# AN EFFICIENT APPROACH FOR DETECTING TINY OBJECTS IN MASSIVE BACKGROUND BASED ON SPLIT-ATTENTION NETWORK

**HOANH NGUYEN**

Faculty of Electrical Engineering Technology, Industrial University of Ho Chi Minh City, Ho Chi Minh

City, Vietnam

E-mail: nguyenhoanh@iuh.edu.vn

## ABSTRACT

With the fast development of deep learning, object detection based on vision has achieved great progress in recent years. Though considerable progress has been made, there still exist challenges for objects with tiny size. One of the main reason is that feature representatives of tiny objects become sparse and weak due to their tiny size in an enormous background. This makes tiny objects difficult to be detected with state-of-the-art object detectors. This paper proposes an efficient method for detecting tiny objects in massive background. First, ResNest architecture is adopted as the backbone to extract features from input images. ResNest captures cross-channel feature correlations, while preserving independent representation in the meta structure. As a result, ResNest architecture achieves better speed-accuracy trade-offs than state-of-the-art deep CNN-based models without incurring excessive computational costs. Next, feature maps generated by the backbone are used to build feature pyramid following FPN network. Finally, this paper proposes an attention network in the detection part to solve problems of occlusion, noise, and blurring and effectively enhance the representations of tiny objects in complex backgrounds based on multi-dimensional attention network and inception module. Experiment results on the AI-TOD dataset show that the proposed method is very efficient in terms of the detection ability of very tiny and tiny objects.

**Keywords:** *Tiny Objects Detection, Convolutional Neural Network, Deep Learning, Object Detection, Split-Attention Network*

## 1. INTRODUCTION

In recent years, the fast development of convolutional neural networks (CNNs) has significantly accelerated the development of object detection. Basically, object detection approaches can be divided into two categories: one-stage approaches and two-stage approaches. Two-stage deep CNNs-based object detection framework usually includes feature extraction subnet, proposal extraction subnet, and detection subnet. The feature extraction subnet [1], [2] applies a convolutional neural network to extract features from input images. The proposal extraction subnet [3], [4] generates regions of interest (ROIs), including foreground positive samples and background negative samples from the feature map. The detection subnet [4], [5] utilizes the pooling feature of ROIs to predict the classification and regression results of detected objects.

Prior works usually feed the single scale features from the last convolution layer into the subsequent subnet to utilize feature information obtained from the feature extraction subnet [4], [6]. Although feature maps are rich in higher-level semantic information, they lack detailed information. This leads to a poor detection performance for small objects and occluded objects. To overcome this problem, some methods proposed network structures to predict separately on feature layers of different resolutions [7], [8] or to merge multi-resolution features at first and then to predict on the merged feature map [9]. Recent methods show that taking advantages of merge multi-resolution features and predicting separately on feature layers can get more accurate results [10], [11]. FPN is a popular network which combines multi-resolution features. FPN constructs a feature pyramid with high-level semantics throughout and
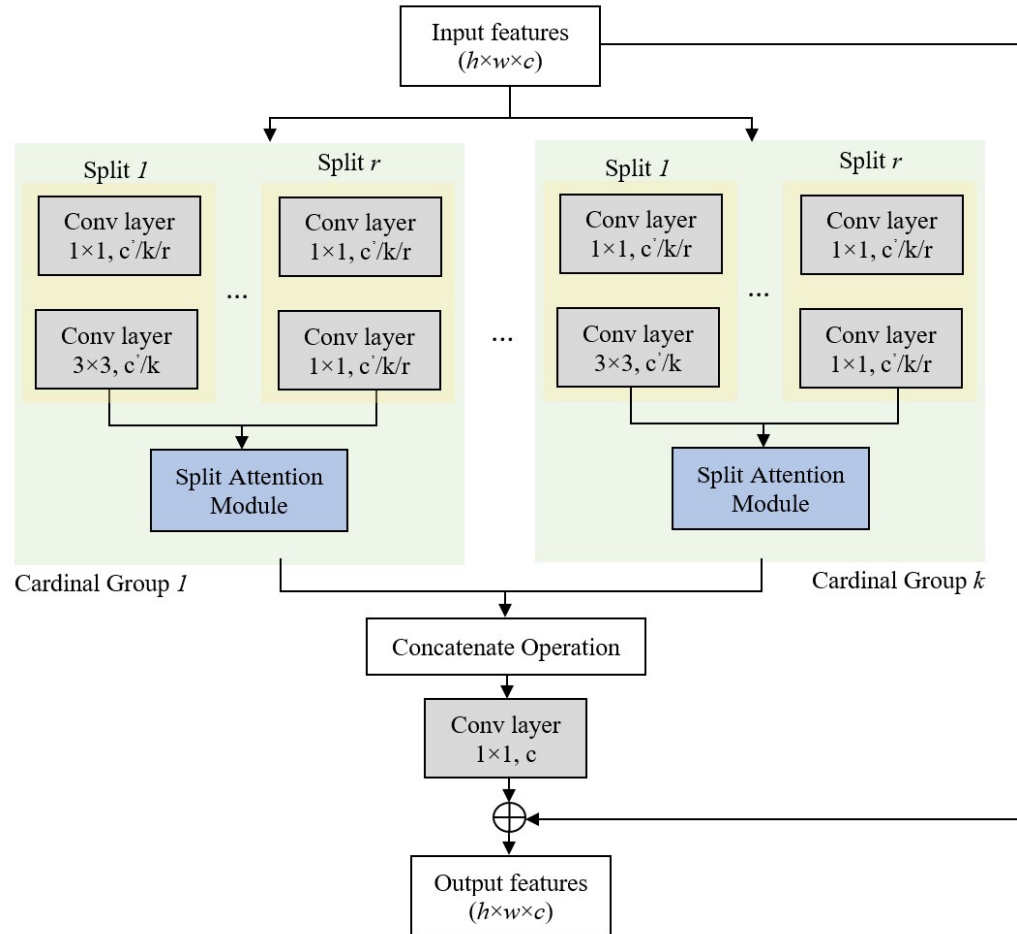
*Figure 1: The Structure of ResNest Block. The Input Feature Map Is First Divided into R Cardinal Groups. Split Attention Module Is Then Applied on Each Feature Group, and The Feature Maps of Each Cardinal Group Are Concatenated to Obtain The Output Feature Map; (h,w,c) Represents The Sizes of The Input Feature Map; (c') Represents The Channels of Inter Forward Blocks*

