
Learning High-Precision Bounding Box for Rotated Object Detection via Kullback-Leibler Divergence

Xue Yang^{1,*}, Xiaojiang Yang¹, Jirui Yang², Qi Ming³, Wentao Wang¹, Qi Tian⁴, Junchi Yan¹

¹Shanghai Jiao Tong University

²University of Chinese Academy of Sciences

³Beijing Institute of Technology ⁴Huawei Inc.

yangxue-2019-sjtu.edu.cn

Abstract

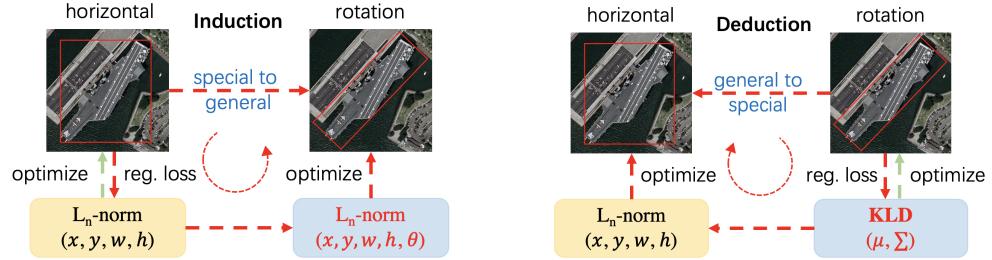
Existing rotated object detectors are mostly inherited from the horizontal detection paradigm, as the latter has evolved into a well-developed area. However, these detectors are difficult to perform prominently in high-precision detection due to the limitation of current regression loss design, especially for objects with large aspect ratios. Taking the perspective that horizontal detection is a special case for rotated object detection, in this paper, we are motivated to change the design of rotation regression loss from induction paradigm to deduction methodology, in terms of the relation between rotation and horizontal detection. We show that one essential challenge is how to modulate the coupled parameters in the rotation regression loss, as such the estimated parameters can influence to each other during the dynamic joint optimization, in an adaptive and synergetic way. Specifically, we first convert the rotated bounding box into a 2-D Gaussian distribution, and then calculate the Kullback-Leibler Divergence (KLD) between the Gaussian distributions as the regression loss. By analyzing the gradient of each parameter, we show that KLD (and its derivatives) can dynamically adjust the parameter gradients according to the characteristics of the object. It will adjust the importance (gradient weight) of the angle parameter according to the aspect ratio. This mechanism can be vital for high-precision detection as a slight angle error would cause a serious accuracy drop for large aspect ratios objects. More importantly, we have proved that KLD is scale invariant. We further show that the KLD loss can be degenerated into the popular l_n -norm loss for horizontal detection. Experimental results on seven datasets using different detectors show its consistent superiority.

1 Introduction

As a fundamental building block for visual analysis across aerial images, scene text etc., rotated object detection has recently been developed rapidly [1, 2, 3, 4, 5], which benefit themselves from the well-established horizontal detection approaches [6, 7, 8, 9, 10]. Specifically, many works [11, 12, 13, 14] build themselves upon the previously established horizontal box detection pipeline from an inductive perspective, as shown in Figure 1(a). However, these detectors are often unable to cope with challenging scenes well due to the limitations of current regression loss, such as large aspect ratio objects, dense scenes, etc., resulting in obvious disadvantages in high-precision detection.

In this paper, we take a step back, and aim to develop (from a deductive perspective) a unified regression framework for rotation detection and its special case: horizontal detection. In fact, our new framework enjoys a coherent property that it can be degenerated into the current commonly used regression loss (e.g. l_n -norm) in special cases (horizontal detection), as shown in Figure 1(b).

*Work done during an internship at Huawei Inc.



(a) Previous methods follow the induction paradigm from special horizontal to general rotated detection. (b) Our proposed method adopts a deduction methodology from general rotated to special horizontal detection.

Figure 1: Methodological road-map difference between horizontal detection (special case) and rotation detection (general case) in the previous methods [1, 11, 12, 13, 14] and the proposed method.

For a devising a rotation regression loss for high-precision rotation detection, one important observation is that the importance of different parameters to different types of objects can vary. For example, the angle parameter (θ) and the center point parameter (x, y) are important for large aspect ratio objects and small objects, respectively. In another word, it is conjectured that regression loss should be self-modulated during the learning process and calls for more dynamic optimization strategy.

Inspired by the above ideas, we first convert the rotated bounding box $\mathcal{B}(x, y, h, w, \theta)$ into a 2-D Gaussian distribution $\mathcal{N}(\mu, \Sigma)$. As a standard distance metric, we then use the Kullback-Leibler Divergence (KLD) [15] to calculate the distribution distance between the predicted bounding box and ground truth as the regression loss. We compare KLD with Smooth L1 loss [6] and another distance metric, Gaussian Wasserstein Distance (GWD) [5, 16], and find that KLD has a more complete parameter optimization mechanism. In particular, by analyzing the gradient of the parameters during learning, we show that the optimization of one parameter will be affected by other parameters (as the gradient weight). It means that the model will adaptively adjust the optimization strategy given a specific configuration of an object for detection, as shown can lead to excellent performance in high-precision detection. In addition, KLD is proven scale invariant, which is an important property that Smooth L1 loss and GWD do not possess. As the horizontal bounding box is a special case of the rotated bounding box, we show that KLD can also be degenerated into the l_n -norm loss as commonly used in existing horizontal detection pipeline. **The highlights of this paper are four-folds:**

- 1) Differing from the dominant existing practices that build rotation detectors heavily upon the horizontal detectors, we develop new rotation detection loss from scratch and show that it is coherent with existing horizontal detection protocol in its degenerated case for horizontal detection.
- 2) To achieve a more principled measurement between the prediction and ground truth, instead of computing the difference for each physically-meaningful parameter related to the bounding box which are in different scales and units, we innovatively convert the regression loss of rotation detection into the KLD of two 2-D Gaussian distributions, leading to a clean and coherent regression loss.
- 3) Through the gradient analysis of each parameter in KLD, we further find that the self-modulated optimization mechanism of KLD greatly promotes the improvement of high-precision detection, which verify the advantage of our loss design. More importantly, we have theoretically shown (in appendix) that KLD is scale invariant for detection, which is crucial for the rotation cases.
- 4) Extensive experimental results on seven public datasets and two popular detectors show the effectiveness of our approach, which achieves new state-of-the-art performance for rotation detection.

2 Background

We first generally discuss the related works on both horizontal and rotated object detection. Then we summarize the current design paradigm of rotation regression loss from two kinds of methodologies, as shown in Figure 1: one is inductive that tries to develop the general rotation detection from the special and classic horizontal detection pipeline. While the other is deductive that aims to devise a general rotation detection pipeline with horizontal detection as its special case.

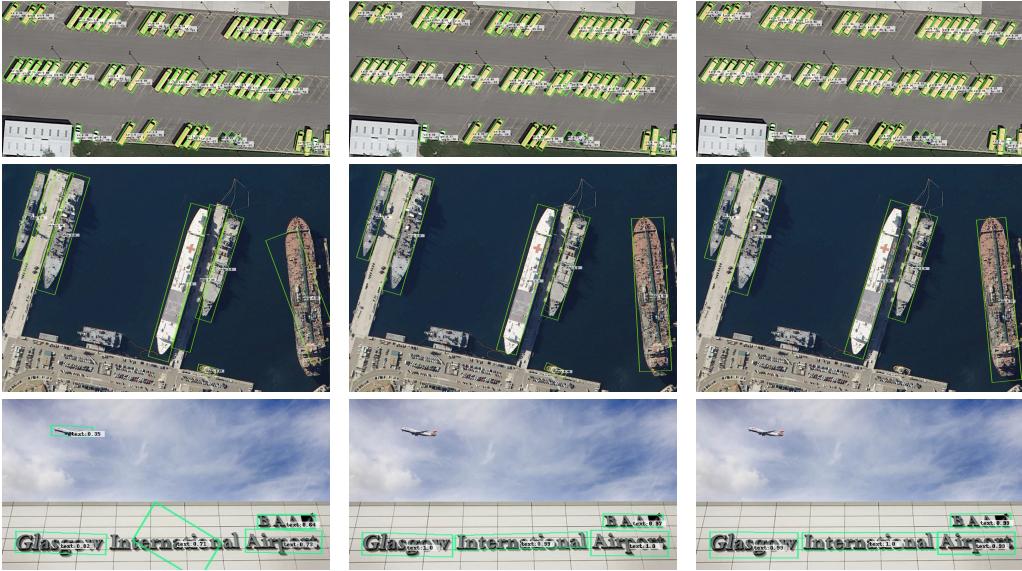


Figure 2: Visual comparison between Smooth L1 loss (left), GWD (middle) and KLD (right).

