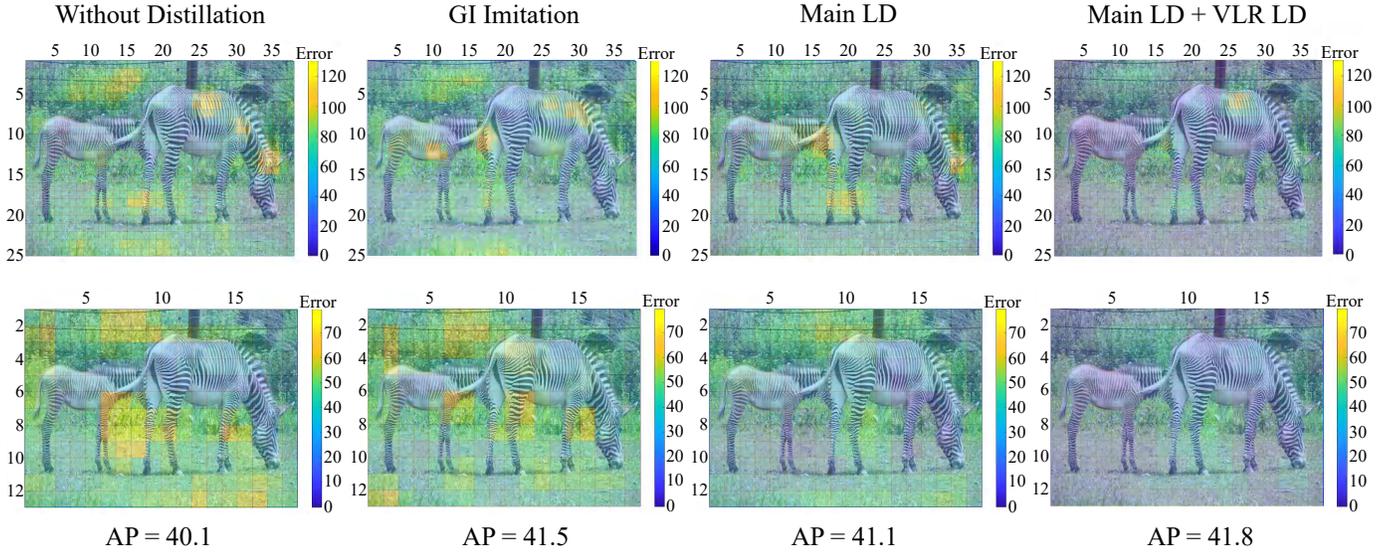
Fig. 5. Visual comparisons between the state-of-the-art feature imitation and our LD. We show the per-location L1 error summation of the localization head logits between the teacher and the student as the P5 (first row) and P6 (second row) FPN levels. The teacher is ResNet-101 and the student is ResNet-50. We can see that compared to the GI imitation [27], our method ("Main LD + VLR LD") can significantly reduce the errors for almost all the locations. Darker is better. Best viewed in color.

TABLE 7
The average Pearson correlation coefficient between the teacher-student pair. 'GI': GI imitation. 'Ours': our logit mimicking scheme with the selective region distillation. The results are evaluated on MS COCO val2017.

|  | w/o distillation | GI | Ours | Ours + GI |
|---|---|---|---|---|
| deep features | -0.0042 | 0.8175 | -0.0031 | 0.8373 |
| bbox logits | 0.9222 | 0.9326 | 0.9733 | 0.9745 |

and the box probability distributions, as shown in Fig 4. One can see that the Fine-Grained feature imitation [9] and GI imitation [27] reduce the two errors as expected, since the classification knowledge and localization knowledge are mixed on feature maps. Our "Main LD" and "Main LD + VLR LD" have comparable or larger classification score average errors than Fine-Grained [9] and GI imitation [27] but lower box probability distribution average errors. This indicates that these two settings with only LD can significantly reduce the box probability distribution distance between the teacher and the student but they cannot reduce this error for the classification head. If we impose the classification KD on the main distillation region, yielding "Main LD + VLR LD + Main KD", both the classification score average error and the box probability distribution average error can be reduced.

We also visualize the L1 error summation of the localization head logits between the student and the teacher for each location at the P5 and P6 FPN levels. As shown in Fig. 5, comparing to "Without Distillation", we can see that the GI imitation [27] does decrease the localization discrepancy between the teacher and the student. Notice that we particularly choose a model ("Main LD + VLR LD") with slightly better AP performance than GI imitation for visualization. Our method can clearly reduce this error and alleviate the localization ambiguity.

In Fig. 6, we plot the average errors between the student

and the teacher in terms of deep feature, class logit and bbox logit, respectively. It can be seen that these three types of errors show an almost consistent trend as the test resolution changes. Interestingly, we find that even though the logit mimicking can shrink the errors of both the bbox logits and the classification ones, it learns complete different feature representations from the teacher's. From the left side of Fig. 6, our method enlarges the distance between the student's feature representations and those of the teacher. Moreover, Tab. 7 shows that the logit mimicking produces a nearly zero Pearson correlation coefficient for the feature representations between the teacher-student pair. This indicates that if the student is only trained with logit mimicking, it produces a far different and nonlinearly correlated feature representation to teacher's. Be that as it may, we can still attain well-performed logits for good generalization. The last column of Tab. 7 and Fig. 6 show that the logit mimicking is able to approach the teacher's logits not only in distance but also in linear correlation.

**AP Landscape.** Distilling an object detector from either the feature level or the logit level is a high-dimensional non-convex optimization problem, which is easy in practice but hard in theory. To better understand the behavior of logit mimicking and feature imitation, we present a new visualization method, termed AP landscape, which is especially designed for object detection to observe the AP changes caused by minute perturbations in the learnt feature representations. A canonical approach was taken in [78], who studied the loss surface visualization by linearly interpolating the parameters of two networks.