independently predicts at each pyramid level. It merges the semantically stronger feature maps from top-down with feature maps which are rich in detail localization information from the same bottom-up level.

Unlike objects in proper scales, detecting objects of tiny scale is much more challenging due to the extremely small size and low signal-to-noise ratio in aerial image [12]. For the CNNs-based object detection methods such as Faster R-CNN with ResNet-50, input images will be down sampled 16 times by pooling layers. Therefore, a number of tiny objects will be filtered out in the final feature map. Despite the fact that lots of methods tackled this problem [13], [14], [15], [16], there is still a large gap between current and the upper bound performances on tiny object detection.

To obtain better performance on tiny object detection, this paper proposes an efficient method for detecting tiny objects in massive background. In the proposed framework, ResNest architecture is first adopted as the backbone to extract features from input images. In addition, this paper proposes an attention network to solve problems of occlusion, noise, and blurring and effectively enhance the representations of tiny objects in complex backgrounds based on multi-dimensional attention network and inception module. According to the numerical results, the detection performance of the proposed approach is significantly better than that of state-of-the-art methods in terms of the detection ability of very tiny and tiny objects.

## 2. METHODOLOGY

This section presents the details of the proposed deep CNN-based framework for tiny object
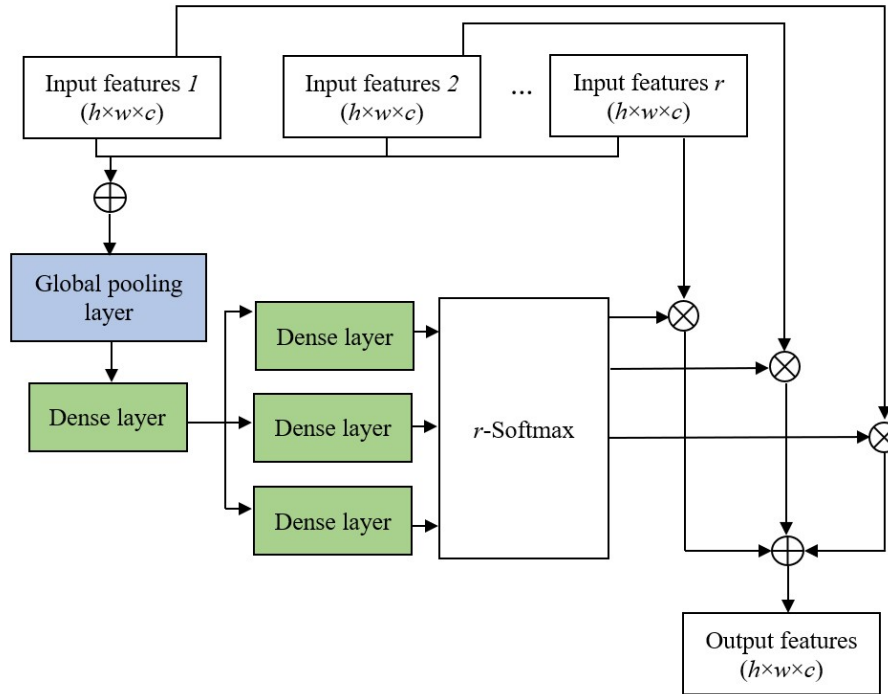
*Figure 2: The Structure of Split Attention Module*

