# Adaptive Period Embedding for Representing Oriented Objects in Aerial Images

Yixing Zhu, Xueqing Wu, Jun Du

National Engineering Laboratory for Speech and Language Information Processing

University of Science and Technology of China

Hefei, Anhui, China

zyxsa@mail.ustc.edu.cn, shirley0@mail.ustc.edu.cn, jundu@ustc.edu.cn

*Abstract*—We propose a novel method for representing oriented objects in aerial images named Adaptive Period Embedding (APE). While traditional object detection methods represent object with horizontal bounding boxes, the objects in aerial images are oriented. Calculating the angle of object is an yet challenging task. While almost all previous object detectors for aerial images directly regress the angle of objects, they use complex rules to calculate the angle, and their performance is limited by the rule design. In contrast, our method is based on the angular periodicity of oriented objects. The angle is represented by two two-dimensional periodic vectors whose periods are different, the vector is continuous as shape changes. The label generation rule is more simple and reasonable compared with previous methods. The proposed method is general and can be applied to other oriented detector. Besides, we propose a novel IoU calculation method for long objects named length independent IoU (LIIoU). We intercept part of the long side of the target box to get the maximum IoU between the proposed box and the intercepted target box. Thereby, some long boxes will have corresponding positive samples. Our method reaches the $1^{st}$ place of DOAI2019 competition task1 (oriented object) held in workshop on Detecting Objects in Aerial Images in conjunction with IEEE CVPR 2019.

*Index Terms*—Oriented object detection, aerial images, IoU.

Fig. 1: Left: horizonal bounding boxes, right: oriented bounding boxes.

## I. INTRODUCTION

Traditional object detections mainly detect objects with horizontal bounding boxes. However, objects in aerial images are oriented and cannot be effectively represented by horizontal bounding boxes. As shown in Fig. 1, detecting oriented objects with horizontal bounding boxes will contain more background and cannot accurately locate the objects. Besides, overlap calculation based on horizontal bounding box is not accurate for oriented objects, as the overlap between horizontal bounding boxes of two oriented objects may be too large, so NMS based on horizontal bounding boxes is not suitable for oriented objects. Thus, representing oriented objects with oriented bounding box is necessary for object detection in aerial images. Regressing oriented bounding box is more challenging than regressing horizontal bounding box. Four variables can represent a horizontal bounding box, such as x,y coordinates of top left corner and bottom right corner. However, oriented bounding box representation needs an extra variable $\theta$ to represent its angle. It is hard to directly regress $\theta$ because the angle is periodic.

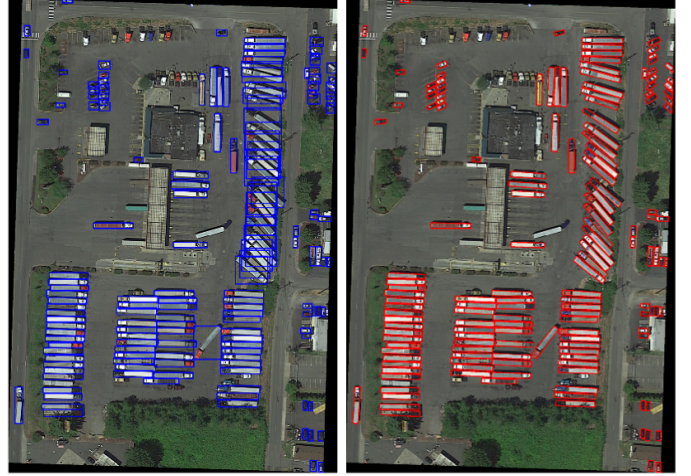Most of previous oriented detectors [1] [2] [3] [4] [5] directly regress $\theta$ or the four vertices of the oriented bounding box and the label is calculated by complicated rules, which is hard for the network to learn. Some methods try to design a simple label calculation rule for oriented objects. For example, [6] adopts Mask R-CNN [7] for detecting oriented text lines. [8] regresses the outline of objects with multiple points on sliding lines. But these methods introduce additional parameters and cannot be adopted by RPN (region proposal networks).

In this study, we propose a novel method for representing oriented objects. Oriented bounding box can be represented by $(x, y, w, h, \theta)$, where $x$, $y$ are the coordinates of the center of the bounding box, and $w$, $h$ are the lengths of the long and short sides, respectively. We do not directly regress $\theta$. Instead, $\theta$ is represented by two two-dimensional periodic vector. The proposed method is different from [9], as in our method, the two vectors' periods are $90°$ and $180°$ respectively. Finally, we calculate the angle with these vectors. Our method is versatile and can be applied in other detectors. Besides, we design a novel cascade R-CNN method for long objects such as harbors. Generally, a two-stage model firstly proposes bounding boxes with RPN, and the output bounding boxes of the second step (R-CNN) are limited by RPN's results. Due to the limited size of the receptive field, some long objects cannot be covered by RPN. With this in mind, we adopt a two-stage cascade R-CNN model with length independent IoU (LIIoU) to detect

long objects. In the first stage, some bounding boxes which only cover part of the objects are also set to positive samples. In this way, the first R-CNN can propose longer bounding boxes. The main contributions of our work are summarized as follows:

1. We present a novel method for representing oriented bounding boxes in aerial images. We do not directly regress $\theta$ of oriented bounding boxes, but instead embed $\theta$ with vectors whose periods are different. In this way, we do not need complex rules to label the angle which avoids ambiguity.

2. We present a novel IoU calculation method named length independent IoU (LIIoU), which is designed for long objects. The presented method makes the detector more robust to long objects.

3. The presented method achieves state-of-the-art on DOTA and wins the first place of Challenge-2019 on Object Detection in Aerial Images task1 (oriented task).

## II. RELATED WORK

### A. Horizontal objects detection

Labels of traditional object detection tasks are horizontal bounding boxes. [10] presents a real-time object detection method based on region proposal networks (RPN) which shares feature maps of RPN and R-CNN and use anchors with different sizes and aspect ratios in RPN stage. Though Faster R-CNN shares feature maps, it still requires much computation in R-CNN's fully connected layer. Region-based fully convolutional networks (R-FCN) [11] presents Position-sensitive score maps and Position-sensitive RoI pooling for saving computation in R-CNN stage. Scale variation is always a very challenging issue in object detection; to help solve this problem, [12] presents Feature Pyramid Network (FPN). FPN generates feature maps of different scales on different layers, and detects large objects on higher layers but detects small objects on lower layers; the parameters of RPN is shared over layers. Based on FPN, Mask R-CNN [7] presents RoIAlign which calculates values in RoI via bilinear interpolation instead of maximum to avoid quantization errors, and add several convolution layers on mask-head to generate instance segmentation maps. [13] improves Mask R-CNN by adding Bottom-up Path Augmentation and feature fusion.