2.1 Related Works

Horizontal object detection. Horizontal object detection which covers most existing detection literature, normally uses a horizontal bounding box to represent the object. The mainstream classical object detection algorithms can be roughly divided according to the following standards: Two-[6, 7, 8, 10] or Single-stage [9, 17, 18] object detection, Anchor-free [19, 20, 21] or Anchor-based [7, 8, 9] object detection and CNN [7, 9, 19] or Transformer-based [22, 23] object detection. Although the pipelines may vary, the mainstream regression loss often uses the popular l_n -norm loss (such as smooth L1 loss) or IoU-based loss (such as GIoU [24], and DIoU [25]). These above-mentioned detectors have also been widely used in other scenarios and have achieved satisfactory performance. However, horizontal detectors do not provide accurate orientation and scale information.

Rotated object detection. Recent advances in rotation detection [3, 4, 11, 13, 26] are mainly driven by adapting the horizontal object detectors with rotated bounding boxes to represent multi-oriented objects. To accurately predict the rotated bounding box, most rotation detection methods extend the l_n -norm [11, 14, 27, 28, 29] used in horizontal detection, or construct a differentiable approximate IoU loss [3, 5, 30]. From scratch, we try to change the design of rotation regression loss from induction paradigm to deduction methodology, which in fact is a generalization to the horizontal case.

In the following, we describe the existing works from the induction and deduction methodologies.

2.2 Inductive Thinking of Loss Design: from Special Horizon to General Rotation Detection

Regression loss is a vital part of most current object detection algorithms. For horizontal bounding box regression, the model [6, 7, 8, 9, 10] mainly outputs four items for location and size:

$$t_x^p = \frac{x_p - x_a}{w_a}, t_y^p = \frac{y_p - y_a}{h_a}, t_w^p = \ln\left(\frac{w_p}{w_a}\right), t_h^p = \ln\left(\frac{h_p}{h_a}\right) \quad (1)$$

to match the four targets from the ground truth

$$t_x^t = \frac{x_t - x_a}{w_a}, t_y^t = \frac{y_t - y_a}{h_a}, t_w^t = \ln\left(\frac{w_t}{w_a}\right), t_h^t = \ln\left(\frac{h_t}{h_a}\right) \quad (2)$$

where x, y, h, w denote the center coordinates, height and width, respectively. Variables x_t, x_a, x_p are for the ground-truth box, anchor box, and predicted box, respectively (likewise for y, w, h).

Extending the above horizontal case, existing rotation detection models [1, 11, 12, 13, 14] also use regression loss which simply involves an extra angle parameter θ :

$$t_\theta^p = f(\theta_p - \theta_a), t_\theta^t = f(\theta_t - \theta_a) \quad (3)$$

where $f(\cdot)$ is used to deal with angular periodicity, such as trigonometric functions, modulo, etc.

The overall regression loss for rotation detection is:

$$L_{reg} = l_n\text{-norm}(\Delta t_x, \Delta t_y, \ln \Delta t_w, \ln \Delta t_h, \Delta t_\theta) \quad (4)$$

where $\Delta t_x = t_x^p - t_x^t = \frac{\Delta x}{w_a}$, $\Delta t_y = t_y^p - t_y^t = \frac{\Delta y}{h_a}$, $\Delta t_w = t_w^p - t_w^t = \frac{w_p}{w_t}$, $\Delta t_h = t_h^p - t_h^t = \frac{h_p}{h_t}$, and $\Delta t_\theta = t_\theta^p - t_\theta^t = \Delta\theta$.

It can be seen that parameters are optimized independently, making the loss (or detection accuracy) sensitive to the under-fitting of any of the parameters. This mechanism is fatal to high-precision detection. Taking the left side of Figure 2 as an example, the detection result based on the Smooth L1 loss often shows the deviation of the center point or angle. Moreover, different types of objects have different sensitivity to these five parameters. For example, the angle parameter is very important for detecting objects with large aspect ratios. This requires to select an appropriate set of weights given a specific single object sample during the training, which is nontrivial or even unrealistic.

2.3 Deductive Thinking of Loss Design: from General Rotation to Special Horizon Detection

To break the original inductive design paradigm, we adopt deductive paradigm to construct more accurate rotation regression loss. Here we rephrase the main idea in the recent work [5], which converts a arbitrary-oriented bounding box $\mathcal{B}(x, y, h, w, \theta)$ into a 2-D Gaussian $\mathcal{N}(\mu, \Sigma)$, as illustrated in Figure 3. Then the distance between two Gaussian is calculated as the final loss. Specifically, the conversion is:

$$\begin{aligned} \boldsymbol{\mu} &= (x, y)^\top \\ \boldsymbol{\Sigma}^{1/2} &= \mathbf{R} \boldsymbol{\Lambda} \mathbf{R}^\top = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} \frac{w}{2} & 0 \\ 0 & \frac{h}{2} \end{pmatrix} \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} \\ &= \begin{pmatrix} \frac{w}{2} \cos^2 \theta + \frac{h}{2} \sin^2 \theta & \frac{w-h}{2} \cos \theta \sin \theta \\ \frac{w-h}{2} \cos \theta \sin \theta & \frac{w}{2} \sin^2 \theta + \frac{h}{2} \cos^2 \theta \end{pmatrix} \end{aligned} \quad (5)$$

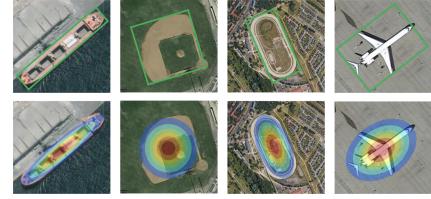


Figure 3: **Top:** rotated box $\mathcal{B}(x, y, h, w, \theta)$. **Bottom:** 2-D Gaussian dist. $\mathcal{N}(\mu, \Sigma)$.

where \mathbf{R} represents the rotation matrix, and $\boldsymbol{\Lambda}$ represents the diagonal matrix of eigenvalues.

The recent work [5] analyzes that the introduction of $\mathcal{N}(\mu, \Sigma)$ can solve the inconsistency between metric and loss, boundary discontinuity and square-like problem. On this basis, we further studies how to design high-precision detection regression loss through new parameter space. Our view is that the self-modulated mechanism is positively correlated with the final high-precision performance.

Gaussian Wasserstein Distance. The Wasserstein distance [5, 16] between two probability measures $\mathbf{X}_p \sim \mathcal{N}_p(\mu_p, \Sigma_p)$ and $\mathbf{X}_t \sim \mathcal{N}_t(\mu_t, \Sigma_t)$ expressed as:

$$\mathbf{D}_w(\mathcal{N}_p, \mathcal{N}_t)^2 = \underbrace{\|\boldsymbol{\mu}_p - \boldsymbol{\mu}_t\|_2^2}_{\text{center distance}} + \underbrace{\text{Tr}(\boldsymbol{\Sigma}_p + \boldsymbol{\Sigma}_t - 2(\boldsymbol{\Sigma}_p^{1/2} \boldsymbol{\Sigma}_t \boldsymbol{\Sigma}_p^{1/2})^{1/2})}_{\text{coupling terms about } h, w \text{ and } \theta} \quad (6)$$

Eq. 6 shows that the Gaussian Wasserstein Distance (GWD) is mainly divided into two parts: the distance between the center points (x, y) and the coupling terms about h, w and θ . Accordingly, the regression loss based on GWD can be regarded as a semi-coupled loss. Although GWD can greatly improve the performance of high-precision rotation detection due to the coupling between part of the parameters, the independent optimization of the center point make the detection result slightly shifted (see Figure 2). Note that GWD is not scale invariant, which is not detection friendly.

When all the boxes are horizontal ($\theta = 0^\circ$), Eq. 6 can be further simplified:

$$\begin{aligned} \mathbf{D}_w^h(\mathcal{N}_p, \mathcal{N}_t)^2 &= \|\boldsymbol{\mu}_p - \boldsymbol{\mu}_t\|_2^2 + \|\boldsymbol{\Sigma}_p^{1/2} - \boldsymbol{\Sigma}_t^{1/2}\|_F^2 \\ &= (x_p - x_t)^2 + (y_p - y_t)^2 + ((w_p - w_t)^2 + (h_p - h_t)^2) / 4 \\ &= l_2\text{-norm}(\Delta x, \Delta y, \Delta w/2, \Delta h/2) \end{aligned} \quad (7)$$

where $\|\cdot\|_F$ is the Frobenius norm. Although Eq. 7 can still be used as the regression loss of horizontal detection, Eq. 4 and 7 are not completely consistent.

Although GWD scheme has played a preliminary exploration of the deductive paradigm, it does not focus on achieving high-precision detection and scale invariance. In the following, we will propose our new approach based on the Kullback-Leibler divergence (KLD) [15].