In our visualization, we are particularly curious about the empirical characterization of the feature representations and how they affect the final performance. Considering two feature representations $M_f$, $M_l$ which are learnt by the detectors trained with feature imitation and logit mimicking, respectively, we visualize the AP landscapes within the 2D projected space $M_f \oplus M_l$. We use two scalar parameters
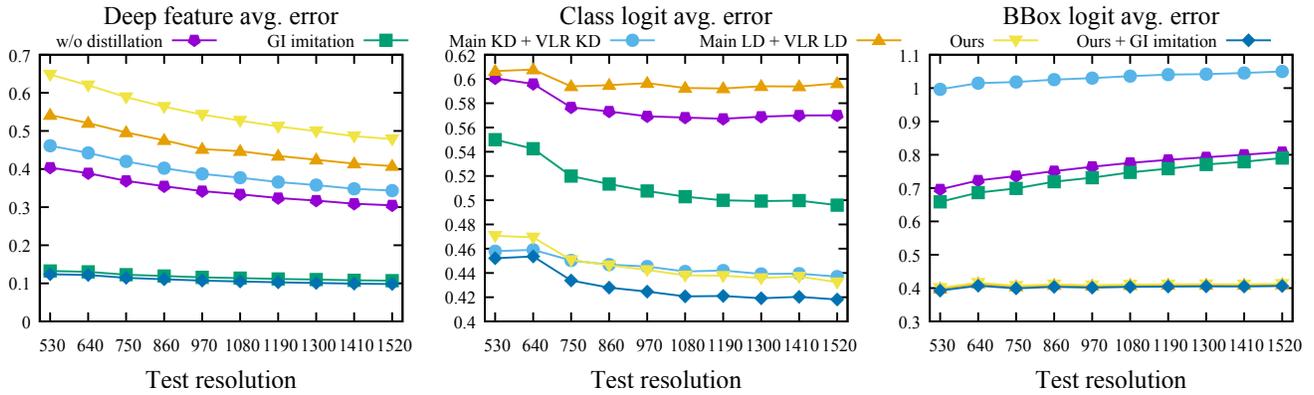
Fig. 6. Average teacher-student error on (left) deep feature representation, (middle) class logits, and (right) bbox logits. "Ours" denotes "Main LD + VLR LD + Main KD". The curves are evaluated on MS COCO val2017.
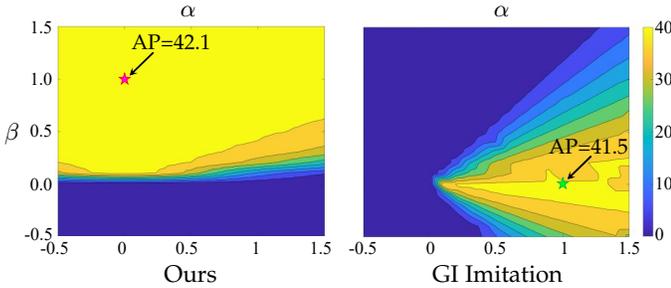


Fig. 7. The 2D contour plots of AP landscapes in feature subspace. The AP landscapes are evaluated on MS COCO val2017.
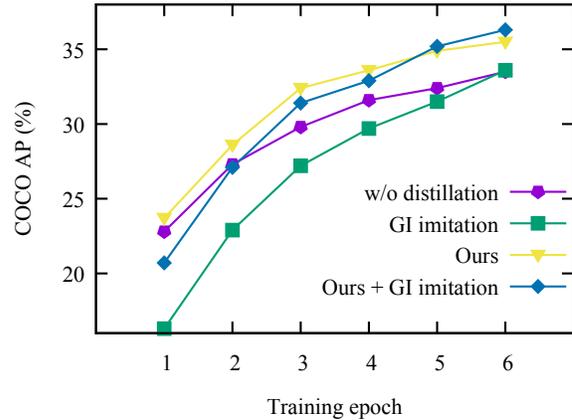


Fig. 8. The average precision (AP) during the early training stage. The feature imitation significantly slows down the convergence and gets a sub-optimal generalization. Logit mimicking (Ours) can reduce the training difficulty in the early training stage.

$\alpha$ and $\beta$ to obtain a new feature representation by using the weighted sum $M(\alpha, \beta) = \alpha M_f + \beta M_l$. Note that when $\alpha = 0$ and $\beta = 1$, it represents that the feature representations are predicted by the logit mimicking method and inversely the feature imitation when $\alpha = 1$ and $\beta = 0$. Then, we feed $M(\alpha, \beta)$ to the downstream heads and plot the final AP score. Due to the computational burden, we set $\alpha, \beta \in [-0.5, 1.5]$ to visualize the 2D AP landscapes.

From Fig. 7, we see that logit mimicking learns robust feature representations, i.e., the red pentagram at $(0, 1)$, which is surrounded by a flat and well-performed region of AP score. Second, we observe that the GI imitation produces a much sharper AP landscape than logit mimicking. We attribute the landscape sharpness of the GI imitation to the hard $l_2$ loss supervision. In this case, it is hard for the student to imitate the high-level and advanced feature representations from the teacher, which corresponds to a heavy detector with a longer training schedule and higher accuracy. On the contrary, the logit mimicking gives the feature representations much more liberty to learn, leading to a better generalization. As shown in Fig. 8, logit mimicking can also reduce the optimization difficulty in the early training stage, while feature imitation converges slower and has a poor generalization in the early training stage.

**Summary.** Based on the above results and observations, we can draw the following conclusions:

- Logit mimicking can outperform feature imitation in object detection when the localization knowledge distillation is explicitly distilled.
- Feature imitation can increase the consistency of the feature representations between the teacher-student

pair, but come some drawbacks such as less feature robustness and slow training convergence. Logit mimicking with the selective region distillation can significantly increase the consistency of the logits between the teacher-student pair, keep the learning liberty of features, and thereby speed up training process and benefit the KD performance more. This indicates that the consistency of feature representations between the teacher-student pair is not the crucial factor of improving the KD performance.

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

We compare our LD with the state-of-the-art dense object detectors by using our LD to further boost GFocalV2 [57]. For COCO val2017, since most previous works use ResNet-50-FPN backbone with the single-scale $1\times$ training schedule (12 epochs) for validation, we also report the results under this setting for a fair comparison. For COCO test-dev 2019, following a previous work [57], the LD models with the $1333 \times [480 : 960]$ multi-scale $2\times$ training schedule (24 epochs) are included. The training is carried on a machine node with 8 GPUs with a batch size of 2 per GPU and initial learning rate 0.01 for a fair comparison. During inference, single-scale testing ($[1333 \times 800]$ resolution) is adopted. For

TABLE 8
Comparison with state-of-the-art methods on COCO *val2017* and *test-dev2019* . **TS**: Traning Schedule. '1×': single-scale training 12 epochs. '2×': multi-scale training 24 epochs.