Two-stage methods require more computation than one-stage methods, so one-stage methods are more suitable for real-time object detection tasks. Single shot multibox detector (SSD) [14] generates multiple layers, and then detect objects with different sizes on different layers. Deconvolutional single shot detector (DSSD) [15] upsamples feature maps and detects small objects on lower layers which improves SSD's performance for small objects. [16] presents focal loss to handle the imbalance between positive and negative samples. Although anchors are widely used in object detection, many models adopt anchor-free method. [17] does not use anchors in RPN, but uses a shrunk segmentation map as the label. [18] also uses segmentation maps as ground truth. GA-RPN [19] combines anchor-free and anchor-based ideas: the label for the first step is generated by a shrunk segmentation map, and the label for the second step is calculated based on the output anchor of the first step.

Traditional object detection in aerial images only focuses on horizontal bounding box. [20] presents local-contextual feature fusion network which is designed for remote sensing images. The RPN includes multiangle, multiscale and multiaspect-ratio anchors which can deal with the oriented objects, but the final output bounding box is still horizontal. [21] presents rotation-invariant matrix (RIM) which can get both the angular spatial information and the radial spatial information. [22] presents an automatic and accurate localization method for detecting objects in high resolution remote sensing images based on Faster R-CNN. [23] presents a method to detect seals in aerial remote sensing images based on convolutional network. [24] presents a hybrid DNN (HDNN), where the last convolutional and max-pooling layer of DNN is divided into multiple blocks, so HDNN can generate multi-scale features which improves the detector's performance for small objects. Unlike the images used for general object detection, the aerial image has a large resolution. However, large models cannot be implemented due to limited memory. So [25] proposes a self-reinforced network named remote sensing region-based convolutional neural network (R2-CNN) including Tiny-Net and intermediate global attention blocks. It adopts a lightweight residual structure, so the network can feedward huge resolution sensing images at high speeds. [26] proposes a novel method for ship detection in synthetic aperture radar (SAR) images. It redesign the network structure, does not pre-train on ImageNet, and specifically design the system for small objects such as ships.

### B. Oriented objects detection

Oriented object detection is firstly presented in text detection field. Textboxes [27] presents a novel SSD-based text detection method, which adapts the size and aspect ratio of anchor and uses $1 \times 5$ convolutional filters for long text lines. Textboxes++ [2] is based on textboxes but directly regresses the 8 vertices of the oriented bounding box. [28] designs rules for calculating the order of the vertices of the oriented bounding box, and proposes parallel IoU computation for timesaver. [3] presents rotation region proposal networks (RRPN) that proposes oriented bounding boxes in RPN stage and uses Rotation Region-of-Interest (RRoI) pooling layer in R-CNN stage. The aspect ratio of text lines varies greatly, and limited anchors cannot cover the size or aspect ratio of all objects; thus, many methods are anchor-free. Both [4] and [1] generate labels with shrunk segmentation maps, and regress the vertices or angles of the bounding box on positive pixels. [29] generates a corner map and a position-sensitive segmentation map, calculates oriented bounding boxes based on the corner map, and calculates the score for each bounding box using the position-sensitive segmentation map. [30] presents anchor-free region proposal network (AF-RPN) based on Faster R-CNN with the same design as FPN [12], and the label is calculated from the shrunk segmentation map instead of the anchors.

Horizontal bounding boxes cannot closely surround objects in aerial images, so the academic community begins to pay attention to oriented bounding box detection in aerial images. [31] labels a large-scale dataset which contains 15
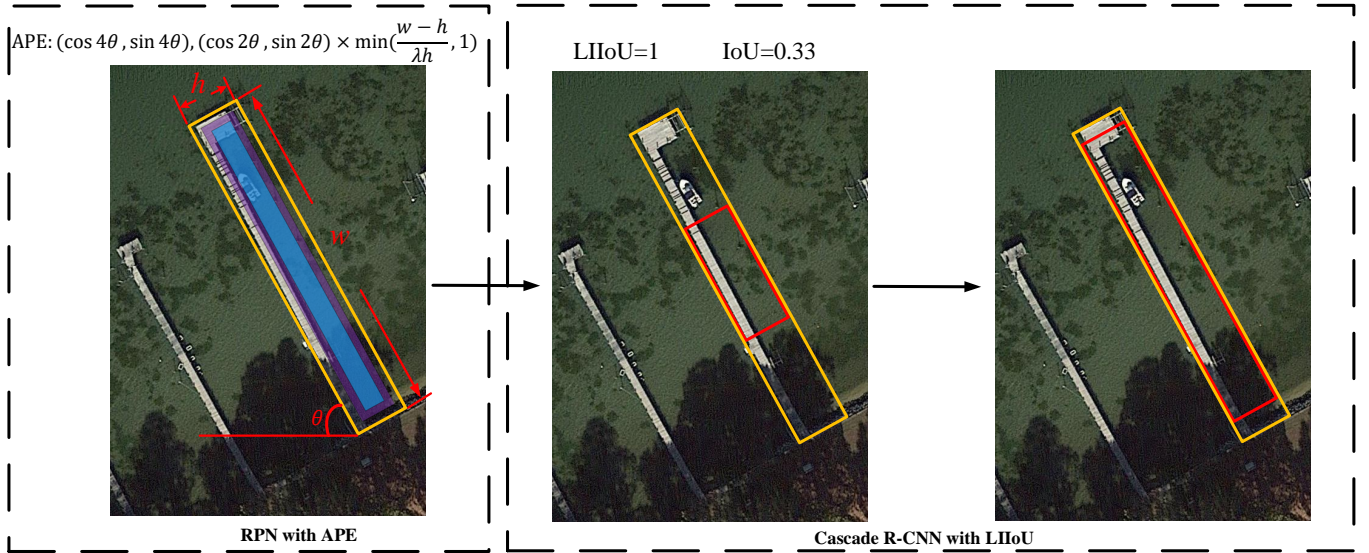
Fig. 2: Illustration of our proposed architecture. From left to right: Anchor-free RPN, Cascade R-CNN. Yellow bounding box is groud truth, red bounding box is proposed box.

categories and 188282 instances, each labeled with an arbitrary quadrilateral (8 vertices). A novel detector which directly regresses 8 vertices based on Faster R-CNN is also presented. ICPR ODAI [32] and CVPR DOTA [33] competitions are organized based on this dataset. [5] presents a two-stage R-CNN method with RoI Transformer, which, in the first step, proposes horizontal bounding boxes. The first R-CNN outputs oriented bounding boxes, and the inputting of the second R-CNN are oriented bounding boxes. [34] proposes a novel method named rotatable region-based residual network (R3-Net) which can detect multi-oriented vehicles in aerial images and videos. The rotatable region proposal network (R-RPN) is adopted to generate rotatable region of interests (R-RoIs) which crops rotated rectangle areas from feature maps.