3 Proposed Approach

3.1 Kullback-Leibler Divergence

To explore the more appropriate regression loss, we adopt the Kullback-Leibler divergence (KLD) [15]. Similarly, the KLD between two 2-D Gaussian is:

$$\mathbf{D}_{kl}(\mathcal{N}_p || \mathcal{N}_t) = \underbrace{\frac{1}{2}(\boldsymbol{\mu}_p - \boldsymbol{\mu}_t)^\top \boldsymbol{\Sigma}_t^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_t)}_{\text{term about } x_p \text{ and } y_p} + \underbrace{\frac{1}{2}\text{Tr}(\boldsymbol{\Sigma}_t^{-1} \boldsymbol{\Sigma}_p) + \frac{1}{2}\ln \frac{|\boldsymbol{\Sigma}_t|}{|\boldsymbol{\Sigma}_p|} - 1}_{\text{coupling terms about } h_p, w_p \text{ and } \theta_p} \quad (8)$$

or

$$\mathbf{D}_{kl}(\mathcal{N}_t || \mathcal{N}_p) = \underbrace{\frac{1}{2}(\boldsymbol{\mu}_p - \boldsymbol{\mu}_t)^\top \boldsymbol{\Sigma}_p^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_t) + \frac{1}{2}\text{Tr}(\boldsymbol{\Sigma}_p^{-1} \boldsymbol{\Sigma}_t) + \frac{1}{2}\ln \frac{|\boldsymbol{\Sigma}_p|}{|\boldsymbol{\Sigma}_t|} - 1}_{\text{chain coupling of all parameters}} \quad (9)$$

It can be seen that each item in $\mathbf{D}_{kl}(\mathcal{N}_t || \mathcal{N}_p)$ is composed of partial parameter coupling, which makes all parameters form a chain coupling relationship. In the optimization process of the KLD-based detector, the parameters influence each other and are jointly optimized which make optimization mechanism of the model is self-modulated. In contrast, $\mathbf{D}_{kl}(\mathcal{N}_p || \mathcal{N}_t)$ and GWD are both semi-coupled, but $\mathbf{D}_{kl}(\mathcal{N}_p || \mathcal{N}_t)$ has a better central point optimization mechanism.

Although KLD is asymmetric, we find that the optimization principles of these two forms are similar by analyzing the gradients of various parameters and experimental results. Take the relatively simple $\mathbf{D}_{kl}(\mathcal{N}_p || \mathcal{N}_t)$ as an example, according to Eq. 5, each item of Eq. 8 can be expressed as

$$(\boldsymbol{\mu}_p - \boldsymbol{\mu}_t)^\top \boldsymbol{\Sigma}_t^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_t) = \frac{4(\Delta x \cos \theta_t + \Delta y \sin \theta_t)^2}{w_t^2} + \frac{4(\Delta y \cos \theta_t - \Delta x \sin \theta_t)^2}{h_t^2} \quad (10)$$

$$\text{Tr}(\boldsymbol{\Sigma}_t^{-1} \boldsymbol{\Sigma}_p) = \frac{h_p^2}{w_t^2} \sin^2 \Delta \theta + \frac{w_p^2}{h_t^2} \sin^2 \Delta \theta + \frac{h_p^2}{h_t^2} \cos^2 \Delta \theta + \frac{w_p^2}{w_t^2} \cos^2 \Delta \theta \quad (11)$$

$$\ln \frac{|\boldsymbol{\Sigma}_t|}{|\boldsymbol{\Sigma}_p|} = \ln \frac{h_t^2}{h_p^2} + \ln \frac{w_t^2}{w_p^2} \quad (12)$$

where $\Delta x = x_p - x_t$, $\Delta y = y_p - y_t$, $\Delta \theta = \theta_p - \theta_t$.

3.2 Advanced Analysis

Analysis of high-precision detection. Without loss of generality, we set $\theta_t = 0^\circ$, then

$$\frac{\partial f_{kl}(\boldsymbol{\mu}_p)}{\partial \boldsymbol{\mu}_p} = \left(\frac{4}{w_t^2} \Delta x, \frac{4}{h_t^2} \Delta y \right)^\top \quad (13)$$

The weights $1/w_t^2$ and $1/h_t^2$ will make the model dynamically adjust the optimization of the object position according to the scale. For example, when the object scale is small or an edge is too short, the model will pay more attention to the optimization of the offset of the corresponding direction. For this kind of object, a slight deviation on the corresponding direction will often cause a sharp drop in IoU. When $\theta_t \neq 0^\circ$, the gradient of the object offset (Δx and Δy) will be dynamically adjusted according to the θ_t for better optimization. In contrast, the gradient of the center point in GWD and L₂-norm are $\frac{\partial f_w(\boldsymbol{\mu}_p)}{\partial \boldsymbol{\mu}_p} = (2\Delta x, 2\Delta y)^\top$ and $\frac{\partial f_{L_2}(\boldsymbol{\mu}_p)}{\partial \boldsymbol{\mu}_p} = (\frac{2}{w_a^2} \Delta x, \frac{2}{h_a^2} \Delta y)^\top$. The former cannot adjust the dynamic gradient according to the length and width of the object. The latter is based on the length and width of the anchor (w_a, h_a) to adjust the gradient instead of the target object (w_t, h_t), which is almost ineffective for those detectors [3, 12, 14, 26, 27, 31, 32] that use horizontal anchors for rotation detection. More importantly, they are not related to the angle of the target object. Therefore, the detection result of the GWD-based and L_n-norm models will show a slight deviation, while the detection result of the KLD-based model is quite accurate, as shown in Figure 2.

For h_p and w_p , we have

$$\frac{\partial f_{kl}(\Sigma_p)}{\partial \ln h_p} = \frac{h_p^2}{h_t^2} \cos^2 \Delta\theta + \frac{h_p^2}{w_t^2} \sin^2 \Delta\theta - 1, \quad \frac{\partial f_{kl}(\Sigma_p)}{\partial \ln w_p} = \frac{w_p^2}{w_t^2} \cos^2 \Delta\theta + \frac{w_p^2}{h_t^2} \sin^2 \Delta\theta - 1 \quad (14)$$

On the one hand, the optimization of the h_p and w_p is affected by the $\Delta\theta$. When $\Delta\theta = 0^\circ$, $\frac{\partial f_{kl}(\Sigma_p)}{\partial \ln h_p} = \frac{h_p^2}{h_t^2} - 1$, $\frac{\partial f_{kl}(\Sigma_p)}{\partial \ln w_p} = \frac{w_p^2}{w_t^2} - 1$, which means that the smaller targeted height or width leads to heavier penalty on its matching loss. This is desirable, as smaller height or width needs higher matching precision. On the other hand, the optimization of $\Delta\theta$ is also affected by h_p and w_p :

$$\frac{\partial f_{kl}(\Sigma_p)}{\partial \theta_p} = \left(\frac{h_p^2 - w_p^2}{w_t^2} + \frac{w_p^2 - h_p^2}{h_t^2} \right) \sin 2\Delta\theta \quad (15)$$

when $w_p = w_t, h_p = h_t$, then $\frac{\partial f_{kl}(\Sigma_p)}{\partial \theta_p} = \left(\frac{h_t^2}{w_t^2} + \frac{w_t^2}{h_t^2} - 2 \right) \sin 2\Delta\theta \geq \sin 2\Delta\theta$, the condition for the equality sign is $h_t = w_t$. This shows that the larger the aspect ratio of the object, the model will pay more attention to the optimization of the angle. This is the main reason why the KLD-based model has a huge advantage in high-precision detection indicators as a slight angle error would cause a serious accuracy drop for large aspect ratios objects. Through the above analysis, we find that when one of the parameters is optimized, the other parameters will be used as its weight to dynamically adjust the optimization rate. In other words, the optimization of parameters is no longer independent, that is, optimizing one parameter will also promote the optimization of other parameters. The optimization of this virtuous circle is the key to KLD as an excellent rotation regression loss. In addition, $\mathbf{D}_{kl}(\mathcal{N}_t || \mathcal{N}_p)$ has similar properties, refer to appendix for details.

Scale invariance. For a full-rank matrix \mathbf{M} , $|\mathbf{M}| \neq 0$, we have $\mathbf{D}_{kl}(\mathcal{N}_p || \mathcal{N}_t) = \mathbf{D}_{kl}(\mathcal{N}'_p || \mathcal{N}'_t)$, where $\mathbf{X}'_p = \mathbf{M}\mathbf{X}_p \sim \mathcal{N}_p(\mathbf{M}\boldsymbol{\mu}_p, \mathbf{M}\boldsymbol{\Sigma}_p\mathbf{M}^\top)$, $\mathbf{X}'_t = \mathbf{M}\mathbf{X}_t \sim \mathcal{N}_t(\mathbf{M}\boldsymbol{\mu}_t, \mathbf{M}\boldsymbol{\Sigma}_t\mathbf{M}^\top)$. Therefore, the affine invariance (including scale invariance when $\mathbf{M} = k\mathbf{I}$, where \mathbf{I} denotes identity matrix) of KLD can be proven (see proof in appendix). Compared with L_n -norm and GWD, KLD is more suitable for replacing the non-differentiable rotated IoU loss for its consistency with detection metric.

