

R³Det: Refined Single-Stage Detector with Feature Refinement for Rotating Object

Xue Yang^{1,4} Qingqing Liu², Junchi Yan¹, Ang Li^{3,4}, Zhiqiang Zhang⁴, Gang Yu⁴

¹Shanghai Jiao Tong University ²Center South University

³Nanjing University of Science and Technology ⁴Megvii Inc. (Face++)

yangxue-2019-sjtu@sjtu.edu.cn

Abstract

Rotation detection is a challenging task due to the difficulties of locating the multi-angle objects and separating them accurately and quickly from the background. Though considerable progress has been made, for practical settings, there still exist challenges for rotating objects with large aspect ratio, dense distribution and category extremely imbalance. In this paper, we propose an end-to-end refined single-stage rotation detector for fast and accurate positioning objects. Considering the shortcoming of feature misalignment in existing refined single-stage detector, we design a feature refinement module to improve detection performance by getting more accurate features. The key idea of feature refinement module is to re-encode the position information of the current refined bounding box to the corresponding feature points through feature interpolation to realize feature reconstruction and alignment. Extensive experiments on two remote sensing public datasets DOTA, HRSC2016 as well as scene text data ICDAR2015 show the state-of-the-art accuracy and speed of our detector. Code is available at https://github.com/Thinklab-SJTU/R3Det_Tensorflow.

1. Introduction

Object detection is one of the fundamental tasks in computer vision, and many high-performance general-purpose object detectors have been proposed. Current popular detection methods can be in general divided into two types: two-stage object detectors [12, 11, 33, 8, 24] and single-stage object detectors [27, 31, 25]. Two-stage methods have achieved promising results on various benchmarks, while the single-stage approach maintains faster detection speeds.

However, current general horizontal detectors have fundamental limitations for many practical applications. For instance, scene text detection and remote sensing object de-

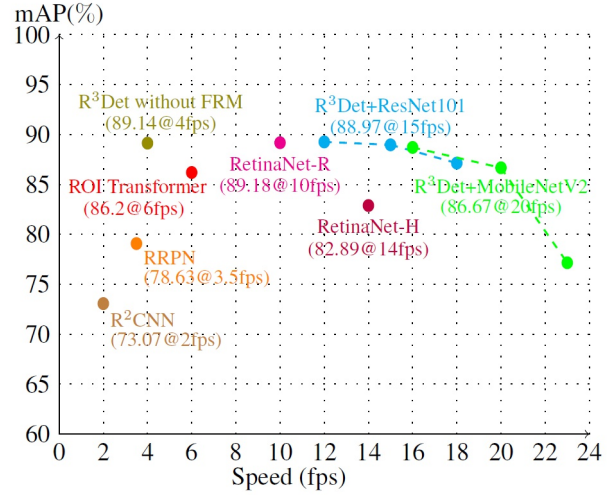


Figure 1: Performance (mAP) versus speed on HRSC2016 [29] dataset. Our detectors (in green and blue) notably surpass competitors in accuracy, whilst running very fast. Detailed results are listed in Table 5 (best viewed in color).

tection whereby the objects can be in any direction and position. Therefore, many rotation detectors based on a general detection framework have been proposed in the field of scene text and remote sensing. In particular, three challenges are pronounced for images in the above two fields, as analyzed as follows:

1) **Large aspect ratio.** The Skew Intersection over Union (SkewIoU) score between large aspect ratio objects is sensitive to change in angle, as sketched in Figure 3.

2) **Densely arranged.** As illustrated in Figure 6, Many objects usually appear in densely arranged forms.

3) **Category unbalance.** Many multi-category rotated datasets are long-tailed datasets whose categories are extremely unbalanced, as shown in Figure 7.

In this paper, we mainly discuss how to design an accu-

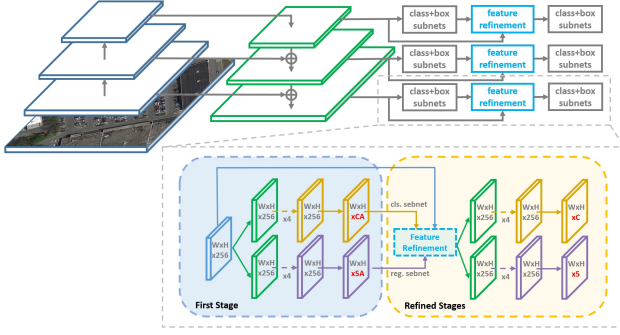


Figure 2: Architecture of the proposed Refined Rotation Single-Stage Detector (RetinaNet [25] as an embodiment). The refinement stage can be repeated multiple times. Only the bounding box with the highest score of each feature point is preserved in the refinement stage to speedup the model. ‘A’ indicates number of anchors on each feature point, and ‘C’ indicates number of categories.

rate and fast rotation detector. To maintain high positioning accuracy and detection speed for large aspect ratio objects, we have adopted a refined single-stage rotation detector. First, we find that rotating anchors can perform better in dense scenes, while horizontal anchors can achieve higher recalls in fewer quantities. Therefore, a combination strategy of two forms of anchors is adopted in the refinement single-stage detector, that is, the horizontal anchors are used in the first stage for faster speed and more proposals, and then the refined rotating anchors are used in the refinement stages to adapt to intensive scenarios. Second, we also notice that existing refined single-stage detectors have feature misalignment problems¹ [42, 7], which greatly limits the reliability of classification and regression during the refined stages. We design a feature refinement module (FRM) that uses the feature interpolation to obtain the position information of the refined anchors and reconstruct the feature map to achieve feature alignment. FRM can also reduce the number of refined bounding box after the first stage, thus speeding up the model. Experimental results have shown that feature refinement is sensitive to location and its improvement in detection results is very noticeable, especially for small sample categories. Combing these two techniques as a whole, our approach achieves state-of-the-art performance with high speed on three public rotating sensitive datasets including DOTA, HRSC2016 and ICDAR2015.

This work makes the following contributions:

- 1) For large aspect ratio objects, an accurate and fast rotation single-stage detector is devised in a refined manner, which enables detector for high-precision detection
- 2) For densely arranged scenes, we consider the advan-

¹Mainly refers to the misalignment between the region of interest (RoI) and the feature.

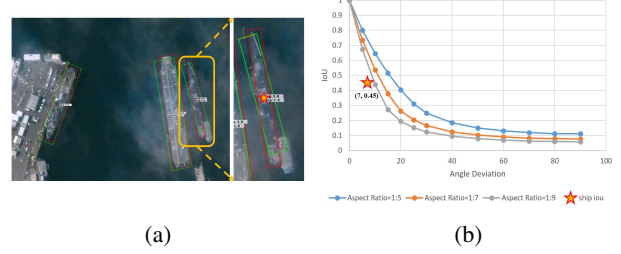


Figure 3: The SkewIoU scores vary with the angle deviation. The red and green rectangles represent the ground truth and the prediction bounding box, respectively.

tages of each of the two forms of anchors, and adopt an anchor combination strategy to enable the detector to cope with intensive scenarios with high efficiency.

3) For category unbalance, we propose an FRM that aims to make the detector features more accurate and reliable during the refinement stages. Experiments show that FRM has greatly improved the category that underfits due to the small number of samples and inaccurate features, such as BD, GTF, BC, SBF, RA, HC (see details in Table 1), which increased by 4.09%, 2.83%, 3.4%, 4.82%, 1.22%, and 19.26%, respectively.

2. Related Work

Two-Stage Object Detectors. Most of the existing two-stage methods are region-based. In a region based framework, category-independent region proposals are generated from an image in the first stage, features are extracted from these regions subsequently, and then category-specific classifiers and regressors are used for classification and regression in the second stage. Finally, the detection results are obtained by using post-processing methods such as non-maximum suppression (NMS). Faster-RCNN [33] is a classic structure in a two-stage approach that can detect object quickly and accurately in an end-to-end manner. Many high-performance detection methods are proposed today, such as R-FCN [8], FPN [24], etc.