### III. METHOD

#### A. Overview

Our method's pipeline is shown in Fig. 2. Recently, anchor-free methods [17] [18] [19] [30] are widely used in object detection. In this study, we also use anchor-free RPN. In particular, the label of RPN is not calculated based on the overlap between the anchor and the ground truth; instead, the label is generated from the shrunk segmentation map of the oriented bounding box. Unlike traditional object detection tasks, the output bounding box of RPN is oriented, so a novel angle embedding method is adopted to better represent oriented bounding boxes. Segmentation maps with 8 channels (x, y, w, h, angle embedding) are generated in RPN stage. Then our model proposes oriented bounding boxs with Rotated RoIAlign in R-CNN stage, where Cascade R-CNN is used. In the first R-CNN, a novel IoU calculation method named length independent IoU (LIIoU) is adopted. To make IoU independent of the length of target box, we intercept part of the long side of the target box to obtain the maximum IoU between the
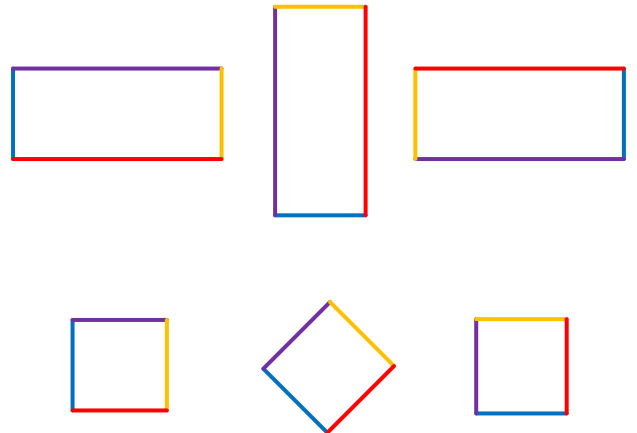


Fig. 3: The peroid of oriented bounding box. Top: rectangle whose peroid is $180°$, down: squre whose peroid is $90°$. (bounding box's sides are in different colors for better visualization)

proposed box and the intercepted target box. In this way, some long boxes will also have corresponding positive samples. Traditional IoU calculation method is used in the second R-CNN. The backbone of our network is based on FPN [12], and we augment the network in the same way as PANet [13] by adopting bottom-up Path Augmentation and feature fusion. Next, we will introduce each component in detail.

#### B. Network Design

Inspired by recent object detection works [12] [7] [13], we use Feature Pyramid Networks (FPN) as our backbone. FPN generates multiple feature maps with different sizes, and detect objects of different sizes on different layers. FPN is robust to

scale variation expecially for small objects, which is suitable for this task. Besides scale variation, aspect ratio variation is another challenging problem. Most traditional object detection methods use anchors of different sizes and aspect ratios to calculate labels in RPN stage. Thus, we have to manually set the hyperparameter of the anchors which is too troublesome, and when the aspect ratio of objects varies greatly, limited anchors cannot cover all the objects. The performance of these detectors highly relies on anchor design. Recently, many methods [17] [18] [19] [30] [4] [1] adopt the anchor-free strategy. In this study, we also adopt anchor-free method and generate the label of RPN from the shrunk segmentation map. Different layers extract different features, and the detector can achieve better performance by combining these features [13]. In R-CNN stage, we fuse features of different layers after the first fully connected layer with max pooling.

### C. Anchor-free label generation

Region proposals network (RPN) is adopted to propose candidate bounding boxes. Most of the previous methods are based on anchors in this stage. Considering the huge difference in aspect ratio, we use anchor-free RPN. The shrunk segmentation label is shown in Fig. 4. The shrinking method is the same as EAST [4]. In particular, $r_1$ is set to 0.1 and $r_2$ is set to 0.25. We shrink the oriented bounding box with $r_2$ ratio and set the pixels in the shrunk bounding box to positive samples (blue area). Next, we shrink the oriented bounding box with $r_1$ ratio, set the pixels in the shrunk bounding box but not set to positive samples (blue area) to "do not care" (purple area), and set the loss weight of these pixels to 0. FPN outputs multi-scale feature maps, and we detect objects of different scales on different layers. We assign a target object whose shorter side is $h$ to the level $p_k$, and $k$ is calculated as follows:

$$k = \lfloor k_0 + \log_2(h/128) \rfloor \quad (1)$$

where $k$ is the layer that objects should be assigned to; $k_0$ is the target layer when the height $h$ of the object is greater than 128 and less than 256, which we set to 4. As the objects of different scales share the regression and classification parameters of RPN, the regression targets should be normalized. An oriented bounding box is labeled as:

$$(x_c, y_c, w, h, \theta) \quad (2)$$

where $(x_c, y_c)$ is the coordinates of the center point, $w, h$ are the lengths of the long side and the short side respectively, and $\theta$ is the angle of the long side. The pixel on $k$ layer is labeled as $x_k, y_k$. First, we normalize the target bounding box with the stride of $k$ layer:

$$x'_c = \frac{x_c}{s_k}, y'_c = \frac{y_c}{s_k}, w' = \frac{w}{s_k}, h' = \frac{h}{s_k}, \theta' = \theta \quad (3)$$

where $s_k$ is the stride of $k$ layer calculated as:

$$s_k = 2 \times 2^k \quad (4)$$

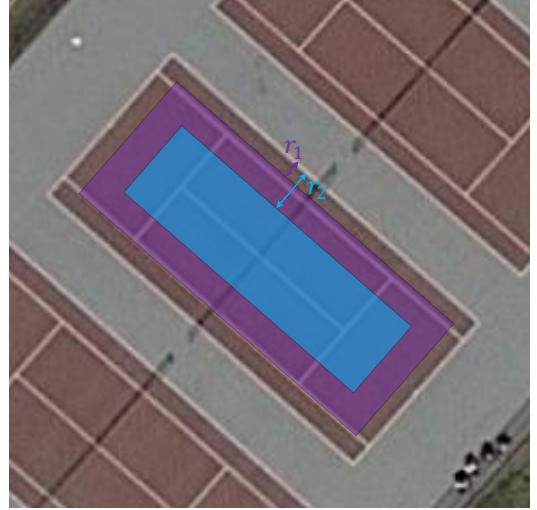The regression targets are calcualted as follows:



Fig. 4: The shrunk segmentation label of anchor-free method. Purple area is ignored area which is shrunk with $r_1$ ratio, blue area is positive area which is shrunk with $r_2$ ratio.

$$t_{x_c} = \frac{x'_c - x_k}{N}, \quad t_{y_c} = \frac{y'_c - y_k}{N}$$
$$t_w = \log \frac{w'}{N}, \quad t_h = \log \frac{h'}{N} \quad (5)$$

where N is a constant and is set to 6 as default.

