

# Lightweight Oriented Object Detection Using Multiscale Context and Enhanced Channel Attention in Remote Sensing Images

Qiong Ran, Qing Wang, Boya Zhao, Yuanfeng Wu , Senior Member, IEEE, Shengliang Pu, Student Member, IEEE, and Zijin Li

**Abstract**—Object detection is a focal point in remote sensing applications. Remote sensing images typically contain a large number of small objects and a wide range of orientations across objects. This results in great challenges to small object detection approaches based on remote sensing images. Methods directly employ channel relations with equal weights to construct information features leads to inadequate feature representation in complex image small object detection tasks. Multiscale detection methods improve the speed and accuracy of detection, while small objects themselves contain limited information, and the features are easily lost following down-sampling. During the detection, the feature images are independent across scales, resulting in a discontinuity at the detection scale. In this article, we propose the multiscale context and enhanced channel attention (MSCCA) model. MSCCA employs PeleeNet as the backbone network. In particular, the feature image channel attention is enhanced and the multiscale context information is fused with multiscale detection methods to improve the characterization ability of the convolutional neural network. The proposed MSCCA method is evaluated on two real datasets. Results show that for  $512 \times 512$  input images, MSCCA was able to achieve 80.4% and 94.4% mAP on the DOTA and NWPU VHR-10, respectively. Meanwhile, the model size of MSCCA is 21% smaller than that of its predecessor. MSCCA can be considered as a practical lightweight oriented object detection model in remote sensing images.

**Index Terms**—Channel attention, lightweight convolutional neural network (CNN), multiscale context, object detection, remote sensing.

## I. INTRODUCTION

THE object detection plays a key role in remote sensing algorithms and applications. They can be roughly divided into traditional and deep learning object detection approaches.

Manuscript received April 3, 2021; revised May 4, 2021; accepted May 10, 2021. Date of publication May 13, 2021; date of current version June 16, 2021. This work was supported by the National Natural Science Foundation of China under Grant 41871245 and Grant 62001455. (Corresponding author: Yuanfeng Wu.)

Qiong Ran and Qing Wang are with the College of Information Science and Technology, Beijing University of Chemical Technology, Beijing, 100029, China. (e-mail: ranqiong@mail.buct.edu.cn; 2019210520@mail.buct.edu.cn).

Boya Zhao and Shengliang Pu are with the Key Laboratory of Digital Earth Science, Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing 100094, China (e-mail: zhaoby@aircas.ac.cn; puysl@aircas.ac.cn).

Yuanfeng Wu and Zijin Li are with the Key Laboratory of Digital Earth Science, Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing 100094, China, and also with the School of Electronic, Electrical and Communication Engineering, University of Chinese Academy of Sciences, Beijing, 100049, China (e-mail: wuyf@radi.ac.cn; 1700012409@pku.edu.cn).

Digital Object Identifier 10.1109/JSTARS.2021.3079968

Traditional object detection methods (HOG [1], SVM [2], DPM [3], etc.) generally include a region proposal, feature extraction, and classification, resulting in a low detection efficiency and poor accuracy due to complex procedures, a large number of redundant windows and the poor robustness of manual feature extraction methods. Thus, traditional detection methods are hardly meeting the object detection performance demands. The emergence of deep learning-based methods has achieved significant breakthroughs in object detection [4]–[6]. Deep learning-based object detection methods mainly can be divided into two types: 1) two-stage detection models, which define detection as a “coarse-to-fine” process; and 2) one-stage detection models, which define detection as a “one-step” process [7].

Two-stage detection approaches are generally region-based and extract a set of object proposals that potentially contain the objects using methods such as selective searching or region proposals. These sets are subsequently fed into a convolutional neural network (CNN) for feature extraction. The classifiers then predict the presence of an object within each region and recognize the object categories. R-CNN [8] is a typical two-stage detector that generates proposals by selective searching and normalizes their size and inputs them to the CNN to extract the features. SVM is then applied to recognize object categories within each region. Fast R-CNN [9] improves R-CNN by using a multitask loss to increase the detection quality. Faster R-CNN [10] introduces the region proposal network, whereby the majority of the individual blocks in the object detection framework (region proposal, feature extraction, bounding box regression, etc.) are gradually integrated into an end-to-end learning framework. Mask R-CNN [11] includes a branch to segment an object based on faster R-CNN and simultaneously performs instance segmentation and object detection. Libra R-CNN [12] integrates IoU-balanced sampling, a balanced feature pyramid and a balanced L1 loss to reduce the imbalance at the sampling, feature extraction, and training procedures, respectively. Although two-stage object detection methods have made a great progress in detection tasks, they are limited by large amounts of parameters and slow detection speeds. HSP [13] considers the utilization and propagation of hierarchical semantic information in the optimized process of the detection network to improve object detection performance in remote sensing imagery.

One-stage detection methods apply a single CNN to divide the image into multiple regions and simultaneously predict

the bounding boxes and category of each region. This process greatly improves the detection speed, yet reduces the detection accuracy compared to two-stage detectors. YOLO [14] is a typical one-stage object detection method that treats object detection as the solution of a regression problem, applying a single CNN to the full image. This network simultaneously predicts the bounding boxes and category for each region. SSD [15] is an additional one-stage detection method that sets default boxes with different aspect ratios in each feature map to perform multiscale detection, significantly improving the one-stage detector detection accuracy. FMSSD [16] leverages the atrous spatial feature pyramid module to integrate the context information into the framework, improving the robustness of features. RetinaNet [17] proposes the focal loss, whereby the detector pays more attention to samples that are difficult to classify during the training process. This maintains a high detection speed while matching the accuracy of two-stage detection methods. RefineDet [18] proposes the anchor refinement module and object detection module (ODM) to improve the detection efficiency without reducing the detection speed. M2Det [19] proposes the multilevel feature pyramid network and constructs more object feature pyramids to detect objects at different scales. Based on FCN, FCOS [20] is an anchor-free detector that abandons the anchor generation process, reducing memory footprints and improving the detection accuracy. MS-VANs [21] proposed a visual attention-based network and simultaneously predict object class at each pixel of the feature maps, and use a visual attention network to highlight the features from the object region and decrease the influence of cluttered backgrounds. S<sup>2</sup>A-Net [44] implemented full feature alignment and alleviates the inconsistency between regression and classification by using feature alignment module and oriented detection module (ODM).

For small object detection tasks, SNIP [23] and SNIPER [24] employ scale normalization and only detect objects with a fixed size for scale-specific feature maps. SNIPER reduces the computation of the multiscale image pyramid generation and accelerates multiscale training. DEFace [25] proposes the extended feature pyramid network (FPN) [26] module with a receptive context module to enhance the distinguishability and robustness of features. TridentNet [27] constructs a parallel multibranch architecture and adopts a scale-aware training scheme to specialize each branch by sampling the object instances of proper scales for training. SCRDet++ [28] introduces the denoising process to object detection, whereby instance-level denoising on the feature map is performed to enhance the detection of small and cluttered objects. Stitcher [29] dynamically generates stitched images to enrich small object samples and adaptively determines whether the input of the next iteration is the original or the stitched image, which improves the small object loss contribution.