| Method | TS | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|---|
| **ResNet-50 backbone on val2017** | | | | | | | |
| RetinaNet [63] | 1× | 36.9 | 54.3 | 39.8 | 21.2 | 40.8 | 48.4 |
| FCOS [17] | 1× | 38.6 | 57.2 | 41.5 | 22.4 | 42.2 | 49.8 |
| SAPD [82] | 1× | 38.8 | 58.7 | 41.3 | 22.5 | 42.6 | 50.8 |
| ATSS [77] | 1× | 39.2 | 57.3 | 42.4 | 22.7 | 43.1 | 51.5 |
| BorderDet [83] | 1× | 41.4 | 59.4 | 44.5 | 23.6 | 45.1 | 54.6 |
| AutoAssign [84] | 1× | 40.5 | 59.8 | 43.9 | 23.1 | 44.7 | 52.9 |
| PAA [85] | 1× | 40.4 | 58.4 | 43.9 | 22.9 | 44.3 | 54.0 |
| OTA [86] | 1× | 40.7 | 58.4 | 44.3 | 23.2 | 45.0 | 53.6 |
| GFocal [12] | 1× | 40.1 | 58.2 | 43.1 | 23.3 | 44.4 | 52.5 |
| GFocalV2 [57] | 1× | 41.1 | 58.8 | 44.9 | 23.5 | 44.9 | 53.3 |
| LD (ours) | 1× | **42.7** | **60.2** | **46.7** | **25.0** | **46.4** | **55.1** |
| **ResNet-101 backbone on test-dev 2019** | | | | | | | |
| RetinaNet [63] | 2× | 39.1 | 59.1 | 42.3 | 21.8 | 42.7 | 50.2 |
| FCOS [17] | 2× | 41.5 | 60.7 | 45.0 | 24.4 | 44.8 | 51.6 |
| SAPD [82] | 2× | 43.5 | 63.6 | 46.5 | 24.9 | 46.8 | 54.6 |
| ATSS [77] | 2× | 43.6 | 62.1 | 47.4 | 26.1 | 47.0 | 53.6 |
| BorderDet [83] | 2× | 45.4 | 64.1 | 48.8 | 26.7 | 48.3 | 56.5 |
| AutoAssign [84] | 2× | 44.5 | 64.3 | 48.4 | 25.9 | 47.4 | 55.0 |
| PAA [85] | 2× | 44.8 | 63.3 | 48.7 | 26.5 | 48.8 | 56.3 |
| OTA [86] | 2× | 45.3 | 63.5 | 49.3 | 26.9 | 48.8 | 56.1 |
| GFocal [12] | 2× | 45.0 | 63.7 | 48.9 | 27.2 | 48.8 | 54.5 |
| GFocalV2 [57] | 2× | 46.0 | 64.1 | 50.2 | 27.6 | 49.6 | 56.5 |
| LD (ours) | 2× | **47.1** | **65.0** | **51.4** | **28.3** | **50.9** | **58.5** |
| **ResNeXt-101-32x4d-DCN backbone on test-dev 2019** | | | | | | | |
| SAPD [82] | 2× | 46.6 | 66.6 | 50.0 | 27.3 | 49.7 | 60.7 |
| GFocal [12] | 2× | 48.2 | 67.4 | 52.6 | 29.2 | 51.7 | 60.2 |
| GFocalV2 [57] | 2× | 49.0 | 67.6 | 53.4 | 29.8 | 52.3 | 61.8 |
| LD (ours) | 2× | **50.5** | **69.0** | **55.3** | **30.9** | **54.4** | **63.4** |

different students ResNet-50, ResNet-101 and ResNeXt-101-32x4d-DCN [79], [80], we also choose different networks ResNet-101, ResNet-101-DCN and Res2Net-101-DCN [81] as their teachers, respectively.

Tab. 8 reports the quantitative results. It can be seen that our LD improves the AP score of the SOTA GFocalV2 by +1.6 and the $AP_{75}$ score by +1.8 when using the ResNet-50-FPN backbone. When using the ResNet-101-FPN and ResNeXt-101-32x4d-DCN with multi-scale 2× training, we achieve the highest AP scores, 47.1 and 50.5 , which outperform all existing dense object detectors under the same backbone, neck and test settings. More importantly, our LD does not introduce any additional network parameters or computational overhead and hence can guarantee exactly the same inference speed as GFocalV2.

## 5 CONCLUSION

In this paper, we propose a flexible localization distillation for dense object detection and a selective region distillation based on a new valuable localization region. We show that 1) logit mimicking can be better than feature imitation; and 2) the selective region distillation for transferring the classification and localization knowledge is important when distilling object detectors. We hope our method could provide

new research intuitions for the object detection community to develop better distillation strategies. In the future, the applications of LD to sparse object detectors (DETR [87] series), the heterogeneous detector pairs, and other relevant fields, e.g., instance segmentation, object tracking and 3D object detection, warrant further research. Besides, since our LD shares the equivalent optimization effect to classification KD, some improved KD methods may also bring gain to LD, e.g., Relational KD [23], Self-KD [88], [89], Teacher Assistant KD [24], and Decoupled KD [90], etc. Cross architecture distillation using recent state-of-the-art classification models [91], [92], [93], [94], [95] as teachers is also an interesting direction to explore.

## APPENDIX

### A.1  Property of LD

Some Notations.

$\boldsymbol{g}^i = [g_1, g_2, \cdots, g_n]$, where $g_i = 1$, and 0 otherwise.

$\boldsymbol{e} = [e_1, e_2, \cdots, e_n] \in \mathbb{R}^n$ is the uniformly discretized variable for the regression range $[e_{\min}, e_{\max}]$.

The gradient of the cross-entropy (CE) loss *w.r.t.* one of the logit $z_i \in \boldsymbol{z}_S, i \in \{1, 2, \cdots, n\}$ can be represented as:

$$\frac{\partial \mathcal{L}_{\text{CE}}}{\partial z_i} = p_i - g_i, \tag{10}$$

where $p_i$ is the predicted class probability at location $i$ and $z_S$ is the logit vector produced by the student network. The gradient of the KD loss along with the CE loss *w.r.t.* one of the logit $z_i \in \boldsymbol{z}_S$ can be represented as:

$$\frac{\partial \mathcal{L}^{KD}}{\partial z_i} = \gamma(p_i - g_i) + \frac{\lambda}{\tau}(p_{\tau i} - q_{\tau i}), \tag{11}$$

where $\gamma$ and $\lambda$ are the CE and KD loss weights and $\tau$ is the temperature. We follow the notations in [70], and denote $\frac{\partial \mathcal{L}_{\text{CE}}}{\partial z_i}$ by $\partial_i$ and $\frac{\partial \mathcal{L}^{KD}}{\partial z_i}$ by $\partial_i^{KD}$. The ratio of Eq. 11 and Eq. 10 indicates that KD performs gradient rescaling to the CE loss in the logits space.