detection in massive background. The proposed framework is based on the Faster R-CNN with FPN backbone.

## 2.1 Backbone

Since AlexNet [17], deep convolutional neural networks have become dominant in image classification, object detection segmentation. ResNet [2] introduced an identity skip connection which alleviates the difficulty of vanishing gradient in deep neural network and allows network to learn improved feature representations. ResNet has become one of the most successful CNN architectures which has been adopted in various computer vision applications. ResNest [18] proposed a simple architecture which combines the channel-wise attention strategy with multipath network layout. ResNest captured cross-channel feature correlations, while preserving independent representation in the meta structure. As a result, ResNest achieves better speed-accuracy trade-offs than state-of-the-art CNN models without incurring excessive computational costs. Inspired by the above networks, this paper adopts ResNest as the backbone network of the model. In the following parts of this section, this paper will introduce the structure of ResNest network in detail.

### 2.1.1 ResNest Block

Figure 1 shows the structure of ResNest block. In ResNest block, the feature map input into the ResNest block can be divided into several groups, and the number of feature groups $K$ is a hyperparameter. The resulting feature groups are regarded as cardinal groups. ResNest introduced a new radix hyperparameter $R$ that indicates the number of splits within a cardinal group, thus the total number of feature groups is $G = KR$. This paper sets $K$=2, $R$=2 as in the origin ResNest model. In each individual group, a series of transformations $\{F_1, F_2, ... F_G\}$ is applied and the intermediate representation of each group $U_i = F_i(X)$ is received. $F_i$ is a 1×1 convolution layer followed by a 3×3 convolution layer. Split attention modules, which is explained in the following section, are assigned to fuse feature maps in each split groups. The output feature maps from each cardinal group are then concatenated along the channel dimension as follow:

$$V = Concat\{V^1, V^2, ... V^K\} \qquad (1)$$

If the input and output feature map share the same shape, the final output feature map of the ResNest block is obtained by using a shortcut connection as follow:
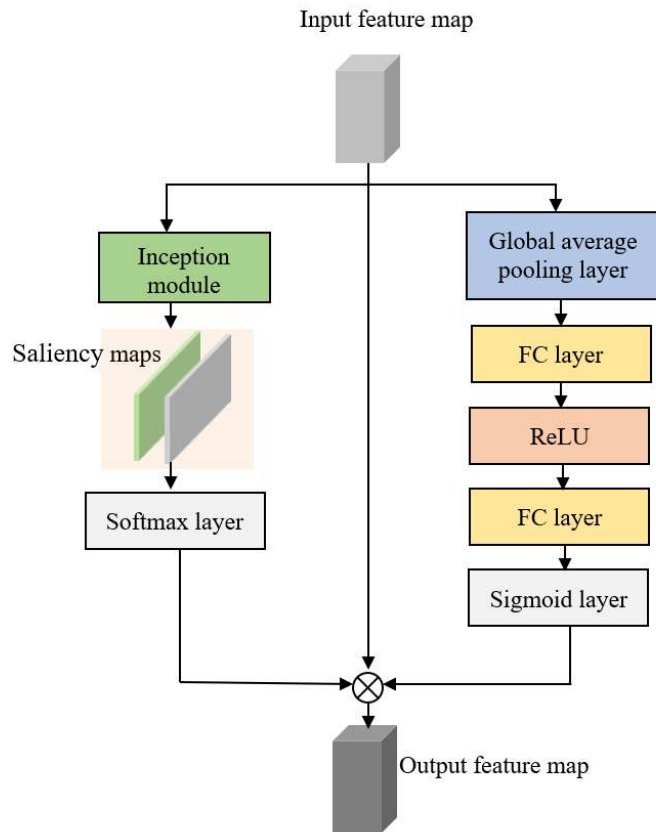
$$Y = V + X \qquad (2)$$