Single-Stage Object Detectors. For their efficiency, single-stage detection methods are receiving more and more attention. OverFeat [35] is one of the first single-stage detectors based on convolutional neural networks. It performs object detection in a multiscale sliding window fashion via a single forward pass through the CNN. Compared with region based methods, Redmon et al. [31] propose YOLO, a unified detector casting object detection as a regression problem from image pixels to spatially separated bounding boxes and associated class probabilities. To preserve real-time speed without sacrificing too much detection accuracy, Liu et al. [27] propose SSD. The work [25] solves the class imbalance problem by proposing RetinaNet with Focal loss and further improves the accuracy of single-stage detector.

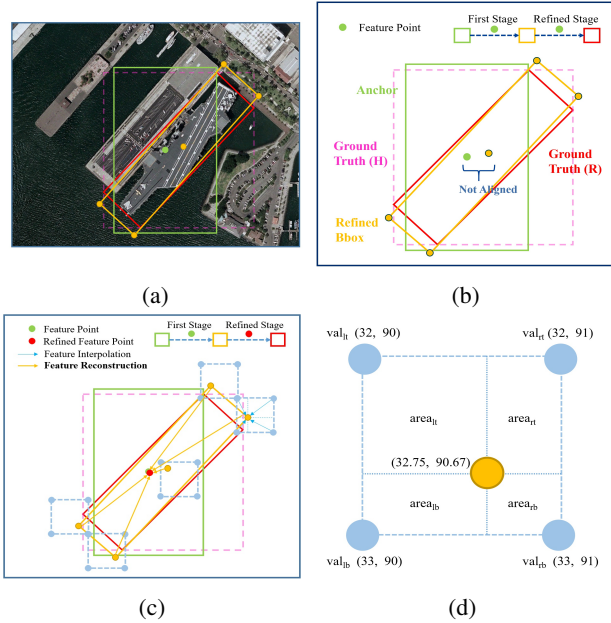


Figure 4: Root cause analysis of feature misalignment and the core idea for feature refinement module. (a) Original image. (b) Refine box without considering the feature misalignment caused by the location changes of the bounding box. (c) Refine box with aligned features by reconstructing the feature map. (d) Feature interpolation.

Rotation Object Detectors. Remote sensing and scene text are the main application scenarios of the rotation detector. Due to the complexity of the remote sensing image scene and the large number of small, cluttered and rotated objects, two-stage rotation detectors are still dominant for their robustness. Among them, ICN [2], ROI-Transformer [10] and SCRDet [40] are state-of-the-art detectors. However, they use a more complicated structure causing speed bottleneck. For scene text detection, there are many efficient rotation detection methods, including both two-stage methods (R²CNN [17], RRPN [30], FOTS [28]), as well as single-stage methods (EAST [44], TextBoxes [22]).

Refined Object Detectors. To achieve better positioning accuracy, many cascaded or refined detectors are proposed. The Cascade RCNN [4], HTC [5], and FSCascade [20] perform multiple classifications and regressions in the second stage, which greatly improved the classification accuracy and positioning accuracy. The same idea is also used in single-stage detectors, such as RefineDet [42]. Unlike the two-stage detectors, which use RoI Pooling [11] or RoI Align [13] for feature alignment. The currently refined single-stage detector is not well resolved in this respect. An important requirement of the refined single-stage detector is to maintain a full convolutional structure, which can retain the advantage of speed, but methods such as RoI Align cannot satisfy it whereby fully-connected layers have to be introduced. Although some works [6, 16, 41] use de-

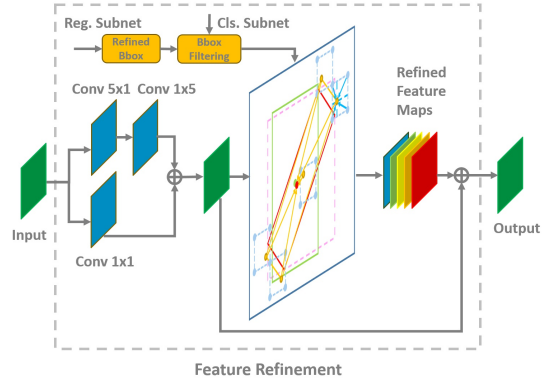


Figure 5: Feature Refinement Module (FRM). It mainly includes three parts: refined bounding box filtering, feature interpolation and feature map reconstruction.

formable convolution [9] for feature alignment, whose offset parameters are often obtained by learning the offset between the pre-defined anchor box and the refined anchor. The essence of these deformable-based feature alignment methods is to expand the receptive field, which is too implicit and can not ensure that features are truly aligned. Feature misalignment still limits the performance of the refined single-stage detector. Compared to these methods, our method can clearly find the corresponding feature area by calculation and achieve the purpose of feature alignment by feature map reconstruction.

3. The Proposed Method

We give an overview of our method as sketched in Figure 2. The embodiment is a single-stage rotation detector based on the RetinaNet [25], namely Refined Rotation RetinaNet (R³Det). The refinement stage (which can be added and repeated by multiple times) is added to the network to refine the bounding box, and the feature refinement module (FRM) is added during the refinement stage to reconstruct the feature map. In a single-stage rotating object detection task, continuous refinement of the predicted bounding box can improve the regression accuracy, and feature refinement is a necessary process for this purpose. It should be noted that FRM can also be used on other single-stage detectors (such as SSD), refer to the discussion section.

3.1. Rotation RetinaNet

RetinaNet is one of the most advanced single-stage detectors available today. It consists of two parts: backbone network, classification and regression subnetwork. RetinaNet adopts the Feature Pyramid Network (FPN) [24] as the backbone network. In brief, FPN augments a convolutional network with a top-down pathway and lateral connections so the network efficiently constructs a rich, multi-scale

feature pyramid from a single resolution input image. Each level of the pyramid can be used for detecting objects at a different scale. Besides, each layer of the FPN is connected to a classification subnet and a regression subnet for predicting categories and locations. Note that the object classification subnet and the box regression subnet, though sharing a common structure, use separate parameters. RetinaNet has proposed focal loss [25] to address the problem caused by category imbalance, which has greatly improved the accuracy of single-stage detector.

To achieve RetinaNet-based rotation detection, we use five parameters (x, y, w, h, θ) to represent arbitrary-oriented rectangle. Ranging in $[-\pi/2, 0)$, θ denotes the acute angle to the x-axis, and for the other side we refer it as w . Therefore, it calls for predicting an additional angular offset in the regression subnet, whose rotation bounding box is:

$$\begin{aligned} t_x &= (x - x_a)/w_a, t_y = (y - y_a)/h_a \\ t_w &= \log(w/w_a), t_h = \log(h/h_a), t_\theta = \theta - \theta_a \end{aligned} \quad (1)$$

$$\begin{aligned} t'_x &= (x' - x_a)/w_a, t'_y = (y' - y_a)/h_a \\ t'_w &= \log(w'/w_a), t'_h = \log(h'/h_a), t'_\theta = \theta' - \theta_a \end{aligned} \quad (2)$$

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

The multi-task loss is used which is defined as follows:

$$\begin{aligned} L &= \frac{\lambda_1}{N} \sum_{n=1}^N t'_n \sum_{j \in \{x, y, w, h, \theta\}} L_{reg}(v'_{nj}, v_{nj}) \\ &+ \frac{\lambda_2}{N} \sum_{n=1}^N L_{cls}(p_n, t_n) \end{aligned} \quad (3)$$

where N indicates the number of anchors, t'_n is a binary value ($t'_n = 1$ for foreground and $t'_n = 0$ for background, no regression for background). v'_{*j} represents the predicted offset vectors, v_{*j} represents the targets vector of ground-truth. t_n represents the label of object, p_n is the probability distribution of various classes calculated by sigmoid function. The hyper-parameter λ_1, λ_2 control the trade-off and are set to 1 by default. The classification loss L_{cls} and regression loss L_{reg} is implemented by focal loss [25] and smooth L1 loss as defined in [11], respectively.