**Definition 1.** *Let $\boldsymbol{p} \in \mathbb{R}^n$ be a predicted probability vector of a network, $M_i > 0, i \in \{1, 2, \cdots, n\}$ are predefined thresholds. $\boldsymbol{p}$ is called $M_i$-well-performed for a task $\mathbf{T}$ if the distance from $\boldsymbol{p}$ to its corresponding ground-truth vector $\boldsymbol{g}^i$ is bounded by $M_i$.*

**Lemma 1.** *If two predicted probability vectors $\boldsymbol{p}, \boldsymbol{q}$ are respectively $M_i$-well-performed and $M_j$-well-performed for the integer position classification with ground-truth vectors $\boldsymbol{g}^i, \boldsymbol{g}^j$, then their linear combination $u_1\boldsymbol{p} + u_2\boldsymbol{q}$ is $M$-well-performed for the float point number position localization with ground-truth value $y = u_1 e_i + u_2 e_j$, where $M = \max\{M_i, M_j\}$ and $u_1 + u_2 = 1$.*

*Proof.* By Def. 1, the two distances satisfy $d(\boldsymbol{p}, \boldsymbol{g}^i) \leqslant M_i$, $d(\boldsymbol{q}, \boldsymbol{g}^j) \leqslant M_j$, where $\boldsymbol{g}^i \neq \boldsymbol{g}^j$. Note that $d(\cdot, \cdot)$ here can be an arbitrary distance metric, e.g., the $l_2$ distance.

A float point number position localization requires a probability, which can be linearly interpolated by $\boldsymbol{l} = u_1\boldsymbol{p} + u_2\boldsymbol{q}$, and its ground-truth vector is $\boldsymbol{g} = u_1\boldsymbol{g}^i + u_2\boldsymbol{g}^j$.

Then we get

$$d(\boldsymbol{l}, \boldsymbol{g}) = d(u_1\boldsymbol{p} + u_2\boldsymbol{q}, u_1\boldsymbol{g}^i + u_2\boldsymbol{g}^j) \tag{12}$$

$$\leqslant d(u_1\boldsymbol{p} + u_2\boldsymbol{q}, u_1\boldsymbol{g}^i + u_2\boldsymbol{q})$$
$$+ d(u_1\boldsymbol{g}^i + u_2\boldsymbol{q}, u_1\boldsymbol{g}^i + u_2\boldsymbol{g}^j) \tag{13}$$

$$= d(u_1\boldsymbol{p}, u_1\boldsymbol{g}^i) + d(u_2\boldsymbol{q}, u_2\boldsymbol{g}^j) \tag{14}$$

$$\leqslant u_1 M_i + u_2 M_j \tag{15}$$

$$\leqslant \max\{M_i, M_j\} \tag{16}$$

$$= M. \tag{17}$$

Hence the network is $M$-well-performed for the float point number position localization. □

**Lemma 2.** *If $\boldsymbol{l}$ is a localization probability vector with ground-truth value $y = u_1 e_i + u_2 e_j$, where $u_1 + u_2 = 1$, then $\boldsymbol{l}$ can be decomposed into two classification probabilities $\boldsymbol{p}$ and $\boldsymbol{q}$ with ground-truth vectors $\boldsymbol{g}^i$ and $\boldsymbol{g}^j$.*

*Proof.* Let $\boldsymbol{l} \in \mathbb{R}^n$ be a predicted localization probability and $\boldsymbol{g}$ be its ground-truth vector. It is easy to decompose $\boldsymbol{g}$ into two integer position ground-truth vectors $\boldsymbol{g}^i$ and $\boldsymbol{g}^j$, satisfying $\boldsymbol{g} = u_1\boldsymbol{g}^i + u_2\boldsymbol{g}^j$.

Existence of the decomposition of $\boldsymbol{l}$:

To decompose $\boldsymbol{l}$ into two classification probabilities $\boldsymbol{p}$ and $\boldsymbol{q}$, satisfying $\boldsymbol{l} = u_1\boldsymbol{p} + u_2\boldsymbol{q}$, we solve the following linear equations,

$$AX = \boldsymbol{b} \Longleftrightarrow \begin{cases} \sum_i p_i = 1, \\ \sum_i q_i = 1, \\ u_1 p_1 + u_2 q_1 = l_1, \\ u_1 p_2 + u_2 q_2 = l_2, \\ \vdots \\ u_1 p_n + u_2 q_n = l_n, \end{cases} \tag{18}$$

where $X = (p_1, p_2, \cdots, p_n, q_1, q_2, \cdots, q_n)^T$, and the augmented matrix $(A, \boldsymbol{b})$ is given by:

$$\begin{pmatrix} 1 & 1 & 1 & \cdots & 1 & 0 & 0 & 0 & \cdots & 0 & 1 \\ 0 & 0 & 0 & \cdots & 0 & 1 & 1 & 1 & \cdots & 1 & 1 \\ u_1 & 0 & 0 & \cdots & 0 & u_2 & 0 & 0 & \cdots & 0 & l_1 \\ 0 & u_1 & 0 & \cdots & 0 & 0 & u_2 & 0 & \cdots & 0 & l_2 \\ 0 & 0 & u_1 & \cdots & 0 & 0 & 0 & u_2 & \cdots & 0 & l_3 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & u_1 & 0 & 0 & 0 & \cdots & u_2 & l_n \end{pmatrix}. \tag{19}$$

By applying Elementary Row Operations, the matrix $(A, \boldsymbol{b})$ is equivalent to

$$\begin{pmatrix} 1 & 0 & 0 & \cdots & 0 & 0 & -\frac{u_2}{u_1} & -\frac{u_2}{u_1} & \cdots & -\frac{u_2}{u_1} & \frac{l_1 - u_2}{u_1} \\ 0 & 1 & 0 & \cdots & 0 & 0 & \frac{u_2}{u_1} & 0 & \cdots & 0 & \frac{l_2}{u_1} \\ 0 & 0 & 1 & \cdots & 0 & 0 & 0 & \frac{u_2}{u_1} & \cdots & 0 & \frac{l_3}{u_1} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 1 & 0 & 0 & 0 & \cdots & \frac{u_2}{u_1} & \frac{l_n}{u_1} \\ 0 & 0 & 0 & \cdots & 0 & 1 & 1 & 1 & \cdots & 1 & 1 \\ 0 & 0 & 0 & \cdots & 0 & 0 & 0 & 0 & \cdots & 0 & 0 \end{pmatrix}, \tag{20}$$

and we obtain the rank of the coefficient matrix $A$ is equal to the rank of the augmented matrix $(A, \boldsymbol{b})$, which is $n + 1$. Note that $n + 1 < 2n$ when $n > 1$. Thus the above linear equations have infinite solutions. □

The following proposition describes the relation between KD and LD that conducting LD to optimize a localization probability is equivalent to conducting KD to optimize two classification probabilities.