Horizontal special case. For horizontal detection, combine Eq. 8 to Eq. 24, we have

$$\begin{aligned} \mathbf{D}_{kl}^h(\mathcal{N}_p || \mathcal{N}_t) &= \frac{1}{2} \left(\frac{w_p^2}{w_t^2} + \frac{h_p^2}{h_t^2} + \frac{4\Delta^2 x}{w_t^2} + \frac{4\Delta^2 y}{h_t^2} + \ln \frac{w_t^2}{w_p^2} + \ln \frac{h_t^2}{h_p^2} - 2 \right) \\ &= 2l_2\text{-norm}(\Delta t_x, \Delta t_y) + l_1\text{-norm}(\ln \Delta t_w, \ln \Delta t_h) + \frac{1}{2}l_2\text{-norm}\left(\frac{1}{\Delta t_w}, \frac{1}{\Delta t_h}\right) - 1 \end{aligned} \quad (16)$$

where the first two terms of Eq. 16 are very similar to Eq. 4, and the divisor part of the two terms x and y is the main difference ($\frac{\Delta x}{w_t}$ vs $\frac{\Delta x}{w_a}$).

Variants of KLD. We have also introduced some variants of KLD to further verify the influence of asymmetry on rotation detection can be ignored. The variants mainly including

$$\begin{aligned} \mathbf{D}_{kl_min(max)}(\mathcal{N}_p || \mathcal{N}_t) &= \min(\max)(\mathbf{D}_{kl}(\mathcal{N}_p || \mathcal{N}_t), \mathbf{D}_{kl}(\mathcal{N}_t || \mathcal{N}_p)) \\ \mathbf{D}_{js}(\mathcal{N}_p || \mathcal{N}_t) &= \frac{1}{2} \left(\mathbf{D}_{kl}\left(\mathcal{N}_t || \frac{\mathcal{N}_p + \mathcal{N}_t}{2}\right) + \mathbf{D}_{kl}\left(\mathcal{N}_p || \frac{\mathcal{N}_p + \mathcal{N}_t}{2}\right) \right) [33] \\ \mathbf{D}_{jef}(\mathcal{N}_p || \mathcal{N}_t) &= \mathbf{D}_{kl}(\mathcal{N}_t || \mathcal{N}_p) + \mathbf{D}_{kl}(\mathcal{N}_p || \mathcal{N}_t) [34] \end{aligned} \quad (17)$$

3.3 Rotation Regression Loss

We normalize the distance function as our final regression loss \mathcal{L}_{reg} :

$$\mathcal{L}_{reg} = 1 - \frac{1}{\tau + f(\mathbf{D})}, \quad \tau \geq 1 \quad (18)$$

where $f(\cdot)$ denotes a non-linear function to transform the distance \mathbf{D} to make the loss more smooth and expressive. In this paper, we mainly use two nonlinear functions, $\text{sqrt}(\mathbf{D})$ and $\ln(\mathbf{D} + 1)$. The hyperparameter τ modulates the entire loss. The multi-task loss is:

$$\mathcal{L} = \frac{\lambda_1}{N_{pos}} \sum_{n=1}^{N_{pos}} \mathcal{L}_{reg}(b_n, gt_n) + \frac{\lambda_2}{N} \sum_{n=1}^N \mathcal{L}_{cls}(p_n, t_n) \quad (19)$$

where N_{pos} and N indicate the number of positive and all anchors. b_n denotes the n -th bounding box, gt_n is the n -th target ground-truth. t_n denotes the label of n -th object, p_n is the n -th probability distribution of various classes calculated by sigmoid function. The hyper-parameter λ_1 , λ_2 control the trade-off and are set to {2, 1} by default. The classification loss L_{cls} is set as focal loss [9].

4 Experiment

4.1 Datasets and Implementation Details

Our experiments are conducted over a variety of datasets, including three large-scale public datasets for aerial images i.e. DOTA [35], UCAS-AOD [36], HRSC2016 [37], as well as scene text dataset ICDAR2015 [38], MLT [39] and MSRA-TD500 [40].

DOTA is one of the largest dataset for oriented object detection in aerial images with three released versions: DOTA-v1.0, DOTA-v1.5 and DOTA-v2.0. DOTA-v1.0 contains 15 common categories, 2,806 images and 188,282 instances. The proportions of the training set, validation set, and testing set in DOTA-v1.0 are 1/2, 1/6, and 1/3, respectively. In contrast, DOTA-v1.5 uses the same images as DOTA-v1.0, but extremely small instances (less than 10 pixels) are also annotated. Moreover, a new category, containing 402,089 instances in total is added in this version. While DOTA-v2.0 contains 18 common categories, 11,268 images and 1,793,658 instances. Compared to DOTA-v1.5, it further includes the new categories. The 11,268 images in DOTA-v2.0 are split into training, validation, test-dev, and test-challenge sets. We divide the images into 600×600 subimages with an overlap of 150 pixels and scale it to 800×800 , in line with the cropping protocol in literature.

UCAS-AOD contains 1,510 aerial images of approximately $659 \times 1,280$ pixels, with two categories of 14,596 instances in total. In line with [29, 35], we randomly select 1,110 for training and 400 for testing. HRSC2016 contains images from two scenarios including ships on sea and ships close inshore. The training, validation and test set include 436, 181 and 444 images.

ICDAR2015, MLT and MSRA-TD500 are commonly used for oriented scene text detection and spotting. ICDAR2015 includes 1,000 training images and 500 testing images. ICDAR2017 MLT is a multi-lingual text dataset, which includes 7,200 training images, 1,800 validation images and 9,000 testing images. MSRA-TD500 dataset consists of 300 training images and 200 testing images.

We use Tensorflow [41] to implement the proposed methods on a server with Tesla V100 and 32G memory. The experiments are all initialized by ResNet50 [42] by default unless otherwise specified. Weight decay and momentum are set 0.0001 and 0.9, respectively. We employ MomentumOptimizer over 8 GPUs with a total of 8 images per minibatch (1 image per GPU).

All the used datasets are trained by 20 epochs in total, and the learning rate is reduced tenfold at 12 epochs and 16 epochs, respectively. The initial learning rate is set to 5e-4. The number of image iterations per epoch for DOTA-v1.0, DOTA-v1.5, DOTA-v2.0, UCAS-AOD, HRSC2016, ICDAR2015, MLT and MSRA-TD500 are 54k, 64k, 80k, 5k, 10k, 10k, 10k and 5k respectively, and doubled if data augmentation or multi-scale training is used.

4.2 Ablation Study and Further Comparison

Regression loss form and hyperparameter. Table 1 compares three forms of KLD-based regression loss on HRSC2016, including D_{kl} , $f(D_{kl})$ and $\mathcal{L}_{reg}(f(D_{kl}), \tau)$. Due to extreme sensitivity to large errors, the performance of D_{kl} is extremely poor, only **0.20%**. Through a simple nonlinear linear transformation, the performance can be increased to **82.96%** and **83.23%** corresponding to *sqrt* and *log*. We further perform a detailed hyperparameter experiment on the loss \mathcal{L}_{reg} proposed in this paper, and the performance reaches the optimal when $\tau = 1$, $f(D_{kl}) = \log(D_{kl} + 1)$, about **85.25%**. Keeping the same loss pattern, we compare six KLD-based distance functions in Table 2, and conclude that the asymmetry of KLD does not have much impact on performance. In subsequent experiments, we use $\mathcal{L}_{reg}(\log(D_{kl}(\mathcal{N}_p || \mathcal{N}_t)), 1)$ as the basic setting.

High-precision detection experiment. We expect that the designed rotation regression loss can show advantages in high-precision detection. Table 3 shows the comparison of the high-precision detection results of three different regression losses using Smooth L1, GWD and KLD on different datasets and different detectors. For the HRSC206 dataset containing a large number of ship with large aspect ratios, GWD-based RetinaNet has a **11.89%** improvement over Smooth L1 on AP₇₅, KLD even gets

Table 1: Ablation study of the loss form and hyperparameter on HRSC2016.

Loss	\mathbf{D}_{kl}	$f(\mathbf{D}_{kl})$	$\mathcal{L}_G(f(\mathbf{D}_{kl}), \tau)$			
			$\tau = 1$	$\tau = 2$	$\tau = 3$	$\tau = 5$
$f(\mathbf{D}_{kl}) = \text{sqrt}(\mathbf{D}_{kl})$	0.20	82.96	84.85	84.15	75.23	73.32
$f(\mathbf{D}_{kl}) = \log(\mathbf{D}_{kl} + 1)$		83.23	85.25	83.63	80.79	73.44

Table 2: Ablation of different KLD-based regression loss form. The based detector is RetinaNet.

Dataset	$\mathbf{D}_{kl}(\mathcal{N}_p \mathcal{N}_t)$	$\mathbf{D}_{kl}(\mathcal{N}_t \mathcal{N}_p)$	$\mathbf{D}_{kl_min}(\mathcal{N}_p \mathcal{N}_t)$	$\mathbf{D}_{kl_max}(\mathcal{N}_p \mathcal{N}_t)$	$\mathbf{D}_{js}(\mathcal{N}_p \mathcal{N}_t)$	$\mathbf{D}_{jeffreys}(\mathcal{N}_p \mathcal{N}_t)$
DOTA-v1.0	70.17	70.64	70.71	70.55	69.67	70.56
HRSC2016	82.83	83.82	83.60	82.70	84.06	83.66

Table 3: High-precision detection experiment under different regression loss. ‘R’, ‘F’ and ‘G’ indicate random rotation, flipping, and graying, respectively.