3.2. Refined Rotation RetinaNet

Refined Detection. The Skew Intersection over Union (SkewIoU) score is sensitive to the change in angle, and a slight angle shift causes a rapid decrease in the IoU score, as shown in Figure 3. Therefore, the refinement of the prediction box helps to improve the recall rate of the rotation detection. We join multiple refinement stages with different IoU thresholds. In addition to using the foreground IoU

Algorithm 1 Feature Refinement Module

Input: original feature map F , the bounding box (B) and confidence (S) of the previous stage

Output: reconstructed feature map F'

```

1:  $B' \leftarrow \text{Filter}(B, S)$ ;
2:  $h, w \leftarrow \text{Shape}(F)$ ,  $F' \leftarrow \text{ZerosLike}(F)$ ;
3:  $F \leftarrow \text{Conv}_{1 \times 1}(F) + \text{Conv}_{1 \times 5}(\text{Conv}_{5 \times 1}(F))$ 
4: for  $i \leftarrow 0$  to  $h - 1$  do
5:   for  $j \leftarrow 0$  to  $w - 1$  do
6:      $P \leftarrow \text{GetFivePoints}(B'(i, j))$ ;
7:     for  $p \in P$  do
8:        $p_x \leftarrow \text{Min}(p_x, w - 1)$ ,  $p_x \leftarrow \text{Max}(p_x, 0)$ ;
9:        $p_y \leftarrow \text{Min}(p_y, h - 1)$ ,  $p_y \leftarrow \text{Max}(p_y, 0)$ ;
10:       $F'(i, j) \leftarrow F'(i, j) + \text{BilinearInte}(F, p)$ ;
11:   end for
12: end for
13: end for
14:  $F' \leftarrow F' + F$ ;
15: return  $F'$ 
```

threshold 0.5 and background IoU threshold 0.4 in the first stage, the thresholds of first refinement stage are set 0.6 and 0.5, respectively. If there are multiple refinement stages, the remaining thresholds are 0.7 and 0.6. The overall loss for refined detector is defined as follows:

$$L_{total} = \sum_{i=1}^N \alpha_i L_i \quad (4)$$

where L_i is the loss value of the i -th refinement stage and trade-off coefficients α_i are set to 1 by default.

Feature Refinement Module. Many refined detectors still use the same feature map to perform multiple classifications and regressions, without considering the feature misalignment caused by the location changes of the bounding box. Figure 4b depicts the box refining process without feature refinement, resulting in inaccurate features, which can be disadvantageous for those categories that have a large aspect ratio or a small sample size. Here we propose to re-encode the position information of the current refined bounding box (orange rectangle) to the corresponding feature points (red point²), thereby reconstructing the entire feature map to achieve the alignment of the features. The whole process is shown in Figure 4c. To accurately obtain the location feature information of the refined bounding box, we adopt the bilinear feature interpolation method, as shown in Figure 4d. Specifically, feature interpolation can

²The red and green points should be totally overlapping to each other, while here the red point is intentionally offset in order to distinguishingly visualize the entire process.

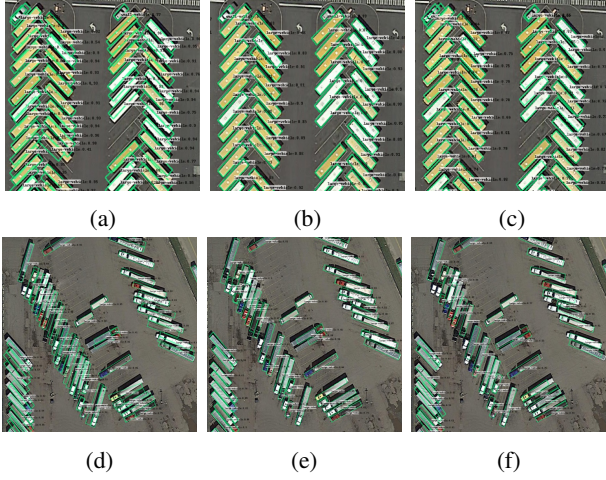


Figure 6: Visualization of three detectors on DOTA. (a)(d) RetinaNet-H. (b)(e) RetinaNet-R. (c)(f) R³Det without FRM.

be formulated as follows:

$$\begin{aligned} val = & val_{lt} * area_{rb} + val_{rt} * area_{lb} \\ & + val_{rb} * area_{lt} + val_{lb} * area_{rt} \end{aligned} \quad (5)$$

Based on the above result, a feature refinement module is devised, whose structure and pseudo code is shown in Figure 5 and Algorithm 1, respectively. Specifically, the feature map is added by two-way convolution to obtain a new feature. Only the bounding box with the highest score of each feature point is preserved in the refinement stage to increase the speed, meanwhile ensuring that each feature point corresponds to only one refined bounding box. For each feature point of the feature map, we obtain the corresponding feature vector on the feature map according to the five coordinates of the refined bounding box (one center point and four corner points). A more accurate feature vector is obtained by bilinear interpolation. We add the five feature vectors and replace the current feature vector. After traversing the feature points, we reconstruct the whole feature map. Finally, the reconstructed feature map is added to the original feature map to complete the whole process.

Discussion for Comparison with RoIAlign. The core to solve feature misalignment for FRM is feature reconstruction. Compared with RoIAlign [13] that has been adopted in many two-stage rotation detectors including R²CNN [17] and RRPN [30], FRM has the following differences that contribute to R³Det’s higher efficiency compared with R²CNN, RRPN as shown in Table 5:

1) RoI Align has more sampling points (the default number is $7 \times 7 \times 4 = 196$), and reducing the sampling point greatly affects the performance of the detector. FRM only samples five feature points, about one-fortieth of RoI Align, which gives FRM a huge speed advantage.

2) Before classification and regression, RoIAlign only

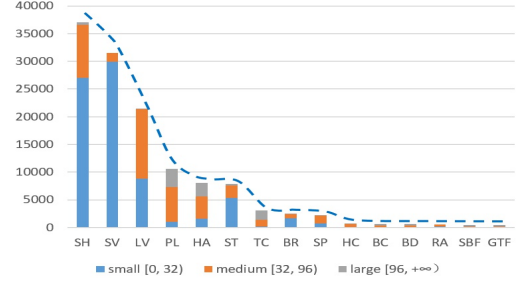


Figure 7: The quantity of each category in the DOTA.

need to obtain the feature corresponding to RoI (**instance level**). In contrast, FRM first obtains the features corresponding to the feature points, and then reconstructs the entire feature map (**image level**). As a result, the FRM based method can maintain a full convolution structure that leads to higher efficiency and fewer parameters, compared with the RoIAlign based method that involves a fully-connected structure.

4. Experiments

Tests are implemented by TensorFlow [1] on a server with GeForce RTX 2080 Ti and 11G memory. We perform experiments on both aerial benchmarks and scene text benchmarks to verify the generality of our techniques.

4.1. Datasets and Protocols

The benchmark DOTA [38] is for object detection in aerial images. It contains 2,806 aerial images from different sensors and platforms. The image size ranges from around 800×800 to $4,000 \times 4,000$ pixels and contains objects exhibiting a wide variety of scales, orientations, and shapes. These images are then annotated by experts using 15 common object categories. The fully annotated DOTA benchmark contains 188,282 instances, each of which is labeled by an arbitrary quadrilateral. There are two detection tasks for DOTA: horizontal bounding boxes (HBB) and oriented bounding boxes (OBB). Half of the original images are randomly selected as the training set, 1/6 as the validation set, and 1/3 as the testing set. We divide the images into 600×600 subimages with an overlap of 150 pixels and scale it to 800×800 . With all these processes, we obtain about 27,000 patches. The model is trained by 135k iterations in total, and the learning rate changes during the 81k and 108k iterations from $5e-4$ to $5e-6$.