### D. Adaptive Period Embedding

A horizontal bounding box can be easily represented by 4 variables $(x, y, w, h)$. But we need an extra variable $\theta$ to represent an oriented bounding box. The primary challenge of oriented bounding box detection is to regress the angle of objects. The property of $\theta$ is different from other variables such as x, y, w, h, as $\theta$ is a periodic variable. As shown in Fig. 3, if the length and width of the rectangle are equal, the rectangle is a square, and the peroid of $\theta$ is $90°$. Otherwise, the peroid of $\theta$ is $180°$. In neural networks, the periodic property cannot be represented by one variable. Even though [9] [35] [36] all use two-dimensional periodic vector $(\cos\theta, \sin\theta)$ for representing angle, they do not adapt vector's period. The proposed Adaptive Period Embedding (APE) uses two two-dimensional vectors to represent the angle. The the first vector has a period of $90°$ and can be formulated as:

$$\mathbf{u}_1 = (\cos 4\theta, \sin 4\theta) \quad (6)$$

where $\theta$ is the angle of rectangle's long side. The period of the second vector is $180°$ which represents the angle of rectangle's long side. It is calculated as follows:

$$\mathbf{v} = (\cos 2\theta, \sin 2\theta) \quad (7)$$

$$\mathbf{u}_2 = \mathbf{v} \times \min(\frac{(w-h)}{\lambda h}, 1) \quad (8)$$

where $\lambda$ is set to 0.5, w is rectangle's long side, h is short side. Each component of $\mathbf{u}_1, \mathbf{u}_2$ is in $[-1, 1]$, so we use sigmoid as activation, and then multiply them by 2 and

subtract 1. Smooth L1 loss [37] is used in all regression tasks of this study which can be formulated as:

$$\text{smooth}_{L_1}(z, z^*) = \begin{cases} 0.5(z - z^*)^2 & \text{if } |z - z^*| < 1 \\ |z - z^*| - 0.5 & \text{otherwise} \end{cases} \quad (9)$$

The final outputs of the neural network are $(x, y, w, h, \mathbf{u}_1, \mathbf{u}_2)$. Next, we calculate the angle of the rectangle's long side based on $(\mathbf{u}_1, \mathbf{u}_2)$. Firstly, $\theta_{90°}$ whose peroid is $90°$ can be calculated as:

$$\theta_{90°} = \frac{atan2(\mathbf{u}_1)}{4} \quad (10)$$

where atan2 function calculates one unique arctangent value from a two-dimensional vector. The $\theta$ of rectangle's long side may be $\theta_{90°}$ or $\theta_{90°} + 90°$. The $\theta_{180°}$ whose peroid is $180°$ can be calculated as:

$$\theta_{180°} = \frac{atan2(\mathbf{u}_2)}{2} \quad (11)$$

Then we calculate the distance between $\theta_{90°}$ and $\theta_{180°}$

$$\text{dis} = |(2\theta_{90°} - 2\theta_{180°} + 180°) \mod 360° - 180°| \quad (12)$$

So the final $\theta$ is calculated as:

$$\theta = \begin{cases} \theta_{90°} & \text{dis} < 90° \\ \theta_{90°} + 90° & \text{otherwise} \end{cases} \quad (13)$$

$\theta_{180°}$ is the angle of the rectangle's long side. When the length and width of the rectangle are equal, the norm of $\mathbf{u}_2$ is nearly zero, so the angle calculated by $\mathbf{u}_2$ is not accurate. $\theta_{90°}$ is accurate but it may be the angle of the long side or the short side. So we find the angle closer to $\theta_{180°}$ from $\theta_{90°}$ and $\theta_{90°} + 90°$, which is the final result.

*E. length independent IoU*

IoU is the evaluation protocol of object detection, the more accurate the regression, the better the performance. But the receptive field of a neural network is limited and thus cannot cover some long objects. The detector proposes candidate bounding boxes in RPN, and then classifies and regresses these boxes again. The result of R-CNN highly relies on the output bounding boxes of RPN. In R-CNN stage, only the proposed bounding boxes whose IoU is higher than 0.5 is set to positive samples. Some target objects that are not well regressed in RPN cannot be detected in R-CNN. One idea is multiple regression [38] in R-CNN stage, but if there are no positive proposed bounding boxes in the first R-CNN, the improvement is limited. Considering this, we propose a novel IoU calculation method named length independent IoU (LIIoU). We intercept part of the target box along its long side, and make the length of the intercepted box the same as the proposed box. The presented method is inspired by Seglink [39], but in our method, the aspect ratio of the proposed bounding box is arbitrary. As shown in Fig. 2, the traditional IoU is only 0.3, but our proposed LIIoU is nearly 1. The details of LIIoU calculation is illustrated in Fig. 5, where AB is the center line of the target box, and point C is the center of the proposed box. First, we find the perpendicular of AB

through point C and label the intersection of the perpendicular and AB as point D. Next, we intercept a rectangle from the target bounding box as follows: if the length of the target box is smaller than the proposed box, we do not intercept; otherwise, the center of the intercepted rectangle is point D and the length is the same with proposed box (green box). Finally, we calculate IoU between the intercepted target box and the proposed box. The procedure is summarized in Algorithm 1. In this way, more bounding boxes will regress targets in R-CNN which can improve the overall quality of the bounding boxes

---

**Algorithm 1** LIIoU calculation

**Input**: $pbbox(x^p, y^p, w^p, h^p, \theta^p)$, $gbbox(x^g, y^g, w^g, h^g, \theta^g)$
$pbbox$ - proposed bounding box
$gbbox$ - ground truth bounding box

**Output**: LIIoU

1: **if** $w^p >= w^g$ **then**
2:     $x'^g = x^g; y'^g = y^g; w'^g = w^g; h'^g = h^g; \theta'^g = \theta^g$
3: **else**
4:     $\mathbf{A}_x = x^g - \cos(\theta^g) \times \frac{w^g}{2}$
5:     $\mathbf{A}_y = y^g - \sin(\theta^g) \times \frac{w^g}{2}$
6:     $\mathbf{B}_x = x^g + \cos(\theta^g) \times \frac{w^g}{2}$
7:     $\mathbf{B}_y = y^g + \sin(\theta^g) \times \frac{w^g}{2}$
8:     $\mathbf{C}_x = x^p; \mathbf{C}_y = y^p$
9:     $z = \frac{(\mathbf{C}-\mathbf{A}) \cdot (\mathbf{B}-\mathbf{A})}{||(\mathbf{B}-\mathbf{A})||}$
10:     $w_1 = z - \frac{w^p}{2}; w_2 = z + \frac{w^p}{2}$
11:     **if** $w_1 <= 0$ **then**
12:       $w_1 = 0; w_2 = w^p$
13:     **else if** $w_2 >= w^g$ **then**
14:       $w_2 = w^g; w_1 = w^g - w^p$
15:     **end if**
16:     $x'^g = A_x + \cos(\theta) \times \frac{w_2+w_1}{2}; y'^g = A_y + \sin(\theta) \times \frac{w_2+w_1}{2}$
17:     $w'^g = w_2 - w_1; h'^g = h^g; \theta'^g = \theta^g$
18: **end if**
19: calulate overlaps between $(x^p, y^p, w^p, h^p, \theta^p)$ and $(x'^g, y'^g, w'^g, h'^g, \theta'^g)$