Method	Dataset	Data Aug.	Reg. Loss	Hmean ₅₀ /AP ₅₀	Hmean ₆₀ /AP ₆₀	Hmean ₇₅ /AP ₇₅	Hmean ₈₅ /AP ₈₅	Hmean _{50:95} /AP _{50:95}
RetinaNet	HRSC2016	R+F+G	Smooth L1	84.28	74.74	48.42	12.56	47.76
			GWD	85.56 (+1.28)	84.04 (+9.30)	60.31 (+11.89)	17.14 (+4.58)	52.89 (+5.13)
			KLD	87.45 (+3.17)	86.72 (+11.98)	72.39 (+23.97)	27.68 (+15.12)	57.80 (+10.04)
R ³ Det			Smooth L1	88.52	79.01	43.42	4.58	46.18
			GWD	89.43 (+0.91)	88.89 (+9.88)	65.88 (+22.46)	15.02 (+10.44)	56.07 (+9.89)
			KLD	89.97 (+1.45)	89.73 (+10.72)	77.38 (+33.96)	25.12 (+20.54)	61.40 (+15.22)
RetinaNet	MSRA-TD500	R+F+G	Smooth L1	70.98	62.42	36.73	12.56	37.89
			GWD	76.76 (+5.78)	68.58 (+6.16)	44.21 (+7.48)	17.75 (+5.19)	43.62 (+5.73)
			KLD	76.96 (+5.98)	70.08 (+7.66)	46.95 (+10.22)	19.59 (+7.03)	45.24 (+7.35)
R ³ Det	ICDAR2015	F	Smooth L1	69.78	64.15	36.97	8.71	37.73
			GWD	74.29 (+4.51)	68.34 (+4.19)	43.39 (+6.42)	10.50 (+1.79)	41.68 (+3.95)
			KLD	75.32 (+5.54)	69.94 (+5.79)	44.46 (+7.49)	10.70 (+1.99)	42.68 (+4.95)
		R+F	Smooth L1	74.83	69.46	42.02	11.59	41.98
			GWD	76.15 (+1.32)	71.26 (+1.80)	45.59 (+3.57)	11.65 (+0.06)	43.58 (+1.60)
			KLD	77.92 (+3.09)	72.77 (+3.31)	43.27 (+1.25)	11.09 (-0.50)	43.65 (+1.67)
		F	Smooth L1	74.28	68.12	35.73	8.01	39.10
			GWD	75.59 (+1.31)	68.36 (+0.24)	40.24 (+4.51)	9.15 (+1.14)	40.80 (+1.70)
			KLD	77.72 (+2.43)	71.99 (+3.87)	43.95 (+8.22)	10.43 (+2.42)	43.29 (+4.19)
		R+F	Smooth L1	75.53	69.69	37.69	9.03	40.56
			GWD	77.09 (+1.56)	71.52 (+1.83)	41.08 (+3.39)	10.10 (+1.07)	42.17 (+1.61)
			KLD	79.63 (+4.63)	73.30 (+3.61)	43.51 (+5.82)	10.61 (+1.58)	43.61 (+3.05)

Table 4: More ablation experiments on other datasets.

Method	Reg. Loss	MLT	UCAS-AOD	DOTA-v1.0	DOTA-v1.5	DOTA-v2.0
RetinaNet	Smooth L1 GWD KLD	48.42 54.58 (+6.16) 57.59 (+9.17)	94.56 95.44 (+0.88) 96.14 (+1.58)	65.73 68.93 (+3.20) 71.28 (+5.55)	58.87 60.03 (+1.16) 62.50 (+3.63)	44.16 46.65 (+2.49) 47.69 (+3.53)

a **23.97%** gain. Even with a stronger R³Det detector, KLD and GWD still increased by **33.96%** and **22.46%** in AP₇₅, and **15.22%** and **9.89%** in AP_{50:95}. The same experimental conclusion are also reflected in the other two scene text datasets MASR-TF500 and ICDAR2015, which is KLD > GWD > Smooth L1. In general, the self-modulation optimization mechanism has a significant help for high-precision detection. For a more intuitive comparison, we visually compare these three regression losses, as shown in Figure 2. Since the center point (x, y) parameters in Smooth L1 Loss and GWD are independently optimized, their prediction results are slightly shifted. In contrast, the KLD-based prediction results are closer to the object boundary and show strong robustness in dense scenes. Similarly, GWD-based or KLD-based model has more accurate angle prediction capabilities than Smooth L1-based model due to their angle parameters (θ) are not independently optimized.

Ablation study on more datasets. To make the results more credible, we continue to verify on the other five datasets, as shown in Table 4. The improvement of KLD on the three data sets of MLT, UCAS-AOD and DOTA-v1.0 is still considerable, with an increase of **9.17%**, **1.58%**, and **5.55%** respectively. Note that for DOTA-v1.5 and DOTA-v2.0, which contain a large number of small objects (less than 10 pixels), KLD has achieved significant gains of **3.63%** and **3.53%**.

Comparison of peer methods. Table 5 compares the six peer techniques, including IoU-Smooth L1 Loss [3], Modulated loss [43], RIL [32], CSL [4], DCL [44], and GWD [5] on DOTA-v1.0. For fairness, these methods are all implemented on the same baseline method, and are trained and tested under the same environment and hyperparameters. We detail the accuracy of the seven categories, including large aspect ratio (e.g. BR, SV, LV, SH, HA) and square-like object (e.g. ST, RD), which can better reflect the real-world challenges and advantages of our method. Without bells and whistles, the combination of RetinaNet and KLD directly surpasses R³Det (**71.28%** vs **70.66%** in AP₅₀ and **69.41%** vs **68.31%** in 7-AP₅₀). Even combined with R³Det, KLD can still further improve performance of the large aspect ratio object (**2.82%** in 7-AP₅₀) and high-precision detection (**6.07%**

Table 5: Accuracy comparison between different rotation detectors on DOTA dataset. \dagger and \ddagger represent the large aspect ratio object and the square-like object, respectively. The bold red and blue fonts indicate the top two performances respectively. D_{oc} and D_{le} represent OpenCV Definition ($\theta \in [-90^\circ, 0^\circ]$) and Long Edge Definition ($\theta \in [-90^\circ, 90^\circ]$) of RBox.

Baseline	Method	Box Def.	v1.0 tranval/test							v1.0 train/val				v1.5	v2.0	
			BR †	SV †	LV †	SH †	HA †	ST †	RA †	7-AP $_{50}$	AP $_{50}$	AP $_{50}$	AP $_{50:95}$	AP $_{50}$	AP $_{50}$	
RetinaNet	-	D_{oc}	42.17	65.93	51.11	72.61	53.24	78.38	62.00	60.78	65.73	64.70	32.31	34.50	58.87	44.16
	-	D_{le}	38.31	60.48	49.77	68.29	51.28	78.60	60.02	58.11	64.17	62.21	26.06	31.49	56.10	43.06
	IoU-Smooth L1 [3]	D_{oc}	44.32	63.03	51.25	72.78	56.21	77.98	63.22	61.26	66.99	64.61	34.17	36.23	59.16	46.31
	Modulated Loss [43]	D_{oc}	42.92	67.92	52.91	72.67	53.64	80.22	58.21	61.21	66.05	63.50	33.32	34.61	57.75	45.17
	Modulated Loss [43]	Quad.	43.21	70.78	54.70	72.68	60.99	79.72	62.08	63.45	67.20	65.15	40.59	39.12	61.42	46.71
	RIL [32]	Quad.	40.81	67.63	55.45	72.42	55.49	78.09	64.75	62.09	66.06	64.07	40.98	39.05	58.91	45.35
	CSL [4]	D_{le}	42.25	68.28	54.51	72.85	53.10	75.59	58.99	60.80	67.38	64.40	32.58	35.04	58.55	43.34
	DCL (BCL) [44]	D_{le}	41.40	65.82	56.27	73.80	54.30	79.02	60.25	61.55	67.39	65.93	35.66	36.71	59.38	45.46
	GWD [5]	D_{oc}	44.07	71.92	62.56	77.94	60.25	79.64	63.52	65.70	68.93	65.44	38.68	38.71	60.03	46.65
	KLD	D_{oc}	44.00	74.45	72.48	80.30	65.54	80.03	65.05	69.41	71.28	68.14	44.48	42.15	62.50	47.69
R^3 Det [26]	-	D_{oc}	44.15	75.09	72.88	86.04	56.49	82.53	61.01	68.31	70.66	67.18	38.41	38.46	62.91	48.43
	DCL (BCL) [44]	D_{le}	46.84	74.87	74.96	85.70	57.72	84.06	63.77	69.70	71.21	67.45	35.44	37.54	61.98	48.71
	GWD [5]	D_{oc}	46.73	75.84	78.00	86.71	62.69	83.09	61.12	70.60	71.56	69.28	43.35	41.56	63.22	49.25
	KLD	D_{oc}	48.34	75.09	78.88	86.52	65.48	82.08	61.51	71.13	71.73	68.87	44.48	42.11	65.18	50.90

Table 6: Performance evaluation of KLD on classic horizontal detection.