**Proposition 1.** *Let $\boldsymbol{s}$ be the student's predicted probability vector, $u_1$ and $u_2$ are two constants and their summation is 1. We have*

1) *If $\boldsymbol{p}$ and $\boldsymbol{q}$ are two classification probabilities, LD effect on the linear combination $\boldsymbol{l} = u_1\boldsymbol{p} + u_2\boldsymbol{q}$ is equal to the linear combination of KD effects on $\boldsymbol{p}, \boldsymbol{q}$.*
2) *If $\boldsymbol{l}$ is a localization probability, LD effect on $\boldsymbol{l}$ is equal to two KD effects on its decomposition $\boldsymbol{p}$ and $\boldsymbol{q}$.*

*Proof.* We first denote the derivatives of the KD loss of two probabilities $\boldsymbol{s}, \boldsymbol{p}$ w.r.t. a given logit $z_i$ by $\partial KD_i^p$, and $\partial LD_i^p$ likewise for the LD loss.

1) According to Lemma 1, the linear combination $\boldsymbol{l} = u_1\boldsymbol{p} + u_2\boldsymbol{q}$ is well defined and the derivatives of the LD loss of $\boldsymbol{s}, \boldsymbol{l}$ w.r.t. a given logit $z_i$ is given by:

$$\partial LD_i^l = s_{\tau i} - l_{\tau i} \tag{21}$$

$$= u_1 s_{\tau i} + u_2 s_{\tau i} - (u_1 p_{\tau i} + u_2 q_{\tau i}) \tag{22}$$

$$= u_1(s_{\tau i} - p_{\tau i}) + u_2(s_{\tau i} - q_{\tau i}) \tag{23}$$

$$= u_1 \partial KD_i^p + u_2 \partial KD_i^q \tag{24}$$

2) According to Lemma 2, the decomposition of $\boldsymbol{l}$ exists, which is written as $\boldsymbol{l} = u_1\boldsymbol{p} + u_2\boldsymbol{q}$. Then Eq. 24 still holds. □

## A.2 Gradient Rescaling

We first give the lemma in [70].

**Lemma 3.** *Let $q_{\tau t} = p_{\tau t} + c_t + \eta$, where $c_t$ is teacher's relative prediction confidence on the ground-truth class $t$ and $\eta$ is a zero-mean random noise. Then the logit's gradient rescaling factor by applying KD is given by:*

$$\mathbb{E}_\eta \left[ \frac{\partial_t^{KD}}{\partial_t} \right] = \mathbb{E}_\eta \left[ \frac{\sum_{i \neq t} \partial_i^{KD}}{\sum_{i \neq t} \partial_i} \right] = \gamma + \frac{\lambda}{\tau} \left( \frac{c_t}{1 - p_t} \right). \tag{25}$$

Next, we give the corollary of Lemma 3, which shows that LD performs gradient rescaling to distribution focal loss (DFL) [12] in the logits space.

**Corollary 1.** *Let $q_{\tau i} = p_{\tau i} + c_i + \eta_i$, where $c_i$ is teacher's relative prediction confidence at position $i$, $\eta_i$ is a zero-mean random noise. Then the logit's gradient rescaling factor to DFL by applying LD is given by:*

$$\mathbb{E}_\eta \left[ \frac{\partial_i^{LD}}{\partial_i} \right] = \mathbb{E}_\eta \left[ \frac{\sum_{s \neq i} \partial_s^{LD}}{\sum_{s \neq i} \partial_s} \right] = \gamma + \frac{\lambda}{\tau} \left( \frac{c_i}{u_i - p_i} \right), \tag{26}$$

*where $\frac{\partial_i^{LD}}{\partial_i}$ denotes the gradient of LD loss along with DFL w.r.t. logits $z_i$, $u_i$ and $u_j$ are two constants and their summation is 1, $\gamma$ and $\lambda$ are the loss weights of the DFL and LD loss respectively, and $\tau$ is temperature.*

*Proof.* Following [12], DFL is defined as the linear combination of two CE loss at position $i$ and $j$,

$$\mathcal{L}_{DFL} = u_i \mathcal{H}(\boldsymbol{p}, \boldsymbol{g}^i) + u_j \mathcal{H}(\boldsymbol{p}, \boldsymbol{g}^j), \tag{27}$$

where $\boldsymbol{g^i} = \{0, 1\}^n$ are ground-truth labels whose value is 1 at position $i$ and 0 otherwise. One can easily get the gradient of DFL $w.r.t.$ the logit $z_i$,

$$\frac{\partial \mathcal{L}_{\text{DFL}}}{\partial z_i} = u_i(p_i - g_i) + u_j p_i = p_i - u_i, \quad (28)$$

and we still use the notation $\partial_i$ to represent $\frac{\partial \mathcal{L}_{\text{DFL}}}{\partial z_i}$. With LD, the total loss is given by:

$$\mathcal{L}^{LD} = \gamma(u_i \mathcal{H}(\boldsymbol{p}, \boldsymbol{g^i}) + u_j \mathcal{H}(\boldsymbol{p}, \boldsymbol{g^j})) + \lambda \mathcal{H}(\boldsymbol{p}_\tau, \boldsymbol{q}_\tau), \quad (29)$$