The HRSC2016 dataset [29] contains images from two scenarios including ships on sea and ships close inshore. All the images are collected from six famous harbors. The image sizes range from 300×300 to $1,500 \times 900$. The training, validation and test set include 436 images, 181 images and 444 images, respectively. For all experiments we use an image scale of 800×800 for training and testing. We train

Method	Backbone	FRM		Data Aug.	PL	BD	BR	GTF	SV	LV	SH	TC	BC	ST	SBF	RA	HA	SP	HC	mAP
		large kernel	reconstruction step																	
RetinaNet-H (baseline)	ResNet50	×	×	×	88.87	74.46	40.11	58.03	63.10	50.61	63.63	90.89	77.91	76.38	48.26	55.85	50.67	60.23	34.23	62.22
RetinaNet-R (baseline)	ResNet50	×	×	×	88.92	67.67	33.55	56.83	66.11	73.28	75.24	90.87	73.95	75.07	43.77	56.72	51.05	55.86	21.46	62.02
R ³ Det (baseline)	ResNet50	✓	×	×	88.78	70.60	41.97	57.05	68.28	67.31	69.35	90.85	74.31	79.96	46.16	57.12	52.57	58.58	24.26	63.14
R ³ Det (proposed)	ResNet50	✓	✓	×	88.78	74.69	41.94	59.88	68.90	69.77	69.82	90.81	77.71	80.40	50.98	58.34	52.10	58.30	43.52	65.73
R ³ Det [†] (proposed)	ResNet50	✓	✓	×	89.01	75.24	43.74	60.57	68.52	74.03	75.55	90.84	79.02	75.37	53.02	57.12	54.23	62.95	48.86	67.20
R ³ Det (proposed)	ResNet50	✓	✓	✓	89.30	80.29	46.21	65.07	70.51	73.38	77.42	90.83	80.59	82.26	59.29	58.25	57.75	65.90	55.31	70.16
R ³ Det (proposed)	ResNet101	✓	✓	✓	89.54	81.99	48.46	62.52	70.48	74.29	77.54	90.80	81.39	83.54	61.97	59.82	65.44	67.46	60.05	71.69
R ³ Det (proposed)	ResNet152	✓	✓	✓	89.24	80.81	51.11	65.62	70.67	76.03	78.32	90.83	84.89	84.42	65.10	57.18	68.10	68.98	60.88	72.81
R ³ Det [†] (proposed)	ResNet152	✓	✓	✓	89.49	81.17	50.53	66.10	70.92	78.66	78.21	90.81	85.26	84.23	61.81	63.77	68.16	69.83	67.17	73.74

Table 1: Ablative study (AP for each category and overall mAP) of each components in our proposed method on the DOTA dataset. The short names for categories are defined as (abbreviation-full name): PL-Plane, BD-Baseball diamond, BR-Bridge, GTF-Ground field track, SV-Small vehicle, LV-Large vehicle, SH-Ship, TC-Tennis court, BC-Basketball court, ST-Storage tank, SBF-Soccer-ball field, RA-Roundabout, HA-Harbor, SP-Swimming pool, and HC-Helicopter. For RetinaNet, ‘H’ and ‘R’ represent the horizontal and rotating anchors, respectively. R³Det[†] indicates that two refinement stages have been added.

mAP	Feature Refinement Interpolation Formula	Feature Extraction
65.73	$val_{lt} * area_{rb} + val_{rt} * area_{lb} + val_{rb} * area_{lt} + val_{lb} * area_{rt}$	Bilinear Interpolation
64.28	$val_{lt} * area_{lt} + val_{rt} * area_{rt} + val_{rb} * area_{rb} + val_{lb} * area_{lb}$	Random Interpolation
64.37	$val_{lt} * area_{lb} + val_{rt} * area_{rb} + val_{rb} * area_{rt} + val_{lb} * area_{lt}$	Random interpolation
64.02	$val_{lt} * 1 + val_{rt} * 0 + val_{rb} * 0 + val_{lb} * 0$	Quantification
64.19	$val_{lt} * 0 + val_{rt} * 0 + val_{rb} * 1 + val_{lb} * 0$	Quantification

Table 2: Experiments with different interpolation formulas. Feature interpolation has position-sensitive properties.

#Stages	Test stage	SV	LV	SH	mAP
1	1	63.10	50.61	63.63	62.22
2	2	68.90	69.77	69.82	65.73
3	3	68.52	74.03	75.55	67.20
4	4	69.54	75.40	76.60	66.99
5	5	68.03	72.36	76.45	66.18
3	$\overline{2-3}$	69.14	74.53	75.98	67.68

Table 3: Ablation study for number of stages on DOTA dataset. $\overline{2-3}$ indicates the ensemble result, which is the collection of all outputs from refinement stages.

the model with 5e-4 learning rate for the first 30k iterations, then 5e-5 and 5e-6 for the other two 10k iterations.

ICDAR2015 is used in Challenge 4 of ICDAR 2015 Robust Reading Competition [18]. It includes a total of 1500 pictures, 1000 of which are used for training and the remaining are for testing. The text regions are annotated by 4 vertices of the quadrangle. We use its origin image size 720×1280 for training and testing. The ICDAR2015 dataset uses the same learning strategy and changes the learning rate size in 15k iterations, 20k iterations, and 25k iterations, respectively.

We experiment with ResNet-FPN and MobileNetv2-FPN [34] backbones. All backbones are pre-trained on ImageNet [19]. Weight decay and momentum are 0.0001 and 0.9, respectively. We employ MomentumOptimizer over 8 GPUs with a total of 8 images per minibatch (1 images per GPU). The anchors have areas of 32^2 to 512^2 on pyramid levels P3 to P7, respectively. At each pyramid level we use anchors at seven aspect ratios $\{1, 1/2, 2, 1/3, 3, 5, 1/5\}$ and three scales $\{2^0, 2^{1/3}, 2^{2/3}\}$. We also add six angles $\{-90^\circ, -75^\circ, -60^\circ, -45^\circ, -30^\circ, -15^\circ\}$ for rotating anchor-based method.

4.2. Robust Baseline Methods

To our knowledge there is no work exactly falling in our idea presented in the paper. While we still believe there are alternatives and thus we meanwhile devise a competitive detector to further verify the advantage of our proposed method. From the perspective of anchor, we analyze the effect of two forms of anchor on the speed and accuracy of the detection method, and finally construct a compromised robust baseline method.

The anchor setting is critical for region-based detection models. Both the horizontal anchor and the rotating anchor can achieve the purpose of rotation detection, but they have their own advantages and disadvantages. The advantage of a horizontal anchor is that it can use less anchor but match more positive samples by calculating the IoU with the horizontal circumscribing rectangle of the ground truth, but it introduces a large number of non-object or regions of other objects. For an object with a large aspect ratio, its prediction rotating bounding box tends to be inaccurate, as shown in Figure 6a. In contrast, in Figure 6b, the rotating anchor avoids the introduction of noise regions by adding angle parameters and has better detection performance in dense scenes. However, the number of anchors has multiplied, making the model less efficient.

The performance of the single-stage detection method based on two forms of anchor (RetinaNet-H and RetinaNet-R) on the DOTA data set OBB task is shown in Table 1. In general, they have similar overall mAP (62.22% versus 62.02%), while with their respective characteristics. The horizontal anchor-based approach clearly has an advantage in speed, while the rotating anchor-based method has better regression capabilities in dense object scenarios, such as small vehicle, large vehicle, and ship. To more effectively verify the validity of the feature refinement module, we also build a refined rotation detector, which does not refine the feature. Since the number of anchors does not decrease before and after the refinement stage, the number of original anchors determines the speed of the model. Taking into account the speed and accuracy, we adopt an anchor combination strategy. Specifically, we first use horizontal anchor to