Detector	Reg. Loss	AP	AP $_{50}$	AP $_{75}$	AP $_s$	AP $_m$	AP $_l$	Detector	Reg. Loss	AP	AP $_{50}$	AP $_{75}$	AP $_s$	AP $_m$	AP $_l$
RetinaNet	Smooth L1	37.2	56.6	39.7	21.4	41.1	41.1	Faster RCNN	Smooth L1	37.9	58.8	41.0	22.4	41.4	49.1
	GIoU	37.4	56.7	39.7	22.2	41.7	48.1		GIoU	38.3	58.7	41.5	22.5	41.7	49.7
	KLD	38.0	56.4	40.6	23.3	43.2	49.3		KLD	38.2	58.7	41.7	22.6	41.8	49.3

Table 7: AP on different objects on DOTA-v1.0. Here R-101 denotes ResNet-101 (likewise for R-50, R-152), and RX-101 and H-104 represent ResNeXt101 [46] and Hourglass-104 [47], respectively. MS indicates that multi-scale training/testing is used. Red and blue indicate the top two performances.

	Method	Backbone	MS	PL	BD	BR	GTF	SV	LV	SH	TC	BC	ST	SBF	RA	HA	SP	HC	AP $_{50}$
Two-stage	ICN [29]	R-101	✓	81.40	74.30	47.70	70.30	64.90	67.80	70.00	90.80	79.10	78.20	53.60	62.90	67.00	64.20	50.20	68.20
	RoI-Trans [11]	R-101	✓	88.64	78.52	43.44	75.92	68.81	73.68	83.59	90.74	77.27	81.46	58.39	53.54	62.83	58.93	47.67	69.56
	SCRDce [3]	R-101	✓	89.98	80.65	52.09	76.38	68.36	60.32	72.41	90.85	87.94	86.86	65.02	66.68	66.25	68.24	65.21	72.61
	Gilding Vertex [48]	R-101	✓	89.64	85.00	52.26	77.34	73.01	73.14	86.82	90.74	79.02	86.81	59.55	70.91	72.94	70.86	57.32	75.02
	Mask OBB [49]	RX-101	✓	89.56	85.95	54.21	72.90	76.52	74.16	85.63	89.85	83.81	86.48	54.89	69.64	73.94	69.06	63.32	75.33
	BBAVectors [50]	R-101	✓	89.83	84.41	54.60	70.25	77.66	78.32	87.19	90.66	84.89	85.27	56.46	69.23	74.13	71.56	66.02	76.03
	FPN-CSL [4]	R-152	✓	90.25	85.53	54.64	75.31	70.44	73.51	77.62	90.84	86.15	86.69	69.60	68.04	73.83	71.10	68.93	76.17
	RSdet-II [43]	R-152	✓	89.93	84.45	53.77	74.35	71.52	78.31	78.12	91.14	87.35	86.93	65.64	65.17	75.35	79.74	63.31	76.34
	SCRDet++ [51]	R-101	✓	90.05	84.39	55.44	73.99	77.54	71.11	86.05	90.67	87.32	87.08	69.62	68.90	73.74	71.29	65.08	76.81
	ReDet [52]	R-50	✓	88.81	82.48	60.83	80.82	78.34	86.06	88.31	90.87	88.77	87.03	68.65	66.90	79.26	79.71	74.67	80.10
Single-stage	PloU [30]	DLA-34 [53]	✓	80.90	69.70	24.10	60.20	38.30	64.40	64.80	90.30	77.20	70.40	46.50	37.10	61.9	64.00	60.50	
	O ² -DNet [54]	H-104	✓	89.31	82.14	47.33	61.21	71.32	74.03	78.62	90.76	82.23	81.36	60.93	60.17	58.21	66.98	61.03	71.04
	DAL [14]	R-101	✓	88.61	79.69	46.27	70.37	65.89	76.10	78.53	90.84	79.98	78.41	58.71	62.02	69.23	71.32	60.65	71.78
	P-RSDet [55]	R-101	✓	88.58	77.83	50.44	69.29	71.10	75.79	78.66	90.88	80.10	81.71	57.92	63.03	66.30	69.77	63.13	72.30
	BBAVectors [56]	R-101	✓	88.35	79.96	50.69	62.18	78.43	78.98	87.94	90.58	83.58	84.35	54.13	60.24	65.22	64.28	55.70	72.32
	DRN [13]	H-104	✓	89.71	82.34	47.22	64.10	76.22	74.43	85.84	90.57	86.18	84.89	57.65	61.93	69.30	69.63	58.48	73.23
	PolarDet [57]	R-101	✓	89.65	87.07	48.14	70.97	78.53	80.34	87.45	90.76	85.63	86.87	61.64	70.32	71.92	73.09	67.15	76.64
	RDD [58]	R-101	✓	89.15	83.92	52.51	73.06	77.81	79.00	87.08	90.62	86.72	87.15	63.96	70.29	76.98	75.79	72.15	77.75
	R-152	✓	89.06	84.32	55.33	77.53	76.95	70.28	83.95	89.75	84.51	86.96	73.47	67.77	72.60	75.76	74.17	77.43	
	KLD	R-50	✓	88.91	83.71	50.10	68.75	78.20	76.05	84.58	89.41	86.15	85.28	63.15	60.98	75.06	71.51	67.45	75.28
Refine-stage	CFC-Net [31]	R-101	✓	89.08	80.41	52.41	70.02	76.28	78.11	87.21	90.89	84.47	85.64	60.51	61.52	67.82	68.02	50.09	73.50
	R ³ Det [26]	R-152	✓	89.80	83.77	48.11	66.77	78.76	83.27	87.84	90.82	85.38	85.51	65.67	62.68	67.53	78.56	72.62	76.47
	DAL [14]	R-50	✓	89.69	83.11	55.03	71.00	78.30	81.90	88.46	90.89	84.97	87.46	64.41	65.65	76.86	72.09	64.35	76.95
	DCL [44]	R-152	✓	89.26	83.60	53.54	72.76	79.04	82.56	87.31	90.67	86.59	86.98	67.69	66.88	73.29	70.56	69.99	77.37
	RIDet [32]	R-50	✓	89.31	80.77	54.07	76.38	79.81	81.99	89.13	90.72	83.58	87.22	64.42	67.56	78.08	79.17	62.07	77.62
	S ² A-Net [12]	R-101	✓	89.28	84.11	56.95	79.21	80.18	82.93	89.21	90.86	84.66	87.61	71.66	68.23	78.58	78.20	65.55	79.15
	R ³ Det-GWD [5]	R-152	✓	89.66	84.99	59.26	82.19	78.97	84.83	87.70	90.21	86.54	86.85	73.04	67.56	76.92	79.22	74.92	80.19
	R ³ Det-KLD	R-50	✓	88.90	84.17	55.80	69.35	78.72	84.08	89.75	84.32	85.73	64.74	61.80	76.62	78.49	70.89	77.36	
	R-152	✓	89.90	84.91	59.21	78.74	78.82	83.95	87.41	89.89	86.63	86.69	70.47	70.87	76.96	79.40	78.62	80.17	
		R-152	✓	89.92	85.13	59.19	81.33	78.82	84.38	87.50	89.80	87.33	87.00	72.57	<b				

5 Discussions

Limitations. Despite the theoretical grounds and the promising experimental justifications, our method has an obvious limitation that it cannot be directly applied to quadrilateral detection [32, 43].

Potential negative societal impacts. Our findings provides a simple regression loss for high-precision rotation detection. However, our research may be applied to some sensitive fields, such as remote sensing, aviation, and unmanned aerial vehicles.

Conclusion. Departure from the vast existing literature in object detection, in this paper we have designed a new regression loss for rotation detection from scratch and consider the popular horizontal detection as its special case. Specifically, we calculate the KLD between the Gaussian distributions corresponding to the rotated bounding box as the regression loss, and we find that in the learning procedure guided by the KLD loss, the gradient of the parameters can be dynamically adjusted according to the characteristics of the object which is a desirable property for robust object detection, regardless its rotation, size and aspect ratio etc. We also proved that KLD has scale invariance, which is crucial for detection tasks. Interestingly, we have shown that KLD can be degenerated into the currently commonly used l_n -norm loss in the horizontal detection task. Extensive experimental results across different detectors and datasets show the effectiveness of our approach.