The gradient of LD loss along with DFL $w.r.t.$ the logit $z_i \in \boldsymbol{z}_S$ can be represented as:

$$\frac{\partial \mathcal{L}^{LD}}{\partial z_i} = \gamma u_i(p_i - g_i) + \gamma u_j p_i + \frac{\lambda}{\tau}(p_{\tau i} - q_{\tau i}), \quad (30)$$

and we still denote $\partial_i^{LD} = \frac{\partial \mathcal{L}_{\text{LD}}}{\partial z_i}$. According to Lemma 3, we have the ratio of Eq. 30 and Eq. 28,

$$\mathbb{E}_\eta \left[ \frac{\partial_i^{LD}}{\partial_i} \right] = \gamma u_i \frac{p_i - g_i}{p_i - u_i} + \gamma \frac{u_j p_i}{p_i - u_i} - \frac{\lambda}{\tau} \frac{c_i}{p_i - u_i} \quad (31)$$

$$= \gamma + \frac{\lambda}{\tau} \frac{c_i}{u_i - p_i}. \quad (32)$$

Thus, the sum of the incorrect position gradients is given by:

$$\begin{aligned}
&\sum_{s \neq i} \partial_s^{LD} \\
&= \gamma u_i \sum_{s \neq i} p_s + \gamma u_j \sum_{s \neq i,j} p_s + \gamma u_j(p_j - g_j) + \frac{\lambda}{\tau} \sum_{s \neq i} (p_{\tau s} - q_{\tau s}) \\
&= \gamma u_i(g_i - p_i) + \gamma u_j(g_i - p_i) - \gamma u_j g_j + \frac{\lambda}{\tau}(q_{\tau s} - p_{\tau s}) \\
&= \gamma u_i(g_i - p_i) - \gamma u_j p_i + \frac{\lambda}{\tau}(q_{\tau s} - p_{\tau s}) \\
&= -\partial_i^{LD}.
\end{aligned} \quad (33)$$

Similarly applies for $\partial_s$, and hence the proof. $\square$

## REFERENCES

[1] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.

[2] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," in *Int. Conf. Learn. Represent.*, 2015.

[3] S. Zagoruyko and N. Komodakis, "Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer," in *Int. Conf. Learn. Represent.*, 2017.

[4] J. Kim, S. Park, and N. Kwak, "Paraphrasing complex network: network compression via factor transfer," in *Adv. Neural Inform. Process. Syst.*, 2018, pp. 2765–2774.

[5] X. Jin, B. Peng, Y. Wu, Y. Liu, J. Liu, D. Liang, J. Yan, and X. Hu, "Knowledge distillation via route constrained optimization," in *Int. Conf. Comput. Vis.*, 2019.

[6] G.-H. Wang, Y. Ge, and J. Wu, "Distilling knowledge by mimicking features," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2021.

[7] G. Chen, W. Choi, X. Yu, T. Han, and M. Chandraker, "Learning efficient object detection models with knowledge distillation," in *Adv. Neural Inform. Process. Syst.*, 2017.

[8] R. Sun, F. Tang, X. Zhang, H. Xiong, and Q. Tian, "Distilling object detectors with task adaptive regularization," *arXiv preprint arXiv:2006.13108*, 2020.

[9] T. Wang, L. Yuan, X. Zhang, and J. Feng, "Distilling object detectors with fine-grained feature imitation," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.

[10] L. Zhang and K. Ma, "Improve object detection with feature-based knowledge distillation: Towards accurate and efficient detectors," in *Int. Conf. Learn. Represent.*, 2020.

[11] Z. Kang, P. Zhang, X. Zhang, J. Sun, and N. Zheng, "Instance-conditional knowledge distillation for object detection," in *Adv. Neural Inform. Process. Syst.*, 2021.

[12] X. Li, W. Wang, L. Wu, S. Chen, X. Hu, J. Li, J. Tang, and J. Yang, "Generalized Focal Loss: learning qualified and distributed bounding boxes for dense object detection," in *Adv. Neural Inform. Process. Syst.*, 2020.

[13] H. Qiu, H. Li, Q. Wu, and H. Shi, "Offset bin classification network for accurate object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.

[14] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2016.

[15] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *Eur. Conf. Comput. Vis.*, 2016.

[16] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Adv. Neural Inform. Process. Syst.*, 2015.

[17] Z. Tian, C. Shen, H. Chen, and T. He, "FCOS: Fully convolutional one-stage object detection," in *Int. Conf. Comput. Vis.*, 2019.

[18] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. Devito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," in *Adv. Neural Inform. Process. Syst.*, 2017.

[19] S.-M. Hu, D. Liang, G.-Y. Yang, G.-W. Yang, and W.-Y. Zhou, "Jittor: A novel deep learning framework with meta-operators and unified graph execution," *Science China Information Sciences*, vol. 63, no. 222103, pp. 1–21, 2020.

[20] Z. Zheng, R. Ye, P. Wang, D. Ren, W. Zuo, Q. Hou, and M. Cheng, "Localization distillation for dense object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022.

[21] S. Zagoruyko and N. Komodakis, "Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer," in *Int. Conf. Learn. Represent.*, 2017.

[22] J.-H. Bae, D. Yeo, J. Yim, N.-S. Kim, C.-S. Pyo, and J. Kim, "Densely distilled flow-based knowledge transfer in teacher-student framework for image classification," *IEEE Transactions on Image Processing*, vol. 29, pp. 5698–5710, 2020.

[23] W. Park, D. Kim, Y. Lu, and M. Cho, "Relational knowledge distillation," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.

[24] S. I. Mirzadeh, M. Farajtabar, A. Li, N. Levine, A. Matsukawa, and H. Ghasemzadeh, "Improved knowledge distillation via teacher assistant," in *Association for the Advancement of Artificial Intelligence*, 2020.

[25] W. Son, J. Na, J. Choi, and W. Hwang, "Densely guided knowledge distillation using multiple teacher assistants," in *Int. Conf. Comput. Vis.*, 2021.

[26] Q. Li, S. Jin, and J. Yan, "Mimicking very efficient network for object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017.

[27] X. Dai, Z. Jiang, Z. Wu, Y. Bao, Z. Wang, S. Liu, and E. Zhou, "General instance distillation for object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.

[28] J. Guo, K. Han, Y. Wang, H. Wu, X. Chen, C. Xu, and C. Xu, "Distilling object detectors via decoupled features," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.

[29] D. Zhixing, R. Zhang, M. Chang, S. Liu, T. Chen, Y. Chen *et al.*, "Distilling object detectors with feature richness," in *Adv. Neural Inform. Process. Syst.*, 2021.

[30] G. Li, X. Li, Y. Wang, S. Zhang, Y. Wu, and D. Liang, "Knowledge distillation for object detection via rank mimicking and prediction-guided feature imitation," in *Association for the Advancement of Artificial Intelligence*, 2022.

[31] S. Gidaris and N. Komodakis, "Locnet: Improving localization accuracy for object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2016.

[32] J. Wang, K. Chen, S. Yang, C. C. Loy, and D. Lin, "Region proposal by guided anchoring," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.

[33] J. Wang, W. Zhang, Y. Cao, K. Chen, J. Pang, T. Gong, J. Shi, C. C. Loy, and D. Lin, "Side-aware boundary localization for more precise object detection," in *Eur. Conf. Comput. Vis.*, 2020.