---

*F. Cascade R-CNN*

As shown in Fig. 2, two R-CNNs are used after RPN. In the first R-CNN, we only refine the center, height and width of the oriented bounding box without regressing the vertices of the target box. This is because the output of the first R-CNN is the input of the second R-CNN, and Rotated RoIAlign can only handle oriented rectangle but not quadrangle. In the second R-CNN, we regress the vertices of the target box.

Rotated RoIAlign is adopted, so the ground truth is calculated in a rotated coordinate system. If the center of a Rotated RoIAlign is $(x_c^p, y_c^p)$ and the angle is $\theta^p$, the affine

Calculate point D                 Intercept target box           Calculate IoU
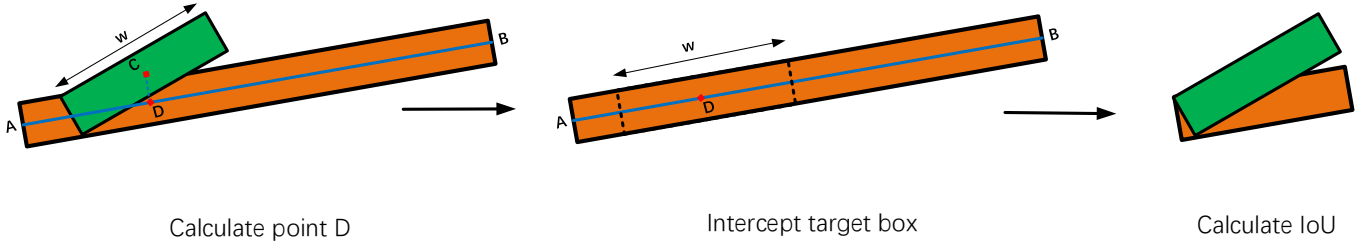
Fig. 5: The details of LIIoU calculation, green bounding box is proposed bounding box, orange bounding box is target box.

transformation can be represented by an affine matrix:

$$
\begin{aligned}
\boldsymbol{M} &= \begin{bmatrix} 1 & 0 & x_c^{\mathrm{p}} \\ 0 & 1 & y_c^{\mathrm{p}} \\ 0 & 0 & 1 \end{bmatrix} * \begin{bmatrix} \cos\theta^{\mathrm{p}} & \sin\theta^{\mathrm{p}} & 0 \\ -\sin\theta^{\mathrm{p}} & \cos\theta^{\mathrm{p}} & 0 \\ 0 & 0 & 1 \end{bmatrix} * \begin{bmatrix} 1 & 0 & -x_c^{\mathrm{p}} \\ 0 & 1 & -y_c^{\mathrm{p}} \\ 0 & 0 & 1 \end{bmatrix} \\
&= \begin{bmatrix} \cos\theta^{\mathrm{p}} & \sin\theta^{\mathrm{p}} & (1-\cos\theta^{\mathrm{p}})x_c^{\mathrm{p}} - y_c^{\mathrm{p}} * \sin\theta^{\mathrm{p}} \\ -\sin\theta^{\mathrm{p}} & \cos\theta^{\mathrm{p}} & (1-\cos\theta^{\mathrm{p}})y_c^{\mathrm{p}} + x_c^{\mathrm{p}} * \sin\theta^{\mathrm{p}} \\ 0 & 0 & 1 \end{bmatrix}
\end{aligned}
\tag{14}
$$

$$
\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = \boldsymbol{M} \begin{pmatrix} x' \\ y' \\ 1 \end{pmatrix}
\tag{15}
$$

We set the coordinate system to rotated coordinate system with Eq. 14 and Eq. 15. The final ground truth in the rotated coordinate system is $(x, y)$, and $(x', y')$ is the coordinates in the original coordinate system. In the first R-CNN, the regression targets are $(t_{x_c}^{\mathrm{rcnn1}}, t_{y_c}^{\mathrm{rcnn1}}, t_w^{\mathrm{rcnn1}}, t_h^{\mathrm{rcnn1}})$ which can be formulated as:

$$
t_{x_c}^{\mathrm{rcnn1}} = \frac{x_c - x_c^{\mathrm{p}}}{w^{\mathrm{p}}}; t_{y_c}^{\mathrm{rcnn1}} = \frac{y_c - y_c^{\mathrm{p}}}{h^{\mathrm{p}}}
\tag{16}
$$

$$
t_w^{\mathrm{rcnn1}} = \log(\frac{w}{w^{\mathrm{p}}}); t_h^{\mathrm{rcnn1}} = \log(\frac{h}{h^{\mathrm{p}}})
\tag{17}
$$

In the second R-CNN, the regression targets are $(t_{x_i}^{\mathrm{rcnn2}}, t_{y_i}^{\mathrm{rcnn2}}), i = 1, 2, 3, 4$ which can be formulated as:

$$
t_{x_i}^{\mathrm{rcnn2}} = \frac{x_i - x_c^{\mathrm{p}}}{w^{\mathrm{p}}}; t_{y_i}^{\mathrm{rcnn2}} = \frac{y_i - y_c^{\mathrm{p}}}{h^{\mathrm{p}}}, i = 1, 2, 3, 4
\tag{18}
$$

where $(x_c, y_c, w, h)$ are the center, width and height of the ground truth, $(x_i, y_i)$ is the vertex of the ground truth bounding box, and $(x_c^{\mathrm{p}}, y_c^{\mathrm{p}}, w^{\mathrm{p}}, h^{\mathrm{p}})$ are the center, width and height of the proposed bounding box.