Method	PL	BD	BR	GTF	SV	LV	SH	TC	BC	ST	SBF	RA	HA	SP	HC	mAP
Two-stage methods																
FR-O [38]	79.09	69.12	17.17	63.49	34.20	37.16	36.20	89.19	69.60	58.96	49.4	52.52	46.69	44.80	46.30	52.93
R-DFPN [39]	80.92	65.82	33.77	58.94	55.77	50.94	54.78	90.33	66.34	68.66	48.73	51.76	55.10	51.32	35.88	57.94
R ² CNN [17]	80.94	65.67	35.34	67.44	59.92	50.91	55.81	90.67	66.92	72.39	55.06	52.23	55.14	53.35	48.22	60.67
RRPN [30]	88.52	71.20	31.66	59.30	51.85	56.19	57.25	90.81	72.84	67.38	56.69	52.84	53.08	51.94	53.58	61.01
ICN [2]	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-Transformer [10]	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
SCRDet [40]	89.98	80.65	52.09	68.36	68.36	60.32	72.41	90.85	87.94	86.86	65.02	66.68	66.25	68.24	65.21	72.61
Single-stage methods																
RetinaNet-H+ResNet50 [25]	88.87	74.46	40.11	58.03	63.10	50.61	63.63	90.89	77.91	76.38	48.26	55.85	50.67	60.23	34.23	62.22
RetinaNet-R+ResNet50 [25]	88.92	67.67	33.55	56.83	66.11	73.28	75.24	90.87	73.95	75.07	43.77	56.72	51.05	55.86	21.46	62.02
R ³ Det+ResNet101	89.54	81.99	48.46	62.52	70.48	74.29	77.54	90.80	81.39	83.54	61.97	59.82	65.44	67.46	60.05	71.69
R ³ Det+ResNet152	89.24	80.81	51.11	65.62	70.67	76.03	78.32	90.83	84.89	84.42	65.10	57.18	68.10	68.98	60.88	72.81
R ³ Det [†] +ResNet152	89.49	81.17	50.53	66.10	70.92	78.66	78.21	90.81	85.26	84.23	61.81	63.77	68.16	69.83	67.17	73.74

Table 4: Detection accuracy on different objects (AP) and overall performance (mAP) evaluation on DOTA. R³Det[†] indicates that two refinement stages have been added.

Method	FRM	Backbone	Image Size	Data Aug.	mAP	Speed
R ² CNN [17]	—	ResNet101	800*800	×	73.07	2fps
RC1 & RC2 [29]	—	VGG16	—	—	75.7	<1fps
RRPN [30]	—	ResNet101	800*800	×	79.08	3.5fps
R ² PN [43]	—	VGG16	—	✓	79.6	<1fps
RetinaNet-H	—	ResNet101	800*800	✓	82.89	14fps
RRD [23]	—	VGG16	384*384	—	84.3	slow
RetinaNet-R	—	ResNet101	800*800	✓	89.18	10fps
RoI-Transformer [10]	—	ResNet101	512*800	×	86.20	6fps
R ³ Det (proposed)	×	ResNet101	800*800	✓	89.14	4fps
	✓	ResNet152	800*800	✓	89.33	10fps
	✓	ResNet101	300*300	✓	87.14	18fps
	✓	ResNet101	600*600	✓	88.97	15fps
	✓	ResNet101	800*800	✓	89.26	12fps
	✓	MobileNetV2	300*300	✓	77.16	23fps
	✓	MobileNetV2	600*600	✓	86.67	20fps
	✓	MobileNetV2	800*800	✓	88.71	16fps

Table 5: Accuracy and speed comparison on HRSC2016.

Method	FRM	Recall	Precision	F-measure	Res.	Device	FPS
CTPN [37]	—	51.56	74.22	60.85	—	—	—
SegLink [36]	—	76.80	73.10	75.00	—	—	—
RRPN [30]	—	82.17	73.23	77.44	—	—	<1
EAST [44]	—	78.33	83.27	80.72	720p	Titan X	13.2
Deep direct regression [14]	—	80.00	82.00	81.00	—	—	<1
R ² CNN [17]	—	79.68	85.62	82.54	720p	K80	0.44
FOTS RT [28]	—	85.95	79.83	82.78	720p	Titan X	24
FOTS [28]	—	91.0	85.17	87.99	1260p	Titan X	7.8
R ³ Det (proposed)	×	81.64	84.97	83.27	720p	2080 Ti	4
	✓	83.54	86.43	84.96	720p	2080 Ti	13.5

Table 6: Accuracy and speed comparison ICDAR2015.

reduce the number of anchors and increase the object recall rate, and then use the rotating refined anchor to overcome the problems caused by dense scenes, as shown in Figure 6c. The refined rotation detector achieves 63.14% performance, better than RetinaNet-H and RetinaNet-R.

4.3. Ablation Study

Feature Refinement Module. It shows that by removing FRM from R³Det, it can improve performance by about 1% which is not significant. We believe that the main reason is that the anchor is inconsistent with the feature map after the box refinement. FRM reconstructs the feature map based on the refined anchor, which increases the overall performance by 2.59% to 65.73% according to Table 1. We count the number of objects for each category, as shown in Figure 7. Coincidentally, FRM has greatly improved the categories that are underfitting due to the small number of samples and inaccurate features, such as BD, GTF, BC, SBF,

RA, HC, which increased by 4.09%, 2.83%, 3.4%, 4.82%, 1.22%, and 19.26%, respectively. In refined detectors, there are two reasons for the poor performance of small sample categories: lack of adequate training and inaccuracies in the refinement stage. The former can be mitigated by focal loss, while the latter can be solved by FRM.

Feature Refinement Interpolation Formula. When we randomly disturb the order of the four weights in the interpolation formula, the final performance of the model will be greatly reduced, rows 3-4 of Table 2. The same conclusions have also appeared in the experiments of quantitative operations, see rows 5-6 of Table 2. This phenomenon reflects the location sensitivity of the feature points and explains why the performance of the model can be greatly improved after the feature is correctly refined.

Number of Refinement Stages. We have known that adding a refinement stage has a significant improvement in rotation detection, especially the introduction of feature refinement. How about joining multiple refinements? R³Det[†] in Table 1 has joined the two refinement stages and bring more gain. To further explore the impact of the number of stages, several experimental results are summarized in Table 8. Experiments show that three or more refinements will not bring additional improvements to overall performance. Despite this, there are still significant improvements in the three large aspect ratio categories (SV, LV and SH). We also find that ensemble multi-stage results can further improve detection performance.

Data Augmentation and Backbone. By data augmentation, we improve the performance from 65.57% to 70.16% by random horizontal, vertical flipping, random graying, and random rotation. In addition, we also explore the gain of the backbone for the model. Under ResNet101 and ResNet152 as the backbone, we observe a reasonable improvement in table 1 (70.16% → 71.69% → 72.81%).

Ablative study for FRM using SSD. We also verify the portability of FRM on different data sets based on SSD, see Table 7 for detailed results. FRM brings 3% and 0.84% gain in DOTA and HRSC2016, respectively. This shows

Model	Backbone	FRM	DOTA	HRSC2016
SSD-H	VGG16	×	60.15	82.21
Refined SSD		×	62.79	88.48
Refined SSD		✓	65.79	89.32

Table 7: Ablative study for FRM in SSD.

the excellent generalization capability of FRM.

4.4. Comparison with the State-of-the-Art

The proposed R³Det with FRM is compared to state-of-the-art object detectors on three datasets: DOTA [10], HRSC2016 [29] and ICDAR2015 [18]. Our model outperforms all other models.

Results on DOTA. We compare our results with the state-of-the-arts in DOTA as depicted in Table 4. The results of DOTA reported here are obtained by submitting our predictions to the official DOTA evaluation server³. The existing two-stage detectors are still dominant in DOTA dataset research, and the latest two-stage detection methods, such as ICN, ROI Transformer, and SCRDet, have performed well. However, they all use complex model structures in exchange for performance improvements, which are extremely low in terms of detection efficiency. The single-stage detection method proposed in this paper achieves comparable performance with the most advanced two-stage method, while maintaining a fast detection speed. The speed analysis is detailed in Section 4.6.

Results on HRSC2016. The HRSC2016 contains lots of large aspect ratio ship instances with arbitrary orientation, which poses a huge challenge to the positioning accuracy of the detector. We use RRPN [30] and R²CNN [17] for comparative experiments, which are originally used for scene text detection. Experimental results show that these two methods under-perform in the remote sensing dataset, only 73.07% and 79.08% respectively. Although RoI Transformer [10] achieves 86.20% results without data augmentation, its detection speed is still not ideal, and only about 6fps without accounting for the post-processing operations. RetinNet-H, RetinaNet-R and R³Det without FRM are the three baseline models used in this paper. RetinaNet-R achieves the best detection results, around 89.14%, which is consistent with the performance of the ship category in the DOTA dataset. This further illustrates that the rotation-based approach has advantages in large aspect ratio target detection. Under ResNet101 backbone, our model achieves state-of-the-art performances.