[34] C. Zhu, Y. He, and M. Savvides, "Feature selective anchor-free module for single-shot object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.

[35] X. Lu, B. Li, Y. Yue, Q. Li, and J. Yan, "Grid R-CNN," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.

[36] T. Kong, F. Sun, C. Tan, H. Liu, and W. Huang, "Deep feature pyramid reconfiguration for object detection," in *Eur. Conf. Comput. Vis.*, 2018.

[37] X. Yang, J. Yang, J. Yan, Y. Zhang, T. Zhang, Z. Guo, X. Sun, and K. Fu, "Scrdet: Towards more robust detection for small, cluttered and rotated objects," in *Int. Conf. Comput. Vis.*, 2019.

[38] X. Yang, J. Yan, Q. Ming, W. Wang, X. Zhang, and Q. Tian, "Rethinking rotated object detection with gaussian wasserstein distance loss," in *International Conference on Machine Learning (ICML)*, 2021.

[39] X. Yang, X. Yang, J. Yang, Q. Ming, W. Wang, Q. Tian, and J. Yan, "Learning high-precision bounding box for rotated object detection via kullback-leibler divergence," in *Adv. Neural Inform. Process. Syst.*, 2021.

[40] H. Zhang, Y. Wang, F. Dayoub, and N. Sünderhauf, "Varifocalnet: An iou-aware dense object detector," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.

[41] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 9, pp. 1627–1645, 2009.

[42] L. Han, P. Tao, and R. R. Martin, "Livestock detection in aerial images using a fully convolutional network," *Computational Visual Media*, vol. 5, no. 2, p. 221 – 228, 2019.

[43] Z. Cai and N. Vasconcelos, "Cascade R-CNN: Delving into high quality object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018.

[44] J. Pang, K. Chen, J. Shi, H. Feng, W. Ouyang, and D. Lin, "Libra R-CNN: Towards balanced learning for object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.

[45] H. Zhang, H. Chang, B. Ma, N. Wang, and X. Chen, "Dynamic R-CNN: Towards high quality object detection via dynamic training," in *Eur. Conf. Comput. Vis.*, 2020.

[46] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017.

[47] ——, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.

[48] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.

[49] C.-Y. Fu, W. Liu, A. Ranga, A. Tyagi, and A. C. Berg, "DSSD: Deconvolutional single shot detector," *arXiv:1701.06659*, 2017.

[50] P. Zhou, B. Ni, C. Geng, J. Hu, and Y. Xu, "Scale-transferrable object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018.

[51] J. Yu, Y. Jiang, Z. Wang, Z. Cao, and T. Huang, "UnitBox: an advanced object detection network," in *ACM Int. Conf. Multimedia*, 2016.

[52] H. Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, and S. Savarese, "Generalized Intersection over Union: A metric and a loss for bounding box regression," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.

[53] Z. Zheng, P. Wang, W. Liu, J. Li, R. Ye, and D. Ren, "Distance-IoU Loss: Faster and better learning for bounding box regression," in *Association for the Advancement of Artificial Intelligence*, 2020.

[54] Z. Zheng, P. Wang, D. Ren, W. Liu, R. Ye, Q. Hu, and W. Zuo, "Enhancing geometric factors in model learning and inference for object detection and instance segmentation," *IEEE Transactions on Cybernetics*, 2021.

[55] Y. He, C. Zhu, J. Wang, M. Savvides, and X. Zhang, "Bounding box regression with uncertainty for accurate object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.

[56] J. Choi, D. Chun, H. Kim, and H.-J. Lee, "Gaussian YOLOv3: An accurate and fast object detector using localization uncertainty for autonomous driving," in *Int. Conf. Comput. Vis.*, 2019.

[57] X. Li, W. Wang, X. Hu, J. Li, J. Tang, and J. Yang, "Generalized focal loss v2: Learning reliable localization quality estimation for dense object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.

[58] B. Jiang, R. Luo, J. Mao, T. Xiao, and Y. Jiang, "Acquisition of localization confidence for accurate object detection," in *Eur. Conf. Comput. Vis.*, 2018.

[59] Z. Huang, L. Huang, Y. Gong, C. Huang, and X. Wang, "Mask scoring R-CNN," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.

[60] E. Xie, P. Sun, X. Song, W. Wang, X. Liu, D. Liang, C. Shen, and P. Luo, "Polarmask: Single shot instance segmentation with polar representation," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.

[61] Y. Zhou, X. Yang, G. Zhang, J. Wang, Y. Liu, L. Hou, X. Jiang, X. Liu, J. Yan, C. Lyu, W. Zhang, and K. Chen, "Mmrotate: A rotated object detection benchmark using pytorch," in *ACM Int. Conf. Multimedia*, 2022.

[62] J. Ma, W. Shao, H. Ye, L. Wang, H. Wang, Y. Zheng, and X. Xue, "Arbitrary-oriented scene text detection via rotation proposals," *IEEE Transactions on Multimedia*, vol. 20, no. 11, pp. 3111–3122, 2018.

[63] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Int. Conf. Comput. Vis.*, 2017.

[64] W. Qian, X. Yang, S. Peng, Y. Guo, and J. Yan, "Learning modulated loss for rotated object detection," in *Association for the Advancement of Artificial Intelligence*, 2021.

[65] X. Yang and J. Yan, "Arbitrary-oriented object detection with circular smooth label," in *Eur. Conf. Comput. Vis.*, 2020.

[66] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard *et al.*, "{TensorFlow}: A system for {Large-Scale} machine learning," in *12th USENIX symposium on operating systems design and implementation (OSDI 16)*, 2016, pp. 265–283.

[67] Z. Chen, K. Chen, W. Lin, J. See, H. Yu, Y. Ke, and C. Yang, "Piou loss: Towards accurate oriented object detection in complex environments," in *Eur. Conf. Comput. Vis.*, 2020.

[68] X. Yang, Y. Zhou, G. Zhang, J. Yang, W. Wang, J. Yan, X. Zhang, and Q. Tian, "The kfiou loss for rotated object detection," *arXiv preprint arXiv:2201.12558*, 2022.

[69] X. Yang, Q. Liu, J. Yan, A. Li, Z. Zhang, and G. Yu, "R3det: Refined single-stage detector with feature refinement for rotating object," in *Association for the Advancement of Artificial Intelligence*, 2021.