### G. Experiment

*1) Datasets:* DOTA [31] is a large dataset which contains 2806 aerial images from different sensors and platforms. The size of the image varies greatly, ranging from about $800 \times 800$ to $4000 \times 4000$ pixels, so it is necessary to crop the image and detect the objects in the cropped images. As the instances in arial images are oriented such as car, ship and bridge, each instance is labeled by an arbitrary (8 d.o.f.) quadrilateral. For the oriented task, the output bounding boxes are quadrilatera; to evaluate the performance of our detector on quadrilateral, we use the evaluation system provided along with this dataset. There are two versions of DOTA dataset, DOTA-v1.0 and DOTA-v1.5; DOTA-v1.5 fixes some errors and is provided for DOAI2019 competition [33]. We use DOTA-v1.5 for this competition, but in the following experiments, we use DOTA-v1.0 for fair comparison.

*2) Implementation Details:* The backbone of our detector is ResNet-50 [40] pre-trained on ImageNet [41]. The number of FPN channels is set to 256. In R-CNN stage, two fully connected (FC) layers are used, the channel of which is set to 1024. Feature fusion is applied after the first FC layer along with maxpooling. Batchnorm is not used in this study. Our network is trained with SGD, where the batchsize of images is 1 and the initial learning rate is set to 0.00125, which is then divided by 10 at $\frac{2}{3}$ and $\frac{8}{9}$ of the entire training. Due to the limited memory, we crop images to $1024 \times 1024$ with the stride of 256 for training and testing. The model is trained and tested at a single scale. Data augmentation is used for better performance; in particular, we randomly rotate images with angle among 0, $\pi/2$, $\pi$, $3\pi/2$, and class balance resampling is adopted to solve class imbalance problem. In default, we train our model with training set and evaluate it on validation set and testing set.

*3) Ablation Study:* In order to evaluate the effect of each component, we conduct abalation experiments on validation set of DOTA. The model is not modified except the component being tested.

**The effect of adaptive period embedding:** We need to propose oriented bounding boxes in RPN stage, but it is challenging to effectively represent a oriented bounding box. Most of previous methods [4] [5] [3] which directly regress the angle do not notice the periodicity of the angle. When the angle is too diverse, the performance of the system will drop significantly. To evaluate whether the proposed APE can well handle the diversity of the angles, we conduct abalation experiments: one model directly regress the angle of the long side of the target box, while the other regresses adaptive period

Fig. 6: The comparison of RPN with APE and without APE. top: without APE, down: with APE.

TABLE I: The experiment of APE on DOTA validation set in RPN stage. (in %)

| Methods | w/o APE | w/ APE |
|---------|---------|--------|
| AP | 70.16 | 72.20 |

TABLE II: The ablation experiments of LIIoU and IoU on DOTA validation set (in %).

| Methods | Faster R-CNN | Cascade R-CNN | Cascade R-CNN+LIIoU |
|---------|--------------|---------------|---------------------|
| mAP | 71.40 | 72.76 | **73.88** |

TABLE III: Results on DOTA testing set (in %). * indicates validation set is also used for training, otherwise only training set is used for training.

| Method | Ours | Ours * | FR-O [31] | RoI Transformer * [5] |
|--------|------|--------|-----------|------------------------|
| Plane | 89.67 | 89.96 | 79.09 | 88.64 |
| BD | 76.77 | 83.62 | 69.12 | 78.52 |
| Bridge | 51.28 | 53.42 | 17.17 | 43.44 |
| GTF | 71.65 | 76.03 | 63.49 | 75.92 |
| SV | 73.11 | 74.01 | 34.2 | 68.81 |
| LV | 77.18 | 77.16 | 37.16 | 73.68 |
| Ship | 79.54 | 79.45 | 36.2 | 83.59 |
| TC | 90.79 | 90.83 | 89.19 | 90.74 |
| BC | 79.01 | 87.15 | 69.6 | 77.27 |
| ST | 84.54 | 84.51 | 58.96 | 81.46 |
| SBF | 66.51 | 67.72 | 49.4 | 58.39 |
| RA | 64.71 | 60.33 | 52.52 | 53.54 |
| Harbor | 73.97 | 74.61 | 57.79 | 62.83 |
| SP | 67.73 | 71.84 | 44.8 | 58.93 |
| HC | 58.40 | 65.55 | 46.3 | 47.67 |
| **mAP** | 73.66 | **75.75** | 52.93 | 69.56 |

embedding (APE) vectors. We evaluate the quality of proposed oriented bounding box in RPN stage. Network only classifies objects into 2 classes (positive sample and negative sample) in RPN stage. As shown in Table I, RPN achieves much better performance with APE. We show the comparison in Fig. 6, where we can see that RPN outputs more accurate angle with APE compared with directly regressing the angle.

**LIIoU vs. IoU:** To evaluate the efficiency of LIIoU, we conduct control experiment. Faster R-CNN means there is only one R-CNN. When Cascade R-CNN is adopted, two R-CNNs are used. In the first model, we calculate the overlap between oriented bounding boxes with traditional IoU in both two R-CNNs. In the second model, the overlap between oriented bounding boxes is calculated with LIIoU in the first R-CNN and with traditional IoU in the last R-CNN, and the threshold is set to 0.5. Results are shown in Table II, where we can see that Cascade R-CNN gains much better performance with LIIoU. We show their comparison in Fig. 7, where we can find that LIIoU can improve the quality of the proposed bounding boxes and the recall rate. Regardless of the aspect ratio and

size, nearly every object has positive samples with LIIoU, so the detector can handle objects with large aspect ratios and lengths well.

*4) Comparing with other state-of-the-art methods:* We compare our method with other state-of-the-art methods. The results are shown in Table III. Our model is trained and tested with the single-scale setting. When our model is only trained with training set ex validation set, our method significantly outperforms other methods, if validation set is also used for training, it achieves better performance. The detection results are shown in Fig. 8. The angle, size, aspect ratio of objects in aerial images vary greatly, but our proposed method can well handle these challenging conditions.

TABLE IV: Task1 - Oriented Leaderboard on DOAI2019 (in %).