In the traditional convolutional pooling process, the convolution operation does not consider the dependence of each feature channel. In addition, the importance of each channel in the generated feature image is considered to be the same, yet in the actual problem, the importance is actually distinct across channels. One-stage detection methods employ multiscale detection that extracts multiscale feature maps from different layers of the network for predictions. Although this does not

increase the number of calculations, the small object itself has less pixel information and is easily lost during downsampling [30].

In this article, we propose the multiscale context and enhanced channel attention (MSCCA) model. MSCCA employs PeleeNet as the backbone network. In particular, the feature image channel attention is enhanced and the multiscale context information is fused with multiscale detection methods to improve the characterization ability of the CNN [31]. The proposed method is evaluated on two real datasets. Results show that for  $512 \times 512$  input images, MSCCA was able to achieve 80.4% and 94.4% mAP on the DOTA and NWPU VHR-10, respectively. Meanwhile, the model size of MSCCA is 21% smaller than that of its predecessor. MSCCA can be considered as a practical lightweight oriented object detection model in remote sensing images.

The rest of this article is structured as follows. In Section II, the MSCCA model is described. In Section III, two real datasets DOTA and NWPU VHR-10 are presented. In Section IV, the datasets are used to evaluate the proposed MSCCA model. Both the detection accuracy and model size are summarized. Section V concludes this article with some remarks and hints at plausible future research lines.

## II. METHODS

MSCCA model employs the PeleeNet [32] as the backbone, while the enhanced channel attention block is added to balance the channel features that have a positive effect on detection and weakens the channel features that have no effect. Then, the multiscale context structure combines high-level and low-level features within the multiscale detection framework. Fig. 1 presents the whole structure of MSCCA. Objects in remote sensing images typically exhibit large-scale changes, arbitrary-orientation, and irregular shapes. Thus, seven different scale feature maps are employed for multiscale objects. Moreover, the quadrilateral representation is used in location loss for objects with arbitrary orientation and irregular shapes.

### A. Backbone

PeleeNet improves employs a large number of dense layers that consist of two branches that extract multiscale features in the receptive field. ResBlock is added prior to the detection of each feature map. Moreover, MSCCA includes the ECA Block following each transition layer of the network structure. Due to the large size of the remote sensing images, in order to ensure the detection accuracy of small objects, the image is not resized, and the input size set to  $512 \times 512$  pixels.

The entire network consists of five stages. Stage0 only contains Stem Block, which is a low-cost and efficient module that can effectively improve the feature extraction ability with a minimal increase in computational cost. Stem Block initially employs a  $3 \times 3$  convolution layer to downsample the image and subsequently divides it into two branches that use 1) the max-pooling layer to downsample the image and 2) one  $1 \times 1$  and one  $3 \times 3$  convolution layer. The two branches are merged to the channel dimension via concat.

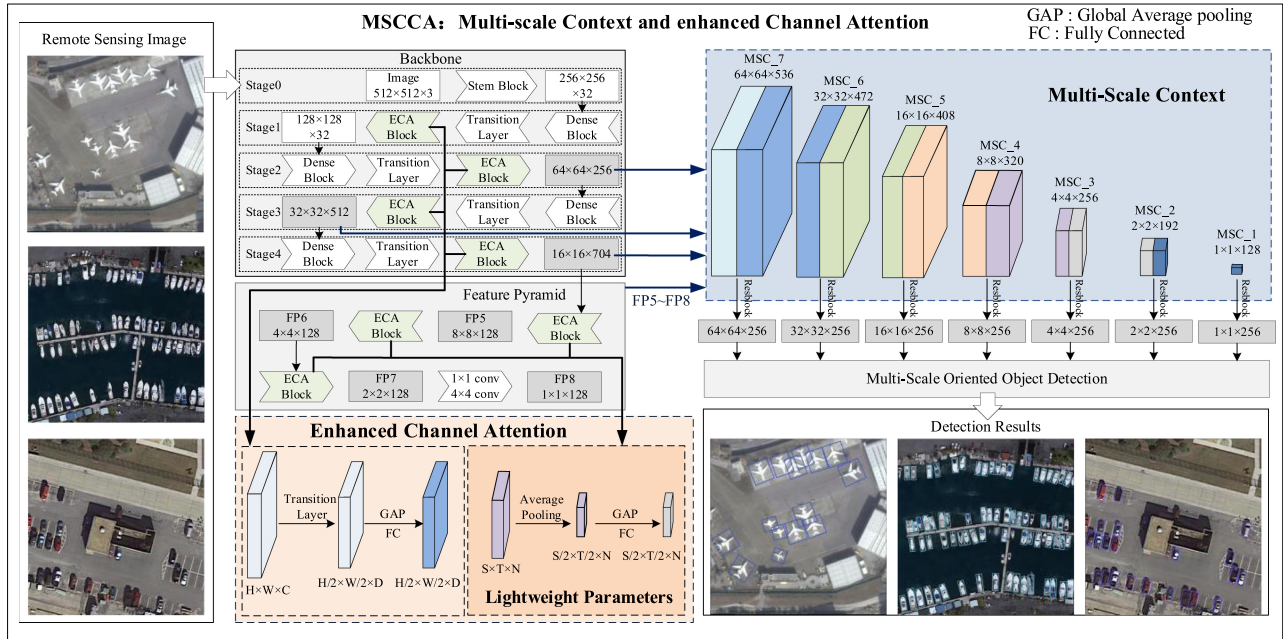


Fig. 1. Structure of MSCCA.

The remaining components consist of dense and transition layers. The dense layer can acquire receptive fields at multiple scales and consists of two branches, one of which employs one  $1 \times 1$  and one  $3 \times 3$  convolution layer, while the other uses one  $1 \times 1$  and two stacked  $3 \times 3$  convolution layers. The two branches are merged with the previous feature to the channel dimension via concat. The transition layer includes a  $1 \times 1$  convolution layer and a  $2 \times 2$  average pooling layer with a stride of 2.

### B. Enhanced Channel Attention

The attention mechanism in the CNN draws on the human visual attention mechanism. Human vision quickly scans a global image to obtain the required object area, generally referred to as the focus of attention. Additional attention is then focused on this area to obtain more detailed information about the target object, while suppressing other useless information. In general, some features learned in the CNN will be redundant for the object detection task [33]. For example, the Relu layer will generate a large number of parameters with a value of 0, while visualizing the intermediate feature image can demonstrate the inability of some channels to detect the object. Thus, during network training, some channels are more important than other channels. In order to emphasize these important channels, we include the channel attention structure ECA Block in the model (see Fig. 2) based on SE Block [34].