[70] J. Tang, R. Shivanna, Z. Zhao, D. Lin, A. Singh, E. H. Chi, and S. Jain, "Understanding and improving knowledge distillation," *arXiv preprint arXiv:2002.03532*, 2020.

[71] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Eur. Conf. Comput. Vis.*, 2014.

[72] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, 2010.

[73] G.-S. Xia, X. Bai, J. Ding, Z. Zhu, S. Belongie, J. Luo, M. Datcu, M. Pelillo, and L. Zhang, "Dota: A large-scale dataset for object detection in aerial images," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018, pp. 3974–3983.

[74] K. Chen, J. Wang, J. Pang, Y. Cao, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Xu, Z. Zhang, D. Cheng, C. Zhu, T. Cheng, Q. Zhao, B. Li, X. Lu, R. Zhu, Y. Wu, J. Dai, J. Wang, J. Shi, W. Ouyang, C. C. Loy, and D. Lin, "MMDetection: Open mmlab detection toolbox and benchmark," *arXiv preprint arXiv:1906.07155*, 2019.

[75] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2016.

[76] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017.

[77] S. Zhang, C. Chi, Y. Yao, Z. Lei, and S. Z. Li, "Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.

[78] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein, "Visualizing the loss landscape of neural nets," *Adv. Neural Inform. Process. Syst.*, 2018.

[79] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017.

[80] X. Zhu, H. Hu, S. Lin, and J. Dai, "Deformable convnets v2: More deformable, better results," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.

[81] S.-H. Gao, M.-M. Cheng, K. Zhao, X.-Y. Zhang, M.-H. Yang, and P. Torr, "Res2net: A new multi-scale backbone architecture," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 2, pp. 652–662, 2021.

[82] C. Zhu, F. Chen, Z. Shen, and M. Savvides, "Soft anchor-point object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.

[83] H. Qiu, Y. Ma, Z. Li, S. Liu, and J. Sun, "Borderdet: Border feature for dense object detection," in *Eur. Conf. Comput. Vis.*, 2020.

[84] B. Zhu, J. Wang, Z. Jiang, F. Zong, S. Liu, Z. Li, and J. Sun, "Autoassign: Differentiable label assignment for dense object detection," *arXiv preprint arXiv:2007.03496*, 2020.

[85] K. Kim and H. S. Lee, "Probabilistic anchor assignment with iou prediction for object detection," in *Eur. Conf. Comput. Vis.*, 2020.

[86] Z. Ge, S. Liu, Z. Li, O. Yoshie, and J. Sun, "OTA: Optimal transport assignment for object detection," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.

[87] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *Eur. Conf. Comput. Vis.*, 2020.

[88] T. Furlanello, Z. Lipton, M. Tschannen, L. Itti, and A. Anandkumar, "Born again neural networks," in *International Conference on Machine Learning (ICML)*, 2018, pp. 1607–1616.

[89] L. Zhang, J. Song, A. Gao, J. Chen, C. Bao, and K. Ma, "Be your own teacher: Improve the performance of convolutional neural networks via self distillation," in *Int. Conf. Comput. Vis.*, 2019, pp. 3713–3722.

[90] B. Zhao, Q. Cui, R. Song, Y. Qiu, and J. Liang, "Decoupled knowledge distillation," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022.

[91] Y.-H. Wu, Y. Liu, X. Zhan, and M.-M. Cheng, "P2T: Pyramid pooling transformer for scene understanding," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2022.

[92] Q. Hou, C.-Z. Lu, M.-M. Cheng, and J. Feng, "Conv2former: A simple transformer-style convnet for visual recognition," *arXiv preprint arXiv:2211.11943*, 2022.

[93] Z. Dai, H. Liu, Q. Le, and M. Tan, "Coatnet: Marrying convolution and attention for all data sizes," *Adv. Neural Inform. Process. Syst.*, vol. 34, 2021.

[94] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A convnet for the 2020s," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022.

[95] M.-H. Guo, J.-X. Cai, Z.-N. Liu, T.-J. Mu, R. R. Martin, and S.-M. Hu, "Pct: Point cloud transformer," *Computational Visual Media*, vol. 7, no. 2, p. 187 – 199, 2021.

**Zhaohui Zheng** received the M.S. degree in computational mathematics from Tianjin University in 2021. He is currently a Ph.D. candidate with the School of Computer Science at Nankai University, Tianjin, China. His research interests include object detection, instance segmentation and knowledge distillation.



**Rongguang Ye** received the B.S. and M.S. degrees from the School of Mathematics, Tianjin University, Tianjin, China, in 2019 and 2022. He is now working at Intel Asia-Pacific Research And Development Ltd as an AI framework engineer. His research interests include object detection and computer vision.



**Qibin Hou** received his Ph.D. degree from the School of Computer Science, Nankai University. Then, he worked at the National University of Singapore as a research fellow. Now, he is an associate professor at School of Computer Science, Nankai University. He has published more than 30 papers on top conferences/journals, including T-PAMI, CVPR, ICCV, NeurIPS, etc. His research interests include deep learning and computer vision.



**Dongwei Ren** received two Ph.D. degrees in computer application technology from Harbin Institute of Technology and The Hong Kong Polytechnic University in 2017 and 2018, respectively. From 2018 to 2021, he was an Assistant Professor with the College of Intelligence and Computing, Tianjin University. He is currently an Associate Professor with the School of Computer Science and Technology, Harbin Institute of Technology. His research interests include computer vision and deep learning.



**Ping Wang** received the B.S., M.S., and Ph.D. degrees in computer science from Tianjin University, Tianjin, China, in 1988, 1991, and 1998, respectively. She is currently a Professor with the School of Mathematics, Tianjin University. Her research interests include image processing and machine learning.



**Wangmeng Zuo** received the Ph.D. degree from the Harbin Institute of Technology in 2007. He is currently a Professor in the School of Computer Science and Technology, Harbin Institute of Technology. His research interests include image enhancement and restoration, image and face editing, object detection, visual tracking, and image classification. He has published over 100 papers in top tier journals and conferences. His publications have been cited more than 30,000 times in literature. He is on the editorial boards of IEEE TPAMI and IEEE TIP.



**Ming-Ming Cheng** received his PhD degree from Tsinghua University in 2012. Then he did 2 years research fellow, with Prof. Philip Torr in Oxford. He is now a professor at Nankai University, leading the Media Computing Lab. His research interests include computer graphics, computer vision, and image processing. He received research awards including National Science Fund for Distinguished Young Scholars and ACM China Rising Star Award. He is on the editorial boards of IEEE TPAMI and IEEE TIP.