Results on ICDAR2015. Scene text detection is also one of the main application scenarios for rotation detection. As shown in Table 6, our method achieves 84.96% while maintaining 13.5fps in the ICDAR2015 dataset, better than most mainstream algorithms except for FOTS, which adds a lot of extra training data and uses large test image. Al-

³<https://captain-whu.github.io/DOTA/>

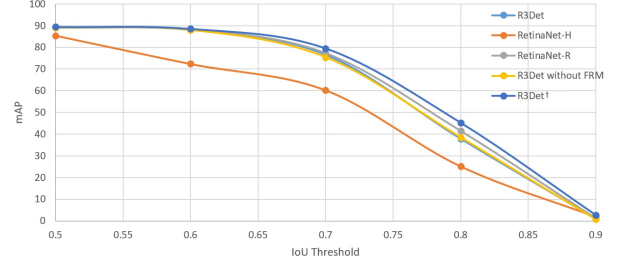


Figure 8: Detection performance (mAP) of the methods for different IoU thresholds on HRSC2016.

though the method in this paper is not a targeted single-class object detection, the experimental results still show that the proposed techniques are general that can be useful for both aerial images and scene text images.

4.5. Speed Comparison

Note we only add a small number of parameters in refinement stage to make the comparison as fair as possible. When we use the FRM, only the bounding box with the highest score of each feature point is preserved in the refinement stage to increase the speed of the model. We compare the speed and accuracy with the other six methods on the HRSC2016 dataset. The time of post process (i.e. R-NMS) is included. At the same time, we also explore the impact of different backbones and image sizes on the per performance of the proposed model. The detailed experimental results are shown in Table 5 and Figure 1. Our method can achieve 86.67% accuracy and 20fps speed, given MobileNetv2 as backbone with input image size 600×600 .

4.6. Discussion

High Precision Detection. Figure 8 shows the detection performance of the five methods for different IoU thresholds on HRSC2016. The rotating anchor-based method obtains a high-quality candidate bounding box, so it has excellent performance in high-precision detection. In contrast, the horizontal anchor-based method does not work well. We also find that the refinement method can achieve higher performance to the rotating anchor-based method.

FRM is More Suitable for Rotation Detection. When applying FRM to horizontal detection, the feature vectors of the four corner points obtained in FRM are likely to be far from the object, resulting in very inaccurate features being sampled. However, in the rotation detection task, the four corner points of the rotating bounding box are very close to the object. We have experimented on COCO dataset and the results are not satisfactory, but we have achieved considerable gains on many rotating datasets. See supplementary material for details.

5. Conclusion

We have presented an end-to-end refined single-stage detector designated for rotating objects with large aspect ratio, dense distribution and category extremely imbalance, which are common in practice especially for aerial and scene text image. Considering the shortcoming of feature misalignment in the current refined single-stage detector, we design a feature refinement module to improve detection performance, which is especially effective in the long tail data set. The key idea of FRM is to re-encode the position information of the current refined bounding box to the corresponding feature points through feature interpolation to achieve feature reconstruction and alignment. We perform careful ablation studies and comparative experiments on multiple rotation detection datasets including DOTA, HRSC2016, and ICDAR2015, and demonstrate that our method achieves state-of-the-art detection accuracy with high efficiency.

6. Appendix

6.1. Two different rotation box definitions

Figure 9a shows the rectangular definition of the 90 degree angle representation range. θ denotes the acute angle to the x-axis, and for the other side we refer it as w . It should be distinguished from another definition in Figure 9b, with 180 degree angular range, whose θ is determined by the long side of the rectangle and x-axis. This article uses the first definition, which is also officially used by OpenCV.

6.2. FRM is Suitable for Rotation Detection

When applying FRM to horizontal detection, the feature vectors of the four corner points (green points in Figure 10a) obtained in FRM are likely to be far from the object, resulting in very inaccurate features being sampled. However, in the rotation detection task, the four corner points (red points in Figure 10b) of the rotating bounding box are very close to the object. We have experimented on COCO dataset and the results are not satisfactory, but we have achieved considerable gains on many rotating datasets.

6.3. Performance evaluation on UCAS-AOD

UCAS-AOD [45] contains 1510 aerial images of approximately 659×1280 pixels, it contains two categories of 14596 instances. As [38, 2] did before, we randomly select 1110 for training and 400 for testing. The model is trained by 20 epoches in total with warm up strategy, and the number of iterations per epoch depends on the total amount of the data set. The initial learning rate is $5e-4$, and the learning rate changes during 12 and 16 epoches from $5e-5$ to $5e-6$. Table 9 illustrates the comparison of performance on UCAS-AOD dataset, we get 96.95% for OBB task. Our results are the best out of all the existing published methods.

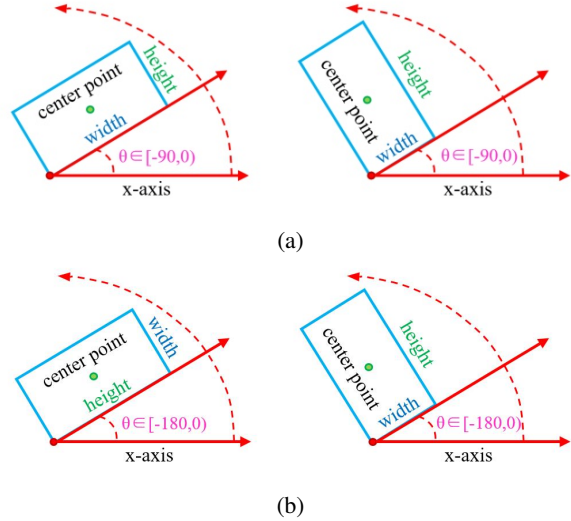


Figure 9: Two different rotation box definitions.

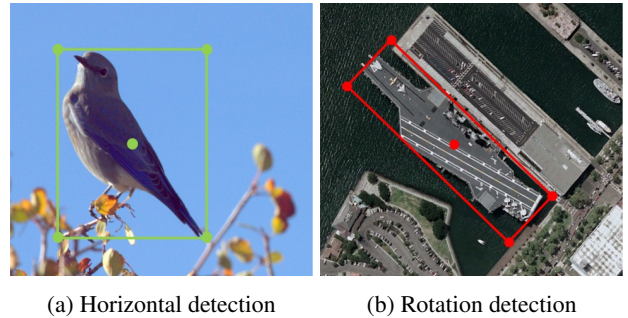


Figure 10: Schematic diagram of sampling points for FRM in horizontal detection and rotation detection. The sample points for horizontal detection are significantly further away from the object, and the sampling points for rotation detection are tighter.

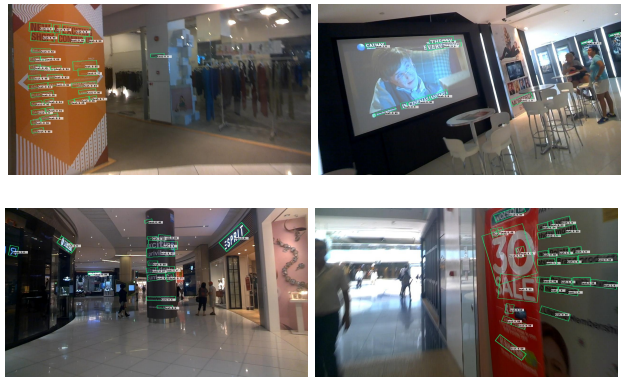


Figure 11: Text detection results on the ICDAR2015 benchmarks.