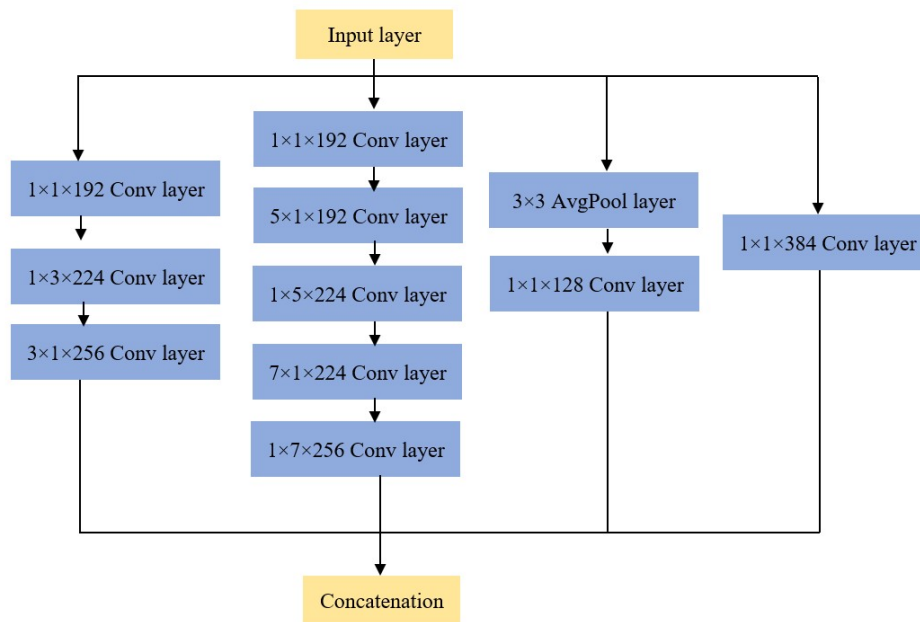*Figure 3: The Structure of The Proposed Attention Network*



*Figure 4: Architecture of The Inception Module*

For blocks with a stride, an appropriate transformation $\Psi$ is applied to the shortcut connection to align the output shapes, the final output feature map is produced as follow:

$$Y = V + \Psi(X) \qquad (3)$$

where $\Psi$ can be strided convolution or combined convolution-with-pooling.

**Split Attention Module**

In ResNest block, split attention module is designed to enables feature map attention across different feature map groups. Figure 2 illustrates the structure of the split attention module. As shown in Figure 2, the feature maps of each split group are first fused via element-wise summation across multiple splits. The feature map for $k_{th}$ cardinal group after fusing operation $\tilde{U}^k$ is calculated as follow:

$$\tilde{U}^k = \sum_{j=R(k-1)+1}^{Rk} \tilde{U}_j \qquad (4)$$

where $\tilde{U}^k \in \mathbb{R}^{H \times W \times C/K}$ for $k \in 1,2,\dots K$, and $(H,W,C)$ represents the block output feature map sizes.

The fused feature maps are then fed into a global average pooling across spatial dimensions $s^k \in \mathbb{R}^{C/K}$ to effectively collect global context information with embedded channel-wise statistics. Here the $c_{th}$ component is calculated as follow:

$$s_c^k = \frac{1}{H \times W} \sum_{i=1}^{H} \sum_{j=1}^{W} \tilde{U}_c^k(i,j) \qquad (5)$$

Two fully connected layers with BN and ReLU activation are used to generate each feature map channel as a weighted combination over splits. The $c_{th}$ channel $V_c^k$ is calculated as follow:

$$V_c^k = \sum_{i=1}^{R} a_i^k(c) U_{R(k-1)+i} \qquad (6)$$

$$a_i^k(c) = \begin{cases} \dfrac{\exp(\phi_i^c(s^k))}{\sum_{j=1}^{R} \exp(\phi_i^c(s^k))} & if\ R > 1, \\ \dfrac{1}{1+\exp(-\phi_i^c(s^k))} & if\ R = 1, \end{cases} \qquad (7)$$

where $a_i^k(c)$ denotes assignment weight, and $\phi_i^c$ determines the weight of each split for the $c_{th}$ channel based on the global context representation $s^k$.