A Appendix

A.1 Proof of Scale Invariance of KLD

Suppose there are two Gaussian distributions, denoted as $\mathbf{X}_p \sim \mathcal{N}_p(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p)$ and $\mathbf{X}_t \sim \mathcal{N}_t(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$. Then, for a full-rank matrix \mathbf{M} , $|\mathbf{M}| \neq 0$, we have $\mathbf{X}_{p'} = \mathbf{M}\mathbf{X}_p \sim \mathcal{N}_p(\mathbf{M}\boldsymbol{\mu}_p, \mathbf{M}\boldsymbol{\Sigma}_p\mathbf{M}^\top)$, $\mathbf{X}_{t'} = \mathbf{M}\mathbf{X}_t \sim \mathcal{N}_t(\mathbf{M}\boldsymbol{\mu}_t, \mathbf{M}\boldsymbol{\Sigma}_t\mathbf{M}^\top)$, denoted as $\mathcal{N}_{p'}$ and $\mathcal{N}_{t'}$. The Kullback-Leibler Divergence (KLD) between $\mathcal{N}_{p'}$ and $\mathcal{N}_{t'}$ is:

$$\begin{aligned} \mathbf{D}_{kl}(\mathcal{N}_{p'} || \mathcal{N}_{t'}) &= \frac{1}{2}(\boldsymbol{\mu}_p - \boldsymbol{\mu}_t)^\top \mathbf{M}^\top (\mathbf{M}^\top)^{-1} \boldsymbol{\Sigma}_t^{-1} \mathbf{M}^{-1} \mathbf{M}(\boldsymbol{\mu}_p - \boldsymbol{\mu}_t) \\ &\quad + \frac{1}{2} \text{Tr} \left((\mathbf{M}^\top)^{-1} \boldsymbol{\Sigma}_t^{-1} \mathbf{M}^{-1} \mathbf{M} \boldsymbol{\Sigma}_p \mathbf{M}^\top \right) \\ &\quad + \frac{1}{2} \ln \frac{|\mathbf{M}| |\boldsymbol{\Sigma}_t| |\mathbf{M}^\top|}{|\mathbf{M}| |\boldsymbol{\Sigma}_p| |\mathbf{M}^\top|} - 1 \\ &= \frac{1}{2}(\boldsymbol{\mu}_p - \boldsymbol{\mu}_t)^\top \boldsymbol{\Sigma}_t^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_t) \\ &\quad + \frac{1}{2} \text{Tr} \left(\mathbf{M}^\top (\mathbf{M}^\top)^{-1} \boldsymbol{\Sigma}_t^{-1} \mathbf{M}^{-1} \mathbf{M} \boldsymbol{\Sigma}_p \right) \\ &\quad + \frac{1}{2} \ln \frac{|\boldsymbol{\Sigma}_t|}{|\boldsymbol{\Sigma}_p|} - 1 \\ &= \mathbf{D}_{kl}(\mathcal{N}_p || \mathcal{N}_t) \end{aligned} \tag{20}$$

Therefore, KLD has affine invariance. Especially when $\mathbf{M} = k\mathbf{I}$ (\mathbf{I} denotes identity matrix), the scale invariance of KLD is proved.

A.2 Analysis of $\mathbf{D}_{kl}(\mathcal{N}_t || \mathcal{N}_p)$'s High-Precision Detection

The $\mathbf{D}_{kl}(\mathcal{N}_t || \mathcal{N}_p)$ between two 2-D Gaussian is:

$$\mathbf{D}_{kl}(\mathcal{N}_t || \mathcal{N}_p) = \frac{1}{2}(\boldsymbol{\mu}_p - \boldsymbol{\mu}_t)^\top \boldsymbol{\Sigma}_p^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_t) + \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_p^{-1} \boldsymbol{\Sigma}_t) + \frac{1}{2} \ln \frac{|\boldsymbol{\Sigma}_p|}{|\boldsymbol{\Sigma}_t|} - 1 \tag{21}$$

each item of Eq. 21 can be expressed as

$$(\boldsymbol{\mu}_p - \boldsymbol{\mu}_t)^\top \boldsymbol{\Sigma}_p^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_t) = \frac{4(\Delta x \cos \theta_p + \Delta y \sin \theta_p)^2}{w_p^2} + \frac{4(\Delta y \cos \theta_p - \Delta x \sin \theta_p)^2}{h_p^2} \tag{22}$$

$$\text{Tr}(\boldsymbol{\Sigma}_p^{-1} \boldsymbol{\Sigma}_t) = \frac{h_t^2}{w_p^2} \sin^2 \Delta\theta + \frac{w_t^2}{h_p^2} \sin^2 \Delta\theta + \frac{h_t^2}{h_p^2} \cos^2 \Delta\theta + \frac{w_t^2}{w_p^2} \cos^2 \Delta\theta \tag{23}$$

$$\ln \frac{|\Sigma_p|}{|\Sigma_t|} = \ln \frac{h_p^2}{h_t^2} + \ln \frac{w_p^2}{w_t^2} \quad (24)$$

where $\Delta x = x_p - x_t$, $\Delta y = y_p - y_t$, $\Delta\theta = \theta_p - \theta_t$.

For the parameter μ_p , we have

$$\frac{\partial f_{kl}(\mu_p)}{\partial \mu_p} = \begin{pmatrix} \frac{4}{w_p^2} (\Delta x \cos \theta_p + \Delta y \sin \theta_p) \cos \theta_p + \frac{4}{h_p^2} (-\Delta x \sin \theta_p + \Delta y \cos \theta_p) (-\sin \theta_p) \\ \frac{4}{w_p^2} (\Delta x \cos \theta_p + \Delta y \sin \theta_p) \sin \theta_p + \frac{4}{h_p^2} (-\Delta x \sin \theta_p + \Delta y \cos \theta_p) \cos \theta_p \end{pmatrix} \quad (25)$$

It is assumed that except for μ_p , other parameters have been optimized to the best. In other words, $h_p = h_t$, $w_p = w_t$, and $\theta_p = \theta_t$. Without loss of generality, we set $\theta_t = 0^\circ$, then

$$\frac{\partial f_{kl}(\mu_p)}{\partial \mu_p} = \left(\frac{4}{w_t^2} \Delta x, \frac{4}{h_t^2} \Delta y \right)^\top \quad (26)$$

The weights $1/w_t^2$ and $1/h_t^2$ will make the model dynamically adjust the optimization of the object position according to the scale.

For h_p and w_p , we have

$$\begin{aligned} \frac{\partial f_{kl}(\Sigma_p)}{\partial \ln h_p} &= 1 - \frac{4(\Delta y \cos \theta_p - \Delta x \sin \theta_p)^2}{h_p^2} - \frac{w_t^2}{h_p^2} \sin^2 \Delta\theta - \frac{h_t^2}{h_p^2} \cos^2 \Delta\theta \\ \frac{\partial f_{kl}(\Sigma_p)}{\partial \ln w_p} &= 1 - \frac{4(\Delta x \cos \theta_p + \Delta y \sin \theta_p)^2}{w_p^2} - \frac{h_t^2}{w_p^2} \sin^2 \Delta\theta - \frac{w_t^2}{w_p^2} \cos^2 \Delta\theta \end{aligned} \quad (27)$$

Similarly, suppose $\Delta x = \Delta y = \Delta\theta = 0$, $\frac{\partial f_{kl}(\Sigma_p)}{\partial \ln h_p} = 1 - \frac{h_t^2}{h_p^2}$, $\frac{\partial f_{kl}(\Sigma_p)}{\partial \ln w_p} = 1 - \frac{w_p^2}{w_t^2}$, which means that the smaller targeted height or width leads to heavier penalty on its matching loss. This is desirable, as smaller height or width needs higher matching precision.

Similarly, suppose $\Delta x = \Delta y = \Delta\theta = 0$ and $w_p = w_t, h_p = h_t$, we have

$$\frac{\partial f_{kl}(\Sigma_p)}{\partial \theta_p} = \left(\frac{h_t^2}{w_t^2} + \frac{w_t^2}{h_t^2} - 2 \right) \sin 2\Delta\theta \geq \sin 2\Delta\theta \quad (28)$$

the condition for the equality sign is $h_t = w_t$. This shows that the larger the aspect ratio of the object, the model will pay more attention to the optimization of the angle.

Compared with $\mathbf{D}_{kl}(\mathcal{N}_p || \mathcal{N}_t)$, $\mathbf{D}_{kl}(\mathcal{N}_t || \mathcal{N}_p)$ has a similar gradient optimization strategy. The difference is that the relationship between the parameters of $\mathbf{D}_{kl}(\mathcal{N}_t || \mathcal{N}_p)$ is tighter.