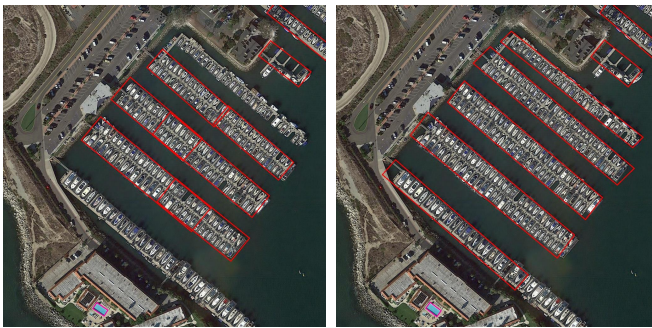| Team Name | USTC-NELSLIP | pca_lab | czh | AICyber | CSULQQ | peijin |
|---|---|---|---|---|---|---|
| Plane | 89.2 | 88.2 | 89.0 | 88.4 | 87.8 | 80.9 |
| BD | 85.3 | 86.4 | 83.2 | 85.4 | 83.6 | 83.6 |
| Bridge | 57.3 | 59.4 | 54.5 | 56.7 | 56.7 | 55.1 |
| GTF | 80.9 | 80.0 | 73.8 | 74.4 | 74.4 | 70.7 |
| SV | 73.9 | 68.1 | 72.6 | 63.9 | 63.2 | 59.9 |
| LV | 81.3 | 75.6 | 80.3 | 72.7 | 71.0 | 76.4 |
| Ship | 89.5 | 87.2 | 89.3 | 87.9 | 87.8 | 88.3 |
| TC | 90.8 | 90.9 | 90.8 | 90.9 | 90.8 | 90.9 |
| BC | 85.9 | 85.3 | 84.4 | 86.3 | 84.6 | 79.2 |
| ST | 85.6 | 84.1 | 85.0 | 85.0 | 84.0 | 78.3 |
| SBF | 69.5 | 73.8 | 68.7 | 68.9 | 67.8 | 59.1 |
| RA | 76.7 | 77.5 | 75.3 | 76.0 | 75.5 | 74.8 |
| Harbor | 76.3 | 76.4 | 74.2 | 74.1 | 67.4 | 74.1 |
| SP | 76.0 | 73.7 | 74.4 | 72.9 | 71.2 | 74.9 |
| HC | 77.8 | 69.5 | 73.4 | 73.4 | 68.8 | 59.8 |
| CC | 57.3 | 49.6 | 42.1 | 37.9 | 22.5 | 39.5 |
| **mAP** | 78.3 | 76.6 | 75.7 | 74.7 | 72.3 | 71.6 |



Fig. 7: The comparison of LIIoU and IoU. From Left to right: calculate overlaps with IoU in the first R-CNN, calculate overlaps with LIIoU in the first R-CNN (overlaps are both calculated with IoU in the last R-CNN).

*5) DOAI2019 competition:* DOAI2019 competition [33] is held in workshop on Detecting Objects in Aerial Images in conjunction with IEEE CVPR 2019. The competition is more difficult and requires detecting all objects including samples labeled as difficult. Based on our proposed methods including APE and LIIoU, we adopt class balance resampling, image rotation, multi-scale training and testing and model assembling for better performance. Three models are used whose backbone is ResNeXt-101(32x4) [42]. Finally, we combine the training set with the validation set for training. The results of the competition are shown in Table IV. Our method wins the first place on oriented task, with a gain of about 1.7% over the most competing competitor.

*6) Conclusion and Future Work:* Detecting oriented objects in arial images is a challenging task. In this study, we make full use of the periodicity of the angle. A novel method named adaptive period embedding (APE) is proposed which can well regress oriented bounding boxes in arial images. The vector with adaptive period can learn the periodicity of the angle, which can not be implemented with one-dimensional vector. The proposed method can be applied to both one-stage methods such as RPN and two-stage methods, and we believe other detectors can also directly adopt APE module.

Besides, we propose a novel length independent IoU (LIIoU). LIIoU set more proposed bounding boxes to positive samples expecially for long objects which can improve the quality of R-CNN regression. Our ablation study proves that each proposed module is effective. Our method achieves state-of-the-art on DOTA. Based on our method, we win the first place on oriented task of DOAI2019. In the future, we will explore more efficient and accurate detector for detecting oriented objects in arial images.

REFERENCES

[1] W. He, X.-Y. Zhang, F. Yin, and C.-L. Liu, "Deep direct regression for multi-oriented scene text detection," *arXiv preprint arXiv:1703.08289*, 2017.
[2] M. Liao, B. Shi, and X. Bai, "Textboxes++: A single-shot oriented scene text detector," *IEEE Transactions on Image Processing*, vol. 27, no. 8, pp. 3676–3690, 2018.
[3] 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*, 2018.
[4] X. Zhou, C. Yao, H. Wen, Y. Wang, S. Zhou, W. He, and J. Liang, "East: An efficient and accurate scene text detector," *arXiv preprint arXiv:1704.03155*, 2017.
[5] J. Ding, N. Xue, Y. Long, G.-S. Xia, and Q. Lu, "Learning roi transformer for detecting oriented objects in aerial images," *arXiv preprint arXiv:1812.00155*, 2018.
[6] Y. Dai, Z. Huang, Y. Gao, and K. Chen, "Fused text segmentation networks for multi-oriented scene text detection," *arXiv preprint arXiv:1709.03272*, 2017.
[7] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," *arXiv preprint arXiv:1703.06870*, 2017.
[8] Y. Zhu and J. Du, "Sliding line point regression for shape robust scene text detection," *arXiv preprint arXiv:1801.09969*, 2018.
[9] S. Long, J. Ruan, W. Zhang, X. He, W. Wu, and C. Yao, "Textsnake: A flexible representation for detecting text of arbitrary shapes," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 20–36.
[10] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 6, pp. 1137–1149, 2017.
[11] 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.
[12] T.-Y. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie, "Feature pyramid networks for object detection." in *IEEE Conference on Computer Vision and Pattern Recognition*, vol. 1, no. 2, 2017, p. 4.

[13] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia, "Path aggregation network for instance segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8759–8768.

[14] 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.

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

[16] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," *arXiv preprint arXiv:1708.02002*, 2017.

[17] L. Huang, Y. Yang, Y. Deng, and Y. Yu, "Densebox: Unifying landmark localization with end to end object detection," *arXiv preprint arXiv:1509.04874*, 2015.

[18] 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.

[19] J. Wang, K. Chen, S. Yang, C. C. Loy, and D. Lin, "Region proposal by guided anchoring," *arXiv preprint arXiv:1901.03278*, 2019.

[20] K. Li, G. Cheng, S. Bu, and X. You, "Rotation-insensitive and context-augmented object detection in remote sensing images," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 4, pp. 2337–2348, 2018.

[21] G. Wang, X. Wang, B. Fan, and C. Pan, "Feature extraction by rotation-invariant matrix representation for object detection in aerial image," *IEEE Geoscience and Remote Sensing Letters*, vol. 14, no. 6, pp. 851–855, 2017.

[22] Y. Long, Y. Gong, Z. Xiao, and Q. Liu, "Accurate object localization in remote sensing images based on convolutional neural networks," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 5, pp. 2486–2498, 2017.

[23] A.-B. Salberg, "Detection of seals in remote sensing images using features extracted from deep convolutional neural networks," in *Geoscience and Remote Sensing Symposium*, 2015, pp. 1893–1896.

[24] X. Chen, S. Xiang, C.-L. Liu, and C.-H. Pan, "Vehicle detection in satellite images by hybrid deep convolutional neural networks," *IEEE Geoscience and remote sensing letters*, vol. 11, no. 10, pp. 1797–1801, 2014.