### 2.1.2 ResNest Network

ResNest architecture is based on the ResNet-D model [19] and ResNest block as described in previous section. In addition to replace Residual block with ResNest block, ResNest model also adopts two effective modifications:

- The first 7×7 convolutional layer is replaced with three consecutive 3×3 convolutional layers, which have the same receptive field size with a similar computation cost as the original design.

- A 2×2 average pooling layer is added to the shortcut connection prior to the 1×1 convolutional layer for the transitioning blocks with stride of two.

In addition, instead of using strided convolution at the transitioning block, ResNest architecture uses an average pooling layer with a kernel size of 3×3. ResNest model captures cross-channel feature correlations, while preserving independent representation in the meta structure. ResNest block performs a set of transformations on low dimensional embeddings and concatenates their outputs as in a multi-path network. Each transformation incorporates channel-wise attention strategy to capture interdependencies of the feature map.

Due to the use of the pooling layer, small and tiny object lose most of their feature information in deep layers. It is a fact that low-level feature maps preserve location information of small and tiny objects, while high-level feature maps contain higher-level semantic information. Feature pyramid networks (FPN) [10] is a common feature fusion method that involves the combination of both high and low-level feature maps generated by the backbone network. Based on FPN, this paper also constructs the FPN with levels $P_2$ through $P_5$ to further fuse the feature maps generated by ResNest backbone.

### 2.2 Attention Network

In real scenarios, object proposals generated by the RPN may contain a large amount of noise information due to the complexity of environment. Uncontrolled noise information can overwhelm object information and blur the boundaries between the objects and backgrounds. These problems lead to missed detection and increasing false alarms. Therefore, it is necessary to enhance the object representations and weaken the background representations. This paper proposes an attention network to solve problems of occlusion, noise, and blurring and effectively enhance the representations of small objects in complex backgrounds based on multi-dimensional attention network [20]. Figure 3

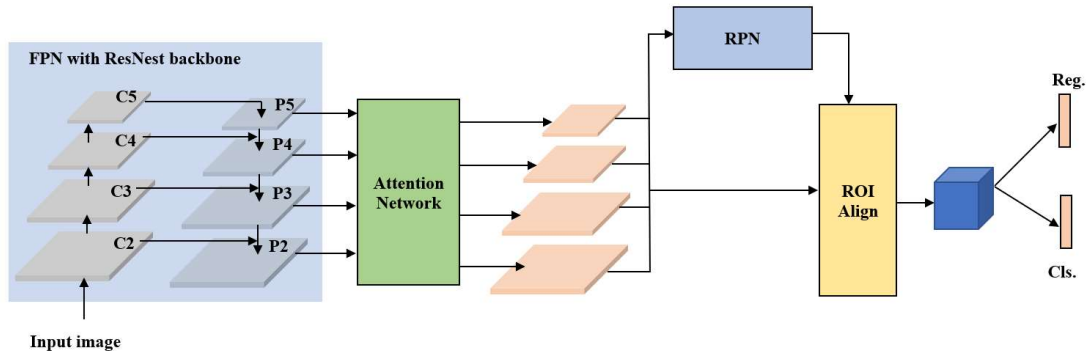illustrates the structure of the proposed attention network. In the pixel attention branch, each feature



*Figure 5: The Structure of The Proposed Model for Tiny Object Detection*

map generated by the FPN backbone is fed into an inception module to expand the receptive field and increase semantic information. The inception module, which is illustrated in Figure 4, contains a variety of ratio convolution kernels to capture the diversity of object shapes. The two-channel saliency map generated by the inception module contains the scores of the foreground objects and background. Softmax operation is then applied on the two-channel saliency map to rescale the value of the saliency map between [0, 1]. This operation can reduce the noise and relatively enhance the object information, especially for small and tiny objects. In the channel attention branch, SENet [21] is used to squeeze global spatial information into a channel descriptor. In this network, global average pooling is adopted to generate channel-wise statistics. Finally, fused feature map is generated by multiplying different input feature maps, including the saliency map, input feature map and feature map generated by the channel attention network.