In ECA Block, for any given feature map  $X \in \mathbb{R}^{C \times H \times W}$ , the global average pooling layer is implemented to generate features  $M \in \mathbb{R}^{C \times 1 \times 1}$

$$M_c = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W X_c(i, j) \quad (1)$$

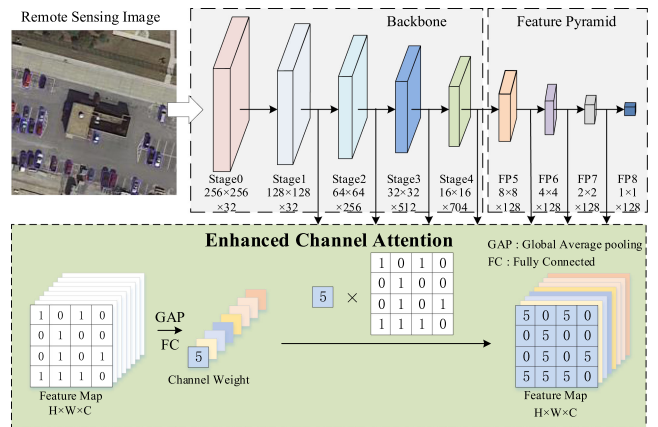


Fig. 2. Enhanced channel attention block.

where  $C$  is the number of channels;  $H$  and  $W$  are the height and width of the image, respectively;  $M_c$  indicates the feature obtained after the current channel performs global average pooling on  $X_c$ ;  $c$  is the current channel; and  $X_c(i, j)$  is the feature value of input feature image  $X_c$  under coordinates  $(i, j)$ . The feature image of each channel accumulates all the values, averages them to generate feature  $M_c$  and subsequently combines features  $M$  of  $C$  channels. Two fully connected (FC) layers are implemented. The first FC layer uses the Relu activation function to generate features of  $C/r \times 1 \times 1$  size, where  $r$  is a hyperparameter and is used to change the ECA block parameter in the network. Here, we set  $r$  to 16 following the previous experience of SE block. The second FC layer uses the Sigmoid activation function to generate feature  $S \in \mathbb{R}^{C \times 1 \times 1}$

$$S = \sigma(W_2 \delta(W_1 M)) \quad (2)$$



TABLE I  
 BACKBONE

Stage	Layer	Feature Map	
Stage0	Input	512×512×3	
	Stem Block	256×256×32	
Stage 1	Dense Block	Dense Layer x 3	
	Transition Layer	1×1 Conv , stride 1 2×2 Ave pool ,stride 2	128×128×32
	ECA Block		
Stage 2	Dense Block	Dense Layer x 3	
	Transition Layer	1×1 Conv , stride 1 2×2 Ave pool ,stride 2	64×64×256
	ECA Block		
Stage 3	Dense Block	Dense Layer x 3	
	Transition Layer	1×1 Conv , stride 1 2×2 Ave pool ,stride 2	32×32×512
	ECA Block		
Stage 4	Dense Block	Dense Layer x 3	
	Transition Layer	1×1 Conv , stride 1 2×2 Ave pool ,stride 2	16×16×704
	ECA Block		

where  $W_1$  and  $W_2$  represent two FC operations;  $\delta$  and  $\sigma$  are two activation functions;  $S$  is the generated feature and represents the importance of each feature channel following feature selection. The normalized weight is multiplied to the feature of each channel to output feature  $U \in \mathbb{R}^{C \times H \times W}$

$$U_c = F_{scale}(X_c, S_c) = S_c X_c \quad (3)$$

where  $U_c = [U_1, U_2, \dots, U_c]$  represents the feature generated following the scale operation for current channel  $c$ . The scale operation multiplies each element in  $S_c$  and  $X_c$  to generate feature  $U$  for each channel.

We add the ECA block to the proposed network to enhance the channel attention. The ECA Block is a simplified structure, which consists of a global average pooling operation, two full connections layers and a scale operation. Therefore, ECA block can be used to replace the complex convolution component of the network in order to reduce the number of network parameters. For example, after replacing the additional convolution layer with ECA block, the amount of network parameters is reduced from 7.06 to 5.08 M. Our results demonstrate that including the ECA Block can generally improve the detection accuracy and reduce the number of parameters (see Section IV-C).

### C. Multiscale Context

The CNN in object detection is associated with a high shallow network resolution and low deep network resolution. Shallow convolution features represent the details of the object, while deep convolution features indicate the semantic information. However, using multiscale feature maps for object detection ignores the detailed features in the shallow convolution features. Such shallow convolution features play a vital role in the detection of small objects. In order to fuse the scale context information [35], [36], we include the FPN-based SC structure to the network. In Table II, we added convolutional layers to the end of the backbone to extract low-scale feature maps.

MSCCA employs feature maps of different sizes to independently detect objects of varying sizes. In our proposed framework, the pyramid is constructed via bottom-up and top-down

 TABLE II  
 FEATURE PYRAMID

Feature Pyramid	Layer	Feature Map
FP5	1×1 conv , stride 1	8×8×256
	3×3 conv , stride 2	
FP6	1×1 conv , stride 1	4×4×256
	3×3 conv , stride 2	
FP7	1×1 conv , stride 1	2×2×256
	3×3 conv , stride 2	
FP8	1×1 conv , stride 1	1×1×256
	4×4 conv , stride 1	

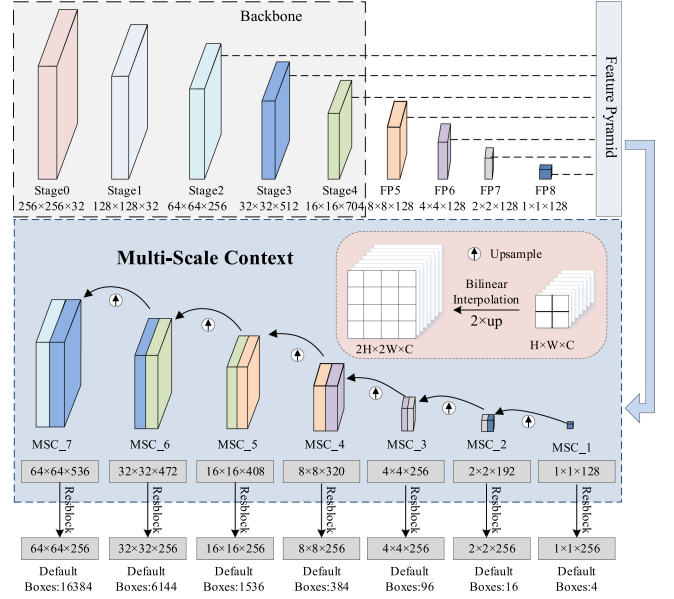


Fig. 3. Multiscale context block.

pathways, and lateral connections (see Fig. 3). For every scale feature image (with the exception of the highest level), we upsample the spatial resolution by a factor of 2 (via bilinear interpolation upsampling) and merge with the same sized feature image convolved by  $1 \times 1$ . Feature maps of other sizes undergo the same procedure until a new feature pyramid is generated. Feature maps that fully integrate the scale context information are then adopted to detect objects of different scales

$$U = [F_{\text{upsample}}(X) \oplus S]. \quad (4)$$

For each feature layer  $X$  of the pyramid,  $X \in \mathbb{R}^{C \times H \times W}$ .  $X$  is upsampled to twice the scale and fused with  $S$  in the channel dimension to generate feature  $U \in \mathbb{R}^{O \times H \times W}$ .  $S \in \mathbb{R}^{L \times H \times W}$  is a feature of the same scale as  $X$ . Fusing the features via the concat operation can make an excessively large feature dimension. Thus, we reduced the number of channels.