References

- [1] X. Yang, H. Sun, K. Fu, J. Yang, X. Sun, M. Yan, and Z. Guo, “Automatic ship detection in remote sensing images from google earth of complex scenes based on multiscale rotation dense feature pyramid networks,” *Remote Sensing*, vol. 10, no. 1, p. 132, 2018.
- [2] X. Yang, H. Sun, X. Sun, M. Yan, Z. Guo, and K. Fu, “Position detection and direction prediction for arbitrary-oriented ships via multitask rotation region convolutional neural network,” *IEEE Access*, vol. 6, pp. 50 839–50 849, 2018.
- [3] 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 *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 8232–8241.
- [4] X. Yang and J. Yan, “Arbitrary-oriented object detection with circular smooth label,” in *Proceedings of the European Conference on Computer Vision*. Springer, 2020, pp. 677–694.
- [5] X. Yang, J. Yan, M. Qi, W. Wang, Z. Xiaopeng, and T. Qi, “Rethinking rotated object detection with gaussian wasserstein distance loss,” in *International Conference on Machine Learning*, 2021.
- [6] R. Girshick, “Fast r-cnn,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1440–1448.

- [7] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [8] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, “Feature pyramid networks for object detection,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2117–2125.
- [9] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal loss for dense object detection,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2980–2988.
- [10] J. Dai, Y. Li, K. He, and J. Sun, “R-fcn: Object detection via region-based fully convolutional networks,” in *Advances in neural information processing systems*, 2016, pp. 379–387.
- [11] J. Ding, N. Xue, Y. Long, G.-S. Xia, and Q. Lu, “Learning roi transformer for oriented object detection in aerial images,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2849–2858.
- [12] J. Han, J. Ding, J. Li, and G.-S. Xia, “Align deep features for oriented object detection,” *IEEE Transactions on Geoscience and Remote Sensing*, 2021.
- [13] X. Pan, Y. Ren, K. Sheng, W. Dong, H. Yuan, X. Guo, C. Ma, and C. Xu, “Dynamic refinement network for oriented and densely packed object detection,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 207–11 216.
- [14] Q. Ming, Z. Zhou, L. Miao, H. Zhang, and L. Li, “Dynamic anchor learning for arbitrary-oriented object detection,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- [15] S. Kullback and R. A. Leibler, “On information and sufficiency,” *The annals of mathematical statistics*, vol. 22, no. 1, pp. 79–86, 1951.
- [16] C. Villani, *Optimal transport: old and new.* Springer Science & Business Media, 2008, vol. 338.
- [17] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.
- [18] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, “Ssd: Single shot multibox detector,” in *European conference on computer vision.* Springer, 2016, pp. 21–37.
- [19] Z. Tian, C. Shen, H. Chen, and T. He, “Fcos: Fully convolutional one-stage object detection,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9627–9636.
- [20] X. Zhou, D. Wang, and P. Krähenbühl, “Objects as points,” *arXiv preprint arXiv:1904.07850*, 2019.
- [21] Z. Yang, S. Liu, H. Hu, L. Wang, and S. Lin, “Repoints: Point set representation for object detection,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9657–9666.
- [22] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-end object detection with transformers,” in *European Conference on Computer Vision.* Springer, 2020, pp. 213–229.
- [23] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, “Deformable detr: Deformable transformers for end-to-end object detection,” *arXiv preprint arXiv:2010.04159*, 2020.
- [24] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 658–666.
- [25] 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 *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020, pp. 12 993–13 000.
- [26] X. Yang, J. Yan, Z. Feng, and T. He, “R3det: Refined single-stage detector with feature refinement for rotating object,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- [27] Y. Jiang, X. Zhu, X. Wang, S. Yang, W. Li, H. Wang, P. Fu, and Z. Luo, “R2cnn: rotational region cnn for orientation robust scene text detection,” *arXiv preprint arXiv:1706.09579*, 2017.
- [28] 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.

- [29] S. M. Azimi, E. Vig, R. Bahmanyar, M. Körner, and P. Reinartz, “Towards multi-class object detection in unconstrained remote sensing imagery,” in *Asian Conference on Computer Vision*. Springer, 2018, pp. 150–165.
- [30] 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,” *Proceedings of the European Conference on Computer Vision*, 2020.
- [31] Q. Ming, L. Miao, Z. Zhou, and Y. Dong, “Cfc-net: A critical feature capturing network for arbitrary-oriented object detection in remote sensing images,” *arXiv preprint arXiv:2101.06849*, 2021.
- [32] Q. Ming, Z. Zhou, L. Miao, X. Yang, and Y. Dong, “Optimization for oriented object detection via representation invariance loss,” *arXiv preprint arXiv:2103.11636*, 2021.
- [33] C. Manning and H. Schütze, *Foundations of statistical natural language processing*. MIT press, 1999.
- [34] H. Jeffreys, “An invariant form for the prior probability in estimation problems,” *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, vol. 186, no. 1007, pp. 453–461, 1946.
- [35] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3974–3983.
- [36] H. Zhu, X. Chen, W. Dai, K. Fu, Q. Ye, and J. Jiao, “Orientation robust object detection in aerial images using deep convolutional neural network,” in *2015 IEEE International Conference on Image Processing*. IEEE, 2015, pp. 3735–3739.
- [37] Z. Liu, L. Yuan, L. Weng, and Y. Yang, “A high resolution optical satellite image dataset for ship recognition and some new baselines,” in *Proceedings of the International Conference on Pattern Recognition Applications and Methods*, vol. 2, 2017, pp. 324–331.
- [38] D. Karatzas, L. Gomez-Bigorda, A. Nicolaou, S. Ghosh, A. Bagdanov, M. Iwamura, J. Matas, L. Neumann, V. R. Chandrasekhar, S. Lu *et al.*, “Icdar 2015 competition on robust reading,” in *2015 13th International Conference on Document Analysis and Recognition*. IEEE, 2015, pp. 1156–1160.
- [39] N. Nayef, F. Yin, I. Bizid, H. Choi, Y. Feng, D. Karatzas, Z. Luo, U. Pal, C. Rigaud, J. Chazalon *et al.*, “Icdar2017 robust reading challenge on multi-lingual scene text detection and script identification-rrc-mlt,” in *2017 14th IAPR International Conference on Document Analysis and Recognition*, vol. 1. IEEE, 2017, pp. 1454–1459.
- [40] C. Yao, X. Bai, W. Liu, Y. Ma, and Z. Tu, “Detecting texts of arbitrary orientations in natural images,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 1083–1090.
- [41] 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.
- [42] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [43] W. Qian, X. Yang, S. Peng, J. Yan, and Y. Guo, “Learning modulated loss for rotated object detection,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- [44] X. Yang, L. Hou, Y. Zhou, W. Wang, and J. Yan, “Dense label encoding for boundary discontinuity free rotation detection,” in *Proceedings of the IEEE Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [45] 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 *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [46] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, “Aggregated residual transformations for deep neural networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1492–1500.
- [47] A. Newell, K. Yang, and J. Deng, “Stacked hourglass networks for human pose estimation,” in *Proceedings of the European Conference on Computer Vision*. Springer, 2016, pp. 483–499.
- [48] Y. Xu, M. Fu, Q. Wang, Y. Wang, K. Chen, G.-S. Xia, and X. Bai, “Gliding vertex on the horizontal bounding box for multi-oriented object detection,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

- [49] J. Wang, J. Ding, H. Guo, W. Cheng, T. Pan, and W. Yang, “Mask obb: A semantic attention-based mask oriented bounding box representation for multi-category object detection in aerial images,” *Remote Sensing*, vol. 11, no. 24, p. 2930, 2019.
- [50] J. Wang, W. Yang, H.-C. Li, H. Zhang, and G.-S. Xia, “Learning center probability map for detecting objects in aerial images,” *IEEE Transactions on Geoscience and Remote Sensing*, 2020.
- [51] X. Yang, J. Yan, X. Yang, J. Tang, W. Liao, and T. He, “Scrdet++: Detecting small, cluttered and rotated objects via instance-level feature denoising and rotation loss smoothing,” *arXiv preprint arXiv:2004.13316*, 2020.
- [52] J. Han, J. Ding, N. Xue, and G.-S. Xia, “Redet: A rotation-equivariant detector for aerial object detection,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021.
- [53] F. Yu, D. Wang, E. Shelhamer, and T. Darrell, “Deep layer aggregation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2403–2412.
- [54] H. Wei, Y. Zhang, Z. Chang, H. Li, H. Wang, and X. Sun, “Oriented objects as pairs of middle lines,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 169, pp. 268–279, 2020.
- [55] L. Zhou, H. Wei, H. Li, W. Zhao, Y. Zhang, and Y. Zhang, “Arbitrary-oriented object detection in remote sensing images based on polar coordinates,” *IEEE Access*, vol. 8, pp. 223 373–223 384, 2020.
- [56] J. Yi, P. Wu, B. Liu, Q. Huang, H. Qu, and D. Metaxas, “Oriented object detection in aerial images with box boundary-aware vectors,” *arXiv preprint arXiv:2008.07043*, 2020.
- [57] P. Zhao, Z. Qu, Y. Bu, W. Tan, Y. Ren, and S. Pu, “Polardet: A fast, more precise detector for rotated target in aerial images,” *arXiv preprint arXiv:2010.08720*, 2020.
- [58] B. Zhong and K. Ao, “Single-stage rotation-decoupled detector for oriented object,” *Remote Sensing*, vol. 12, no. 19, p. 3262, 2020.