[25] J. Pang, C. Li, J. Shi, Z. Xu, and H. Feng, "R-cnn: Fast tiny object detection in large-scale remote sensing images," *IEEE Transactions on Geoscience and Remote Sensing*, 2019.

[26] Z. Deng, H. Sun, S. Zhou, and J. Zhao, "Learning deep ship detector in sar images from scratch," *IEEE Transactions on Geoscience and Remote Sensing*, 2019.

[27] M. Liao, B. Shi, X. Bai, X. Wang, and W. Liu, "Textboxes: A fast text detector with a single deep neural network." in *AAAI*, 2017, pp. 4161–4167.

[28] Y. Liu and L. Jin, "Deep matching prior network: Toward tighter multi-oriented text detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3454–3461.

[29] P. Lyu, C. Yao, W. Wu, S. Yan, and X. Bai, "Multi-oriented scene text detection via corner localization and region segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7553–7563.

[30] Z. Zhong, L. Sun, and Q. Huo, "An anchor-free region proposal network for faster r-cnn based text detection approaches," *arXiv preprint arXiv:1804.09003*, 2018.

[31] G.-S. Xia, X. Bai, J. Ding, Z. Zhu, S. Belongie, J. Luo, M. Datcu, M. Pelillo, and L. Zhang, "Dota: A large-scale dataset for object detection in aerial images," in *IEEE Conference on Computer Vision and Pattern Recognition*, June 2018.

[32] "Icpr-odai," https://captain-whu.github.io/ODAI/.

[33] "Cvpr-dota," https://captain-whu.github.io/DOAI2019/challenge.html.

[34] Q. Li, L. Mou, Q. Xu, Y. Zhang, and X. X. Zhu, "R-net: A deep network for multioriented vehicle detection in aerial images and videos," *IEEE Transactions on Geoscience and Remote Sensing*, 2019.

[35] Y. Zhu and J. Du, "Textmountain: Accurate scene text detection via instance segmentation," *arXiv preprint arXiv:1811.12786*, 2018.

[36] Y. Xu, Y. Wang, W. Zhou, Y. Wang, Z. Yang, and X. Bai, "Textfield: Learning a deep direction field for irregular scene text detection," *IEEE Transactions on Image Processing*, 2019.

[37] R. Girshick, "Fast r-cnn," in *IEEE international conference on computer vision*, 2015, pp. 1440–1448.

[38] Z. Cai and N. Vasconcelos, "Cascade r-cnn: Delving into high quality object detection," *arXiv preprint arXiv:1712.00726*, 2017.

[39] B. Shi, X. Bai, and S. Belongie, "Detecting oriented text in natural images by linking segments," *arXiv preprint arXiv:1703.06520*, 2017.

[40] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[41] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

[42] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2017, pp. 5987–5995.
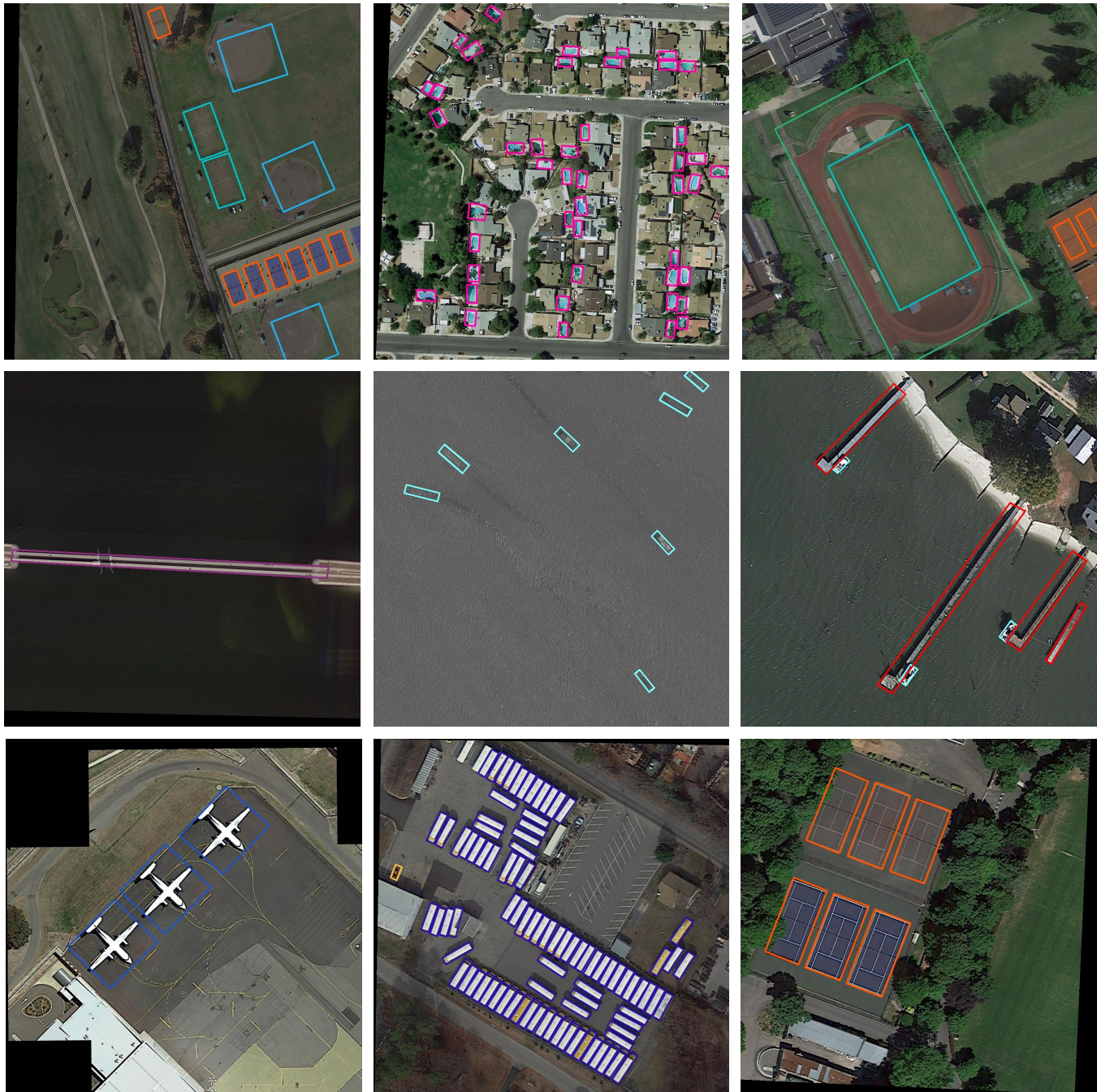
Fig. 8: Some results of our method on DOTA. The image's size is $1024 \times 1024$.