### D. Loss Function

Remote sensing images typically exhibit arbitrary object orientations. Thus, MSCCA employs quadrilateral bounding boxes to detect objects across different directions. The location information of the bounding box is expressed as  $(x, y, w, h)$ ,

whereby  $(x, y)$  represents the center point coordinates of the bounding box, and  $w$  and  $h$  are the width and height of the bounding box, respectively. If we define the default box as  $b = (x, y, w, h)$ , then the corresponding quadrilateral is represented as  $q = (x_1, y_1, x_2, y_2, x_3, y_3, x_4, y_4)$ , whereby  $(x_i, y_i)$  are the coordinates of the four vertices of the quadrilateral frame.

The loss function is divided into the confidence loss  $L_{\text{conf}}$  and the location loss  $L_{\text{loc}}$

$$L(x, c, l, g) = \frac{1}{N} (L_{\text{conf}}(x, c) + \alpha L_{\text{loc}}(x, l, g)) \quad (5)$$

where  $\alpha$  is the weight;  $N$  is the number of matched default boxes;  $x \in \{1, 0\}$  is the matching value indicating whether the default box matches the ground truth;  $c$  is the confidence;  $l$  is the predicted bounding box; and  $g$  is the ground truth. Positioning loss  $L_{\text{loc}}$  is a smooth L1 loss between the predicted bounding box and ground truth. If the overlap between the default box and ground truth exceeds the threshold (0.5), then it is considered as a positive sample

$$L_{\text{loc}}(x, l, g) = \sum_{i \in \text{Pos}} \sum_m x_{ij}^k \text{smooth}_{L1}(l_i^m - \hat{g}_j^m) \quad (6)$$

where  $m \in \{x, y, w, h, x_1, y_1, x_2, y_2, x_3, y_3, x_4, y_4\}$ ,  $(x, y)$  represents the center coordinates of the box;  $x_{ij}^k \in \{1, 0\}$  is an indicator of the match between the  $i$ th predicted bounding box and the  $j$ th ground truth;  $\text{Pos}$  is a positive sample;  $k$  is a ground truth object category; and  $\hat{g}_j^m$  represents the coded ground truth, which ensures that the weight of the ground truth center position and weakens the width and height widths

$$\hat{g}_j^x = (g_j^x - d_i^x)/d_i^w \quad \hat{g}_j^y = (g_j^y - d_i^y)/d_i^h \quad (7)$$

$$\hat{g}_j^w = \log\left(\frac{g_j^w}{d_i^w}\right) \quad \hat{g}_j^h = \log\left(\frac{g_j^h}{d_i^h}\right) \quad (8)$$

$$\hat{g}_j^{x_1} = (g_j^{x_1} - d_i^{x_{\min}})/d_i^w \quad \hat{g}_j^{y_1} = (g_j^{y_1} - d_i^{y_{\min}})/d_i^h \quad (9)$$

$$\hat{g}_j^{x_2} = (g_j^{x_2} - d_i^{x_{\max}})/d_i^w \quad \hat{g}_j^{y_2} = (g_j^{y_2} - d_i^{y_{\min}})/d_i^h \quad (10)$$

$$\hat{g}_j^{x_3} = (g_j^{x_3} - d_i^{x_{\max}})/d_i^w \quad \hat{g}_j^{y_3} = (g_j^{y_3} - d_i^{y_{\max}})/d_i^h \quad (11)$$

$$\hat{g}_j^{x_4} = (g_j^{x_4} - d_i^{x_{\min}})/d_i^w \quad \hat{g}_j^{y_4} = (g_j^{y_4} - d_i^{y_{\max}})/d_i^h \quad (12)$$

where  $d$  represents the default box;  $d_i^w$  and  $d_i^h$  are the width and height of the default box, respectively;  $(x_{\min}, y_{\min})$  and  $(x_{\max}, y_{\max})$  represent the coordinates of the upper left and lower right points of the horizontal default box, respectively; and is the  $\text{smooth}_{L1}$  loss defined as

$$\text{Smooth}_{L1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases} \quad (13)$$

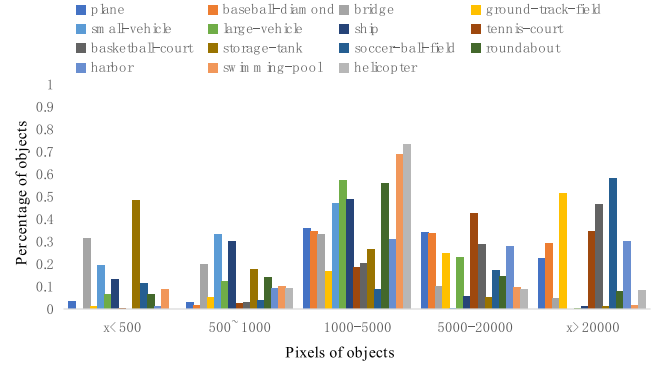


Fig. 4. Area distribution map of the objects contained in the DOTA dataset.

Confidence loss  $L_{\text{conf}}$  is described in formula (17) and can be divided into the cross-entropy loss of the positive and negative samples

$$L_{\text{conf}}(x, c) = - \sum_{i \in \text{pos}} x_{ij}^p \log(\hat{c}_i^p) - \sum_{i \in \text{neg}} \log(\hat{c}_i^0) \quad (14)$$

where  $\hat{c}_i^p = \frac{\exp(c_i^p)}{\sum_p \exp(c_i^p)}$

where  $c_i^p$  is the multcategory output; confidence  $\hat{c}_i^p$  is obtained following the activation of the Softmax function;  $p$  represents the  $p$ th category; the 0th category is the background;  $i$  is the  $i$ th predicted bounding box; and  $\text{Pos}$  and  $\text{Neg}$  indicate positive and negative samples, respectively. In order to ensure a balance, the ratio of the positive to negative sample is set to 3:1.

### III. DATASETS

#### A. DOTA

The DOTA dataset [37] was published on CVPR by Wuhan University. DOTA is a large-scale dataset used for the object detection of aerial images. It contains 2806 aerial images from different sensors and platforms. The images in the DOTA-v1.0 dataset were collected from Google Earth, some of which were taken by the satellite JL-1, and others were taken by the satellite GF-2 of the China Resources Satellite Data and Application Center. The size of each image ranges from approximately  $800 \times 800$  to  $4000 \times 4000$  pixels and contains objects of various proportions, orientations, and shapes. Current object detection methods generally divide small objects into two categories: 1) objects smaller than  $32 \times 32$  pixels; and 2) objects with a width and height less than one-tenth of the original image. Fig. 4 presents the area distribution of all object types in the DOTA dataset, where the horizontal axis represents the object pixel area size and the vertical axis is the percentage of each category in a certain scale range. The DOTA dataset contains a large number of small objects, the majority of which are aircrafts, cars, and boats. The objects are divided into the following 15 categories: plane, ship, storage tank, baseball field, tennis court, basketball court, ground track field, harbor, bridge, large vehicle, small

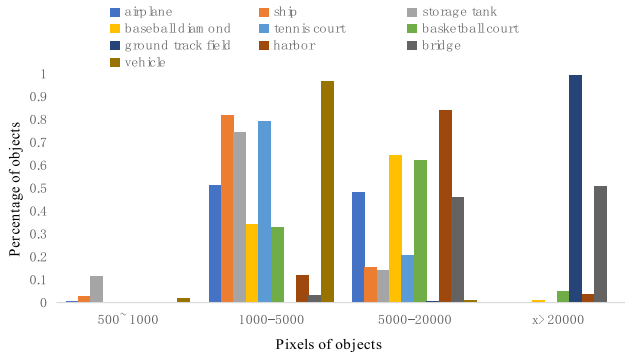


Fig. 5. Area distribution map of the objects contained in the NWPU VHR-10 dataset.

vehicle, helicopter, roundabout, soccer ball field, and swimming pool.