#Stages	Test stage	PL	BD	BR	GTF	SV	LV	SH	TC	BC	ST	SBF	RA	HA	SP	HC	mAP
1	1	88.87	74.46	40.11	58.03	63.10	50.61	63.63	90.89	77.91	76.38	48.26	55.85	50.67	60.23	34.23	62.22
2	2	88.78	74.69	41.94	59.88	68.90	69.77	69.82	90.81	77.71	80.40	50.98	58.34	52.10	58.30	43.52	65.73
3	3	89.01	75.24	43.74	60.57	68.52	74.03	75.55	90.84	79.02	75.37	53.02	57.12	54.23	62.95	48.86	67.20
4	4	88.87	74.38	43.40	61.34	68.55	75.40	76.60	90.87	80.03	75.73	53.01	60.49	53.62	60.44	41.18	66.99
5	5	88.98	72.30	41.70	62.11	68.03	72.36	76.45	90.84	76.11	79.20	51.01	57.56	51.10	59.71	45.22	66.18
3	$\overline{2-3}$	88.96	73.89	44.27	60.81	69.14	74.53	75.98	90.83	79.83	78.03	54.75	56.52	54.80	64.66	48.17	67.68

Table 8: Ablation study for number of stages on DOTA dataset. $\overline{2-3}$ indicates the ensemble result, which is the collection of all outputs from refinement stages.

Method	mAP	Plane	Car
YOLOv2 [32]	87.90	96.60	79.20
R-DFPN [39]	89.20	95.90	82.50
DRBox [26]	89.95	94.90	85.00
S ² ARN [3]	94.90	97.60	92.20
RetinaNet-H	95.47	97.34	93.60
ICN [2]	95.67	-	-
FADet [21]	95.71	98.69	92.72
R ³ Det	96.17	98.20	94.14

Table 9: Performance evaluation on UCAS-AOD datasets.

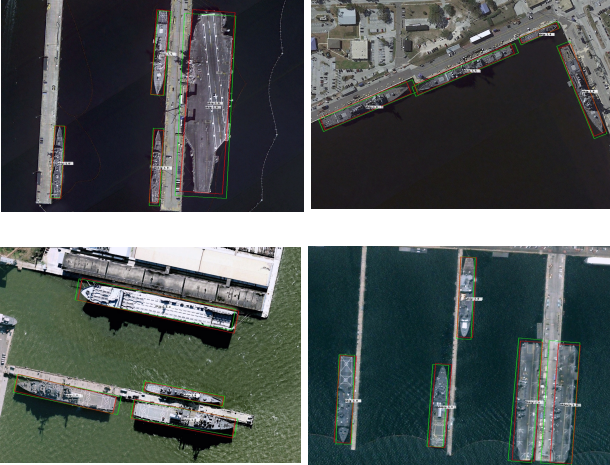


Figure 12: Ship detection results on the HRSC2016 benchmarks. The red and green bounding box indicate the ground truth and prediction box, respectively.

6.4. Ablation study for anchor Setting

We also briefly explore the impact of the number of anchors on the accuracy of detection. Table 10 shows that although more anchors will bring more gains on small backbone, these increases are not obvious when using large backbone and data augmentation. Therefore, considering the trade-off of speed and precision, the scale and ratio of the anchor in this paper are set to $\{2^0, 2^{1/3}, 2^{2/3}\}$ and $\{1, 1/2, 2, 1/3, 3, 5, 1/5\}$, respectively.



Figure 13: Face detection results on the Fddb [15] benchmarks.

Base Model	Backbone	Anchor Scales	Anchor Ratios	mAP
RetinaNet-H	ResNet50	$\{2^0\}$	$\{1, 1/2, 2, 1/3, 3, 5, 1/5\}$	59.48
	ResNet50	$\{2^0, 2^{1/3}, 2^{2/3}\}$	$\{1, 1/2, 2, 1/3, 3, 5, 1/5\}$	62.22
R ³ Det	ResNet50	$\{2^0, 2^{1/3}, 2^{2/3}\}$	$\{1, 1/2, 2, 1/3, 3, 5, 1/5\}$	65.73
	ResNet50	$\{0.5, 2^0, 2^{1/3}, 2^{2/3}\}$	$\{1, 1/2, 2, 1/3, 3, 5, 1/5, 7, 1/7, 9, 1/9\}$	67.54
	ResNet152	$\{2^0, 2^{1/3}, 2^{2/3}\}$	$\{1, 1/2, 2, 1/3, 3, 5, 1/5\}$	72.81
	ResNet152	$\{0.5, 2^0, 2^{1/3}, 2^{2/3}\}$	$\{1, 1/2, 2, 1/3, 3, 5, 1/5, 7, 1/7, 9, 1/9\}$	72.87

Table 10: Ablation study for anchor scale and ratio on DOTA dataset.

6.5. Number of Refinement Stages

Table 8 shows in detail that three or more refinements will not bring additional improvements to overall performance. Despite this, there are still significant improvements in the three large aspect ratio categories (SV, LV and SH). We also find that ensemble multi-stage results can further improve detection performance.

6.6. Visualization on Different Datasets

We visualize the detection results of R³Det on different types of data sets, including remote sensing datasets (Figure 14 and Figure 12), scene text datasets (Figure 11), and face datasets (Figure 13).

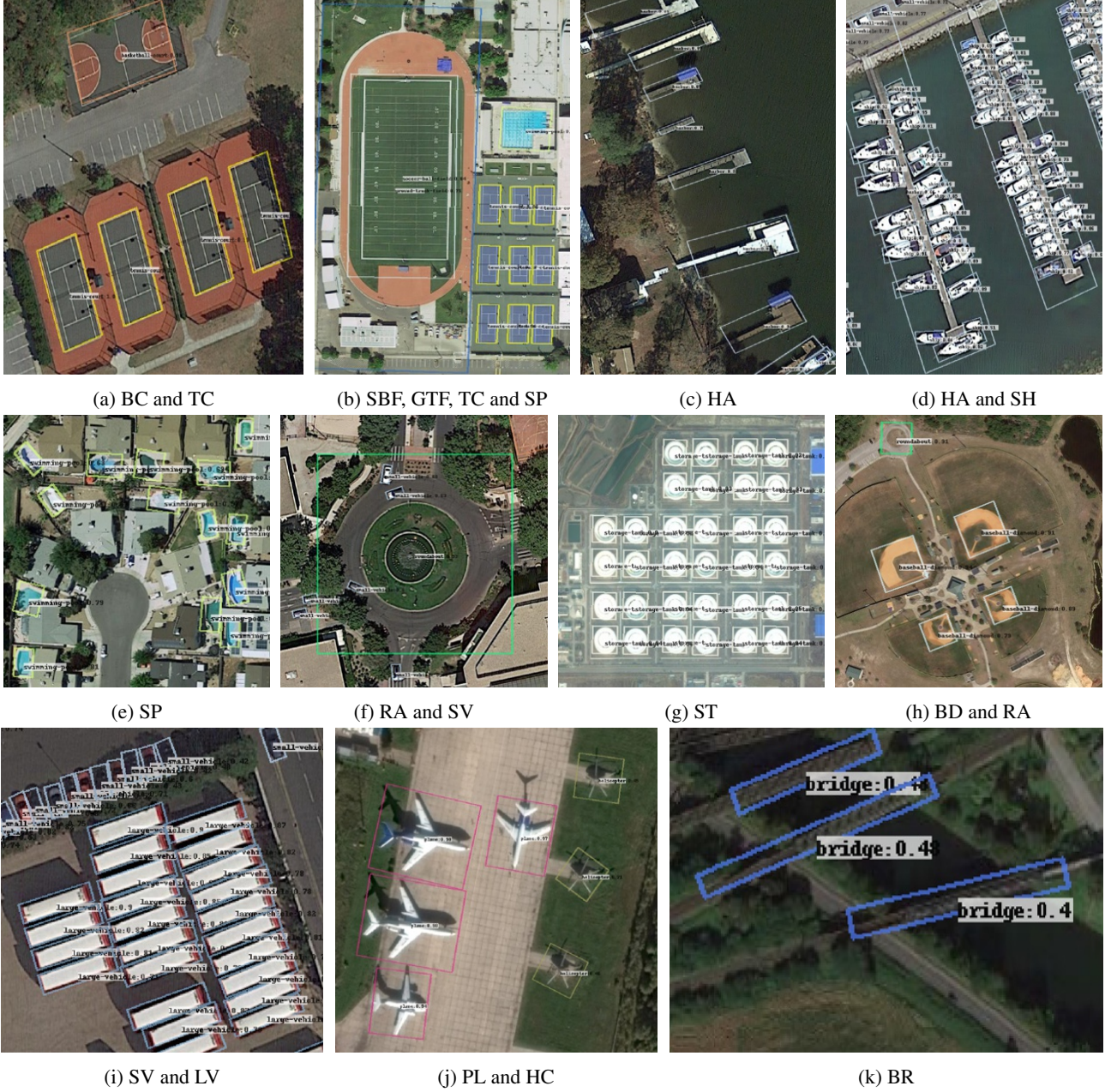


Figure 14: Detection results on the OBB task on DOTA. Our method performs better on those with large aspect ratio , in arbitrary direction, and high density.