Due to the large number of pictures in the dataset and the large scale changes, we crop the pictures to a size of  $512 \times 512$  pixels and randomly select 3/5 of the samples as the training set, 1/5 as the verification set, and 1/5 as the test set.

### B. NWPU VHR-10

The NWPU VHR-10 dataset [38] is derived from a 10-level geographic remote sensing dataset for space object detection. The dataset includes 650 images containing objects and 150 background images. The image content and object types/characteristics are similar to those of the DOTA dataset (see Fig. 5). Although the dataset contains many object types, the number of samples is small, and the number and proportion of small objects is much less than that of the DOTA dataset. In particular, the NWPU VHR-10 dataset has almost no objects with an area of less than 1000 pixels. The 10 types of objects are: airplane, ship, storage tank, baseball field, tennis court, basketball court, ground track field, harbor, bridge, and vehicle.

## IV. RESULTS

### A. DOTA Dataset Results

The multiscale object detection method generates candidate regions of different scales on the feature maps, which are of different receptive field sizes and the size of the default box is based on these receptive field sizes. The default box setting contains two features: the scale and aspect ratio. The scale of each feature image default box is set as follows:

$$S_k = S_{\min} + \frac{S_{\max} - S_{\min}}{m - 1} (k - 1), k \in [1, m] \quad (15)$$

where  $S_k$  is the scale of the default box to the image;  $S_{\min}$  and  $S_{\max}$  represent the ratio of the lowest and highest scales, set to 0.15 and 0.9, respectively; and  $m$  is the number of feature maps of different sizes. Once the default box scale  $S_k$  of each feature image layer is determined, the specific default box is calculated according to the predefined aspect ratio. When the aspect ratio is 1, the side lengths of the two square default boxes are equal to  $S_k$  and  $S'_k = \sqrt{S_k S_{k+1}}$ , where  $S_{k+1}$  is the default box scale

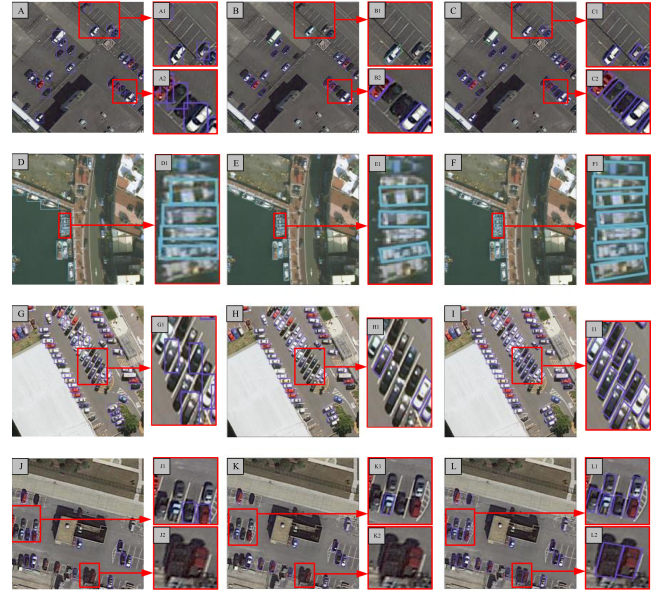


Fig. 6. DOTA results of SSD (left), Pelee (middle), and MSCCA (right).

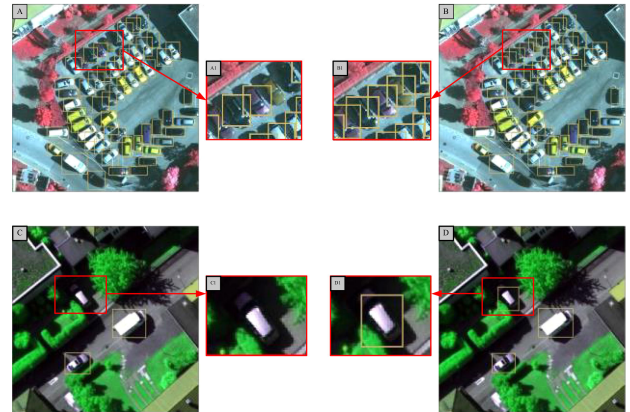


Fig. 7. NWPU VHR-10 results of Pelee (left) and MSCCA (right).

of the feature image in the next layer. If the aspect ratio does not equal 1, the default box is calculated as follows:

$$w_k^a = S_k \sqrt{a} h_k^a = S_k / \sqrt{a} \quad (16)$$

where  $w_k^a$  and  $h_k^a$  are the width and height of the candidate region of the  $k$ th feature image; and  $a$  is the value of the aspect ratio.

For the size of input image is  $512 \times 512$ , we select seven feature image scales to cover the different object sizes, as same as SSD [15], DSSD [42], FSSD [52], Rainbow SSD [53], Pelee [32], etc. The default box aspect ratios set to  $[[1, 2, 1/2], [1, 2, 3, 1/2, 1/3], [1, 2, 3, 1/2, 1/3], [1, 2, 3, 1/2, 1/3], [1, 2, 3, 1/2, 1/3], [1, 2, 1/2], \text{ and } [1, 2, 1/2]]$ .

During training, the pretrained model is employed to initialize the parameters. The learning rate is set at 0.005 for the first 120 000 iterations and is subsequently reduce by an order of magnitude after every 40 000 iterations computation, with 200 000 iterations in total. The momentum, weight decay, and



TABLE III  
DOTA DATASET DETECTION RESULTS

Method	PL	BD	BR	GTF	SV	LV	SH	TC	BC	ST	SBF	RA	HA	SP	HC	mAP	Fps
SSD[15]	41.0	24.3	4.55	17.1	15.9	7.72	13.2	39.9	12.0	46.8	9.09	30.8	1.36	3.50	0	17.8	59
YOLOv2[39]	76.9	33.8	22.7	34.8	38.7	32.0	52.3	61.6	48.5	33.9	29.2	36.8	36.4	38.2	11.6	39.2	30
RetinaNET[17]	78.2	53.4	26.3	42.2	63.6	52.6	73.1	87.1	44.6	57.9	18.0	51.0	43.3	56.5	7.4	50.3	14
R-FCN[40]	81.0	58.9	31.6	58.9	49.7	45.0	49.2	68.9	52.0	67.4	41.8	51.4	45.1	53.3	33.8	52.5	9
YOLOv3 [41]	79.0	77.1	33.9	68.1	52.8	52.2	49.8	89.9	74.8	59.2	55.5	49.0	61.5	55.9	41.7	60.0	13
DSSD[42]	<b>91.1</b>	71.8	54.6	66.4	79.0	77.2	87.5	87.6	52.1	69.7	38.0	72.6	75.4	59.4	28.9	67.4	9
DYOLO[43]	86.0	71.4	54.6	52.5	79.2	80.6	87.8	82.2	54.1	75.0	51.0	69.2	66.4	59.2	51.3	68.1	17
FPN[26]	88.7	75.1	52.6	59.2	69.4	78.8	84.5	90.6	81.3	82.6	52.5	62.1	76.7	66.3	60.1	72.0	6
FMSSD[16]	89.1	81.5	48.2	67.9	69.2	73.5	76.8	90.7	82.6	73.3	52.6	67.5	72.3	80.5	60.1	72.4	16
DRN[44]	89.7	82.3	47.2	64.1	76.2	74.4	85.8	90.5	86.1	84.8	57.6	61.9	69.3	69.6	58.4	73.2	9.8
R <sup>3</sup> Det [45]	89.4	81.1	50.5	66.1	70.9	78.6	78.2	<b>90.8</b>	85.2	84.2	61.8	63.7	68.1	69.8	67.1	73.7	10
SCRDET++[28]	90.0	84.3	55.4	73.9	77.5	71.1	86.0	90.6	87.3	87.0	<b>69.6</b>	68.9	73.7	71.2	65.0	76.8	13
FR-EST[46]	89.7	<b>85.2</b>	55.4	77.7	<b>80.2</b>	83.7	87.5	90.8	<b>87.6</b>	86.9	65.6	68.7	71.6	79.9	66.2	78.4	—
S <sup>2</sup> A-NET[22]	89.2	84.1	56.9	79.2	80.1	82.9	<b>89.2</b>	90.8	84.6	<b>87.6</b>	71.6	68.2	<b>78.5</b>	78.2	65.5	79.1	34
Pelee	87.6	72.9	52.8	73.7	73.5	77.9	76.3	90.0	80.7	74.4	40.7	68.0	71.7	79.6	83.7	74.0	29.2
MSCCA	89.7	84.9	<b>64.5</b>	<b>81.3</b>	77.3	<b>83.9</b>	84.8	90.4	86.2	77.1	54.8	<b>79.7</b>	78.0	<b>84.4</b>	<b>89.1</b>	<b>80.4</b>	<b>32.3</b>

TABLE IV  
NWPU VHR-10 DATASET DETECTION RESULTS

Method	PL	SH	ST	BD	TC	BC	GT	HA	BR	VH	mAP
RICNN[47]	88.3	77.3	85.2	88.1	40.8	58.4	86.7	68.6	61.5	71.1	72.6
COPD[48]	89.1	81.7	97.3	89.3	73.2	73.4	82.9	73.3	62.8	83.3	80.6
Faster R-CNN[10]	94.6	82.3	65.3	95.5	81.9	89.7	92.4	72.4	57.5	77.8	80.9
HyperNet[49]	99.4	89.7	<b>98.6</b>	90.9	90.6	90.3	89.2	80.3	68.9	88.6	88.7
Pelee	99.5	<b>93.4</b>	90.8	<b>97.2</b>	90.7	96.0	95.9	88.9	<b>88.9</b>	90.7	93.2
MSCCA	<b>99.7</b>	90.4	90.8	90.8	<b>90.8</b>	<b>98.6</b>	<b>98.3</b>	<b>90.3</b>	88.2	<b>98.3</b>	<b>94.4</b>

batch size are set to 0.9, 0.0005, and 16, respectively. The model is trained using the stochastic gradient descent method on four Nvidia Titan Xp GPUs.

Table III is the test results of the MSCCA model on the DOTA dataset, while Fig. 6 depicts the results of the model leaflet test. The result proves that MSCCA has higher detection accuracy than Pelee in detecting various objects. Pelee achieves a detection accuracy of 74% on the DOTA dataset, while that of the proposed MSCCA is 80.4%. This demonstrates the ability of the ECA Block and scale context features to improve the detection accuracy. The proposed MSCCA has a higher detection accuracy than S2A-NET.

The following are the detection results of SSD, Pelee, and MSCCA on DOTA. As shown in Fig. 6, the information in Pelee is not enough to detect the objects. The prediction result of MSCCA outperforms the Pelee by a large margin. And the boxes of objects are regressed more accurately.

### B. NWPU VHR-10 Dataset Results

The cropping, sample selection, and settings of the NWPU VHR-10 dataset [38] follow those of the DOTA dataset. Seven feature maps of different sizes are used, and six default boxes of varying ratios are generated for each pixel and scale feature layer. However, in contrast to the DOTA dataset, the NWPU VHR-10 dataset only contains a horizontal manual annotation box, and thus the results are maintained in the horizontal box.

The pretrained PeleeNet model is employed to initialize the parameters during training, with a 0.005 learning rate for the first 60 000 iterations that is subsequently reduced by an order of magnitude until the total 80 000 iterations are complete. The momentum, weight decay, and batch size are set to 0.9,

0.0005, and 16, respectively, and training is performed using the stochastic gradient descent method using four Nvidia Titan Xp GPUs.

Table IV is the detection results of the MSCCA model and other methods on the NWPU VHR-10 dataset. MSCCA outperforms the HyperNet by 5.7%. The following are the detection results of Pelee and MSCCA on NWPU VHR-10. Since there is no manual quadrilateral annotation in NWPU VHR-10 dataset, we use a horizontal default box to detect. As shown in Fig. 7, as same as the results on DOTA, the detection effect of MSCCA is better than that of Pelee.

### C. Ablation Study

In order to investigate the impact of the ECA Block and MSC structure on the detection results, we created several training models for the DOTA dataset and NWPU VHR-10 dataset to test using Nvidia Titan Xp and applied on Jetson TX2. We then evaluated the model size, detection speed, and computational complexity of the proposed method.

Tables V and VI are the impact of each structure in terms of the detection accuracy, parameter file size, and detection speed under a single Nvidia Titan Xp GPU. In Table V, without any structure, Pelee achieves a detection result of 74.0% mAP. Following the addition of the ECA block after each network stage, the accuracy improves to 75.8% mAP. This demonstrates the ability of the ECA block to strengthen the characterization performance of the network, thus improving the detection results. The inclusion of the MSC structure fusion scale context further improves the detection accuracy to 80.2% mAP. We then evaluate the impact of the ECA block, replacing the complex convolutional layer in the network. Replacing the ECA block

TABLE V  
INFLUENCE OF THE ECA BLOCK AND MSC STRUCTURE ON THE DETECTION PERFORMANCE FOR DOTA DATASET

Resblock-> ECA	Transition layers-> ECA	Additional layers-> ECA	Add ECA after each stage	MSC	mAP	Model size	Speed (Nvidia Titan Xp)	Parameters	FLOPs
					74.0	26.5 MB	29.2 fps	6.56M	4.58G
√			√	√	75.4	22.1 MB	32.2 fps	5.48M	3.93G
			√		75.8	26.9 MB	30.0 fps	6.67M	4.58G
	√			√	78.9	25.2 MB	33.2 fps	6.26M	4.43G
				√	79.2	28.2 MB	31.9 fps	7.00M	5.36G
			√	√	80.2	28.6 MB	31.9 fps	7.10M	5.36G
		√	√	√	80.4	20.6 MB	32.3 fps	5.12M	5.02G