References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: a system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283, 2016.
- [2] Seyed Majid Azimi, Eleonora Vig, Reza Bahmanyar, Marco Körner, and Peter Reinartz. Towards multi-class object detection in unconstrained remote sensing imagery. In *Asian Conference on Computer Vision*, pages 150–165. Springer, 2018.
- [3] Songze Bao, Xing Zhong, Ruifei Zhu, Xiaonan Zhang, Zhuqiang Li, and Mengyang Li. Single shot anchor refinement network for oriented object detection in optical remote sensing imagery. *IEEE Access*, 7:87150–87161, 2019.
- [4] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162, 2018.

- [5] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, et al. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4974–4983, 2019.
- [6] Xingyu Chen, Junzhi Yu, Shihan Kong, Zhengxing Wu, and Li Wen. Dual refinement networks for accurate and fast object detection in real-world scenes. *arXiv preprint arXiv:1807.08638*, 2018.
- [7] Cheng Chi, Shifeng Zhang, Junliang Xing, Zhen Lei, Stan Z Li, and Xudong Zou. Selective refinement network for high performance face detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8231–8238, 2019.
- [8] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in neural information processing systems*, pages 379–387, 2016.
- [9] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 764–773, 2017.
- [10] Jian Ding, Nan Xue, Yang Long, Gui-Song Xia, and Qikai Lu. Learning roi transformer for oriented object detection in aerial images. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [11] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.
- [12] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014.
- [13] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [14] Wenhao He, Xu Yao Zhang, Fei Yin, and Cheng Lin Liu. Deep direct regression for multi-oriented scene text detection. 2017.
- [15] Vedit Jain and Erik Learned-Miller. Fddb: A benchmark for face detection in unconstrained settings. 2010.
- [16] Ho-Deok Jang, Sanghyun Woo, Philipp Benz, Jinsun Park, and In So Kweon. Propose-and-attend single shot detector. *arXiv preprint arXiv:1907.12736*, 2019.
- [17] Yingying Jiang, Xiangyu Zhu, Xiaobing Wang, Shuli Yang, Wei Li, Hua Wang, Pei Fu, and Zhenbo Luo. R2cnn: rotational region cnn for orientation robust scene text detection. *arXiv preprint arXiv:1706.09579*, 2017.
- [18] Dimosthenis Karatzas, Lluís Gomez-Bigorda, Angelos Nicolaou, Suman Ghosh, Andrew Bagdanov, Masakazu Iwamura, Jiri Matas, Lukas Neumann, Vijay Ramaseshan Chandrasekhar, Shijian Lu, et al. Icdar 2015 competition on robust reading. In *2015 13th International Conference on Document Analysis and Recognition (ICDAR)*, pages 1156–1160. IEEE, 2015.
- [19] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [20] Ang Li, Xue Yang, and Chongyang Zhang. Rethinking classification and localization for cascade r-cnn. *arXiv preprint arXiv:1907.11914*, 2019.
- [21] Chengzheng Li, Chunyan Xu, Zhen Cui, Dan Wang, Tong Zhang, and Jian Yang. Feature-attended object detection in remote sensing imagery. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 3886–3890. IEEE, 2019.
- [22] Minghui Liao, Baoguang Shi, Xiang Bai, Xinggang Wang, and Wenyu Liu. Textboxes: A fast text detector with a single deep neural network. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [23] Minghui Liao, Zhen Zhu, Baoguang Shi, Gui-song Xia, and Xiang Bai. Rotation-sensitive regression for oriented scene text detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5909–5918, 2018.
- [24] Tsung-Yi Lin, Piotr Dollár, Ross B Girshick, Kaiming He, Bharath Hariharan, and Serge J Belongie. Feature pyramid networks for object detection. In *CVPR*, volume 1, page 4, 2017.
- [25] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [26] Lei Liu, Zongxu Pan, and Bin Lei. Learning a rotation invariant detector with rotatable bounding box. *arXiv preprint arXiv:1711.09405*, 2017.
- [27] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.
- [28] Xuebo Liu, Ding Liang, Shi Yan, Dagui Chen, Yu Qiao, and Junjie Yan. Fots: Fast oriented text spotting with a unified network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5676–5685, 2018.
- [29] Zikun Liu, Liu Yuan, Lubin Weng, and Yiping Yang. A high resolution optical satellite image dataset for ship recognition and some new baselines. In *Proc. ICPRAM*, volume 2, pages 324–331, 2017.
- [30] Jianqi Ma, Weiyuan Shao, Hao Ye, Li Wang, Hong Wang, Yingbin Zheng, and Xiangyang Xue. Arbitrary-oriented scene text detection via rotation proposals. *IEEE Transactions on Multimedia*, 2018.
- [31] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [32] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7263–7271, 2017.
- [33] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: towards real-time object detection with region

- proposal networks. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (6):1137–1149, 2017.
- [34] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang Chieh Chen. Inverted residuals and linear bottlenecks: Mobile networks for classification, detection and segmentation. 2018.
 - [35] Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*, 2013.
 - [36] Baoguang Shi, Xiang Bai, and Serge Belongie. Detecting oriented text in natural images by linking segments. 2017.
 - [37] Zhi Tian, Weilin Huang, Tong He, Pan He, and Yu Qiao. Detecting text in natural image with connectionist text proposal network. In *European conference on computer vision*, pages 56–72. Springer, 2016.
 - [38] Gui-Song Xia, Xiang Bai, Jian Ding, Zhen Zhu, Serge Belongie, Jiebo Luo, Mihai Datcu, Marcello Pelillo, and Liangpei Zhang. Dota: A large-scale dataset for object detection in aerial images. In *Proc. CVPR*, 2018.
 - [39] Xue Yang, Hao Sun, Kun Fu, Jirui Yang, Xian Sun, Menglong Yan, and Zhi Guo. Automatic ship detection in remote sensing images from google earth of complex scenes based on multiscale rotation dense feature pyramid networks. *Remote Sensing*, 10(1):132, 2018.
 - [40] Xue Yang, Jirui Yang, Junchi Yan, Yue Zhang, Tengfei Zhang, Zhi Guo, Xian Sun, and Kun Fu. Scrdet: Towards more robust detection for small, cluttered and rotated objects. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
 - [41] Hongkai Zhang, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen. Cascade retinanet: Maintaining consistency for single-stage object detection. *arXiv preprint arXiv:1907.06881*, 2019.
 - [42] Shifeng Zhang, Longyin Wen, Xiao Bian, Zhen Lei, and Stan Z. Li. Single-shot refinement neural network for object detection. 2018.
 - [43] Zenghui Zhang, Weiwei Guo, Shengnan Zhu, and Wenxian Yu. Toward arbitrary-oriented ship detection with rotated region proposal and discrimination networks. *IEEE Geoscience and Remote Sensing Letters*, 15(11):1745–1749, 2018.
 - [44] Xinyu Zhou, Cong Yao, He Wen, Yuzhi Wang, Shuchang Zhou, Weiran He, and Jiajun Liang. East: An efficient and accurate scene text detector. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017.
 - [45] Haigang Zhu, Xiaogang Chen, Weiqun Dai, Kun Fu, Qixiang Ye, and Jianbin Jiao. Orientation robust object detection in aerial images using deep convolutional neural network. In *2015 IEEE International Conference on Image Processing (ICIP)*, pages 3735–3739. IEEE, 2015.