TABLE VI  
INFLUENCE OF THE ECA BLOCK AND MSC STRUCTURE ON THE DETECTION PERFORMANCE FOR NWPU-VHR DATASET

Resblock-> ECA	Transition layers-> ECA	Additional layers-> ECA	Add ECA after each stage	MSC	mAP	Model size	Speed (Nvidia Titan Xp)	Parameters	FLOPs
					93.2	26.4 MB	29.1 fps	6.52M	4.55G
√			√	√	93.5	21.9 MB	32.0 fps	5.43M	3.91G
			√		93.8	26.8 MB	30.0 fps	6.63M	4.55G
	√			√	93.8	25.0 MB	32.8 fps	6.22M	4.40G
				√	94.0	28.2 MB	32.0 fps	6.95M	5.33G
			√	√	94.3	28.4 MB	32.0 fps	7.06M	5.33G
		√	√	√	94.4	20.5 MB	32.4 fps	5.08M	4.99G

with Resblock or the transition layer reduces the network parameters yet the detection accuracy is also weakened. Following this, we include a convolutional layer to provide small-scale feature maps for the multiscale detection framework. We use a pooling layer to replace the convolutional layer downsampling, and subsequently add the ECA Block to enhance the channel attention.

For lightweight network, flops, model parameters, and memory access cost (MAC) [50], [51] is widely used to measure the computational cost. Follow the design guide of lightweight network, in the proposed structure MSC and ECA block, we balanced the number of input and output channels for  $1 \times 1$  convolution and make their ratio approach 1:1. This operation has been proved to reduce the MAC of the network. With the addition of MSC structure, the parameters and flops of the model are increased, but the inference time of the model is accelerated. Compared with Pelee, this structure was able to achieve an mAP of 80.4% and 6.4% higher than Pelee. The model size reduced from 26.5 to 20.6 MB and detection speed increased from 30.0 to 32.3 fps. Thus, the MSCCA can be considered as a lightweight oriented object detection model in remote sensing images.

## V. CONCLUSION

A lightweight MSCCA model was proposed in this article. It employs PeleeNet as the backbone network. The feature image channel attention is enhanced and the multiscale context information is fused with multiscale detection methods to improve the characterization ability of the CNN. Results show that for  $512 \times 512$  input images, MSCCA was able to achieve 80.4% and 94.4% mAP on the DOTA and NWPU VHR-10, respectively. Meanwhile, the model size of MSCCA is 21% smaller than that of its predecessor. MSCCA can be considered as a practical lightweight oriented object detection model in remote sensing images. In the future, the proposed MSCCA model will be applied to edge devices for object detection application in remote sensing images. Moreover, computing optimization methods

(like TensorRT) will be used to improve the processing efficiency of model inference procedures.

## REFERENCES

- [1] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2005, pp. 886–893.
- [2] T. Malisiewicz, A. Gupta, and A. A. Efros, "Ensemble of exemplar-svms for object detection and beyond," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2011, pp. 89–96.
- [3] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 9, pp. 1627–1645, Sep. 2009.
- [4] H. Ghanbari, M. Mahdianpari, S. Homayouni, and F. Mohammadimanesh, "A meta-analysis of convolutional neural networks for remote sensing applications," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 14, pp. 3602–3613, Mar. 2021.
- [5] D. Hong *et al.*, "More diverse means better: Multimodal deep learning meets remote sensing imagery classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 5, pp. 4340–4354, May 2021.
- [6] D. Hong, N. Yokoya, G. S. Xia, J. Chanussot, and X. X. Zhu, "X-ModalNet: A semi-supervised deep cross-modal network for classification of remote sensing data," *ISPRS J. Photogramm. Remote Sens.*, vol. 167, pp. 12–23, Sep. 2020.
- [7] Z. Zou, Z. Shi, Y. Guo, and J. Ye, "Object detection in 20 years: A survey," 2018, *arXiv:1905.05055v1*.
- [8] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and se-mantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2014, pp. 580–587.
- [9] R. Girshick, "Fast R-CNN," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015, pp. 1440–1448.
- [10] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Adv. Neural Inf. Process. Syst.*, vol. 39, no. 6, pp. 91–99, Jun. 2016.
- [11] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-CNN," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017, pp. 2980–2988.
- [12] J. Pang, K. Chen, J. Shi, H. Feng, W. Ouyang, and D. Lin, "Libra R-CNN: Towards balanced learning for object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 821–830.
- [13] C. Xu, C. Li, Z. Cui, T. Zhang, and J. Yang, "Hierarchical semantic propagation for object detection in remote sensing imagery," *Proc. IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 6, pp. 4353–4364, Jun. 2020.
- [14] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 779–788.



- [15] W. Liu *et al.*, "SSD: Single shot multibox detector," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 21–37.
- [16] P. Wang, X. Sun, W. Diao, and K. Fu, "FMSSD: Feature-merged single-shot detection for multiscale objects in large-scale remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 9, pp. 3377–3390, Dec. 2019.
- [17] T. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 2, pp. 2999–3007, Jul. 2018.
- [18] S. Zhang, L. Wen, X. Bian, Z. Lei, and S. Z. Li, "Single-shot refinement neural network for object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 4203–4212.
- [19] Q. Zhao *et al.*, "M2Det: A single-shot object detector based on multi-level feature pyramid network," in *Proc. IEEE AAAI*, 2018, pp. 9259–9266.
- [20] Z. Tian, C. Shen, H. Chen, and T. He, "FCOS: Fully convolutional one-stage object detection," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017, pp. 9627–9636.
- [21] C. Wang, X. Bai, S. Wang, J. Zhou, and P. Ren, "Multiscale visual attention networks for object detection in VHR remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 16, no. 2, pp. 310–314, Feb. 2019.
- [22] Q. Ming, Z. Zhou, L. Miao, H. Zhang, and L. Li, "Dynamic anchor learning for arbitrary oriented object detection," 2020, *arXiv:2012.04150*.
- [23] B. Singh and L. S. Davis, "An analysis of scale invariance in object detection–snip," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 3578–3587.
- [24] B. Singh, M. Najibi, and L. S. Davis, "Sniper: Efficient multi-scale training," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 9333–9343.
- [25] T. M. Hoang, G. P. Nam, J. Cho, and I.-J. Kim, "DEFace: Deep efficient face network for small scale variations," *IEEE Access*, vol. 8, pp. 142423–142433, Jul. 2020.
- [26] T. Lin, P. Dollár, R. Girshick, K. He, and S. Belongie, "Feature pyramid networks for object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 936–944.
- [27] Y. Li, Y. Chen, N. Wang, and Z. Zhang, "Scale-aware trident networks for object detection," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 6053–6062.
- [28] X. Yang, J. Yan, X. Yang, J. Tang, W. Liao, and T. He, "SCRDet++: Detecting small, cluttered and rotated objects via instance-level feature denoising and rotation loss smoothing," 2020, *arXiv:2004.13316*.
- [29] Y. Chen *et al.*, "Stitcher: Feedback-driven data provider for object detection," *CVPR*, Apr. 2020.
- [30] J. Su, J. Liao, D. Gu, Z. Wang, and G. Cai, "Object detection in aerial images using a multiscale keypoint detection network," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 14, pp. 1389–1398, Dec. 2021.
- [31] Y. Li, Q. Huang, X. Pei, Y. Chen, L. Jiao, and R. Shang, "Cross-Layer attention network for small object detection in remote sensing imagery," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 14, pp. 2148–2161, Dec. 2021.
- [32] R. J. Wang, X. Li, and C. X. Ling, "Pele: A real-time object detection system on mobile devices," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 1967–1976.
- [33] J. Wang, Y. Zheng, M. Wang, Q. Shen, and J. Huang, "Object-Scale adaptive convolutional neural networks for high-spatial resolution remote sensing image classification," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 14, pp. 283–299, Dec. 2021.
- [34] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 8, pp. 2011–2023, Apr. 2019.
- [35] Z. Tian, W. Wang, R. Zhan, Z. He, J. Zhang, and Z. Zhuang, "Cascaded detection framework based on a novel backbone network and feature fusion," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 12, no. 9, pp. 3480–3491, Sep. 2019.
- [36] Z. Gao, H. Ji, T. Mei, B. Ramesh, and X. Liu, "EOVNet: Earth-observation image-based vehicle detection network," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 12, no. 9, pp. 3552–3561, Sep. 2019.
- [37] G. Xia *et al.*, "DOTA: A large-scale dataset for object detection in aerial images," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 3974–3983.
- [38] G. Cheng, J. Han, P. Zhou, and L. Guo, "Multi-class geospatial object detection and geographic image classification based on collection of part detectors," *ISPRS J. Photogramm. Remote Sens.*, vol. 98, no. 12, pp. 119–132, Dec. 2014.
- [39] J. Redmon and A. Farhadi, "YOLO9000: Better, faster, stronger," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 6517–6525.
- [40] Y. Li, K. He, and J. Sun, "R-fcn: Object detection via region-based fully convolutional networks," in *Proc. 30th Conf. Neural Inf. Process. Syst.*, 2016, pp. 1–9.
- [41] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," 2018, *arXiv:1804.02767*.
- [42] C. Y. Fu, W. Liu, A. Ranga, A. Tyagi, and A. C. Berg, "DSSD: Deconvolutional single shot detector," 2017.
- [43] O. Acatay, L. Sommer, A. Schumann, and J. Beyerer, "Comprehensive evaluation of deep learning based detection methods for vehicle detection in aerial imagery," in *Proc. IEEE Int. Conf. Adv. Video Signal Based Surveill.*, 2018, pp. 1–6.
- [44] X. Pan *et al.*, "Dynamic refinement network for oriented and densely packed object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 11204–11213.
- [45] X. Yang, J. Yan, Z. Feng and T. He, "R3det: Refined single-stage detector with feature refinement for rotating object," in *Proc. IEEE AAAI*, 2021, pp. 11207–11216.
- [46] K. Fu, Z. Chang, Y. Zhang, and X. Sun, "Point based estimator for arbitrary-oriented object detection in aerial images," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 5, pp. 4370–4387, May 2021.
- [47] G. Cheng, P. Zhou, and J. Han, "Learning rotation-invariant convolutional neural networks for object detection in VHR optical remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 12, pp. 7405–7415, Sep. 2016.
- [48] G. Cheng, J. Han, P. Zhou, and L. Guo, "Multi-class geospatial object detection and geographic image classification based on collection of part detectors," *ISPRS J. Photogramm. Remote Sens.*, vol. 98, no. 1, pp. 119–132, Apr. 2014.
- [49] T. Kong, A. Yao, Y. Chen, and F. Sun, "HyperNet: Towards accurate region proposal generation and joint object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 845–853, Jun. 2016.
- [50] Y. Lee, J. Hwang, S. Lee, Y. Bae, and J. Park, "An energy and GPU-Computation efficient backbone network for real-time object detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops*, 2019, pp. 752–760.
- [51] N. Ma, X. Zhang, H. T. Zheng, and J. Sun, "ShuffleNet V2: Practical guidelines for efficient CNN architecture design," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 122–138.
- [52] Z. Li and F. Zhou, "FSSD: Feature fusion single shot multibox detector," 2017, *arXiv:1712.00960*.
- [53] J. Jeong, H. Park, and N. Kwak, "Enhancement of SSD by concatenating feature maps for object detection," *BMVC*, 2017.



**Qiong Ran** received the Ph.D. degree from the Institute of Remote Sensing Applications, Chinese Academy of Sciences, Beijing, China, in 2009.

She has authored or coauthored more than ten papers in China and abroad. She is currently with the College of Information Science and Technology, Beijing University of Chemical Technology, Beijing, China. Her research interests include image acquisition, image processing, hyperspectral image analysis, and applications.



**Qing Wang** received the bachelor's degree from Qingdao Institute of Technology, China, in 2017. He is currently working toward the master's degree with College of Information Science and Technology, Beijing University of Chemical Technology, Beijing, China.



**Boya Zhao** received the B.Sc. degree from the School of Electrical Engineering and Information, Hebei University of Technology, Tianjin, China, in 2013, and the Ph.D. degree from the School of Electrical and Information Engineering, Beijing Institute of Technology, Beijing, China, in 2019.

He is currently an Assistant Professor with the Key Laboratory of Digital Earth Science, Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing, China. His research interests include the object detection in complex background and on-board real-time information processing.



**Shengliang Pu** (Student Member, IEEE) received the B.S. degree in geodetic engineering from the School of Geodesy and Geomatics, Wuhan University, Wuhan, China, in 2009, the M.S. degree in software engineering from the College of Computer Science, Inner Mongolia University, Hohhot, China, in 2013, and the Ph.D. degree in photogrammetry and remote sensing from the School of Geodesy and Geomatics, Wuhan University, in 2019.

He is with the Key Laboratory of Digital Earth Science, Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing, China. He is currently a Postdoc Researcher with the Key Laboratory of Digital Earth Science, Aerospace Information Research Institute, Chinese Academy of Sciences; and a Lecturer with the Faculty of Geomatics, East China University of Technology, Nanchang, China. He has a wide variety of research interests in remote sensing. His currently research is mainly focused on hyperspectral remote sensing image processing.



**Yuanfeng Wu** (Senior Member, IEEE) received the B.S. and M.S. degrees in computer science from China University of Mining and Technology, Beijing, China, and the Ph.D. degree in cartography and geographical information system from the Graduate University of Chinese Academy of Sciences, Beijing, China, in 2010.

He is currently an Associate Professor with the Key Laboratory of Digital Earth Science, Aerospace Information Research Institute, Chinese Academy of Sciences. His research interests include the development of onboard real-time algorithms, high-performance computing implementation, and computer software in hyperspectral image processing.



**Zijin Li** received the B.S degree from the School of Earth and Space Sciences, Peking University, Beijing, China, in 2021. She is currently working toward the M.S. degree at the Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing, China.

Her research interest includes hyperspectral image processing and object detection.