### 2.3 Proposed Network

Figure 5 shows the overall structure of the proposed model. Following FPN, ResNest model is used as the backbone network, and feature maps {C2, C3, C4, C5} generated by the backbone are adopted to generate feature pyramid {P2, P3, P4, P5}. The attention network is then applied to all different scale levels {P2, P3, P4, P5} generated by the feature pyramid to generate corresponding fused feature layers. These fused feature maps are then fed into the RPN to generate object proposals. As in FPN, this paper also assigns anchors of a single scale and multiple aspect ratios at each level of the fused feature maps. To be more specific, this paper defines the anchors to have areas of {32×32, 64×64, 128×28, 256×256} pixels and box ratios of {1:1, 1:2, 2:1}.

Thus, there are total 12 anchors over the feature pyramid. Since there are many ROIs heavily overlapping with each other, non-maximum suppression (NMS) algorithm is adopted to filter the number of ROIs before feeding them into the ROI align layer. This paper sets the intersection-over-union (IoU) threshold at 0.5 for NMS. Then, this paper assigns anchors training labels based on their IoU ratios with ground truth bounding boxes. To be more specific, if the anchor has IoU over 0.5 with any ground truth box, it will be set as positive anchor. In addition, anchors which have the highest IoU for each ground truth box will also be assigned as positive anchor. Otherwise, if anchors have IoU less than 0.3 with all ground truth boxes, they will be set as negative anchor. The parameters of the RPN are shared across all fused feature levels. The fused feature maps are also fed into the ROI align layer to generate fixed-size proposals. Fixed-size proposals are finally fed into the R-CNN subnet for final prediction.

### 3. RESULTS AND DISCUSSION

### 3.1 Dataset and Evaluation Metrics

This paper conducts all experiments on the AI-TOD dataset [22]. AI-TOD is a dataset for tiny object detection in aerial images. There are 700,621 annotated object instances of eight categories across 28,036 aerial images with sizes of 800x800 pixels in this dataset, including airplane, bridge, storage tank, ship, swimming-pool, vehicle, person, and windmill. Objects in AI-TOD appear in various sizes. The largest object in AI-TOD is smaller than 64 pixels, and 86% of objects are smaller than 16 pixels as shown in Figure 6. Objects are classified based on their size. Objects in the range 2 to 8 pixels are considered as very tiny, 8 to 16 pixels as tiny, 16 to

32 as small, 32 to 64 as medium, and no large objects. The percentages of very tiny, tiny, small and

medium objects in AI-TOD are 13:3%, 72:3%, 12:3% and 2:1%, respectively. For dataset splits, 2/5,



*Figure 6: Example Images in the AI-TOD Dataset*

*Table 1: Number of Object Instances in the AI-TOD Dataset*

| Dataset | Train Set | Validation Set | Test Set |
|---|---|---|---|
| vehicle | 248042 | 59904 | 306665 |
| person | 14126 | 3841 | 15443 |
| ship | 13539 | 3791 | 17633 |
| storage-tank | 5269 | 2477 | 5860 |
| bridge | 512 | 140 | 689 |
| airplane | 623 | 170 | 745 |
| swimming-pool | 293 | 34 | 292 |
| windmill | 176 | 67 | 290 |
| **Total** | 282580 | 70424 | 347617 |

1/10 and 1/2 of the images are used to form training set, validation set and test set. For each object category and image set, the number of object instances is illustrated in Table 1.

For evaluation metrics, this paper employs the Average Precision (AP) metric, which has been widely used to assess various detection algorithms. In addition, $AP_{vt}$, $AP_t$, $AP_s$, $AP_m$ denote APs for very tiny, tiny, small, medium scales, respectively as in [22].

**3.2 Main Results on AI-TOD Dataset**

This paper compares the proposed method with state-of-the-art object detectors, including Faster R-CNN [4], Cascade R-CNN [23], YOLOv3 [24], RetinaNet [25], SSD-512 [8], FCOS [26], and Grid R-CNN [27]. Faster R-CNN and Cascade R-CNN represent the anchor-based two-stage detectors. YOLOv3, RetinaNet, and SSD-512 represent the anchor-based one-stage detectors. FCOS represents the anchor-free center-based detectors. The detection results are shown in Table 2. The AP

obtained by the proposed model is 15.2, and the $AP_{vt}$, $AP_t$, $AP_s$, and $AP_m$ are 5.5, 16.4, 24.8, and 24.1, respectively. Compared with other state-of-the-art



*Figure 7: Examples of Detection Results of The Proposed Method on the AI-TOD Dataset*

models and generic object detectors, the accuracy of the proposed model has been significantly improved, and the AP, $AP_{vt}$, and $AP_t$ obtained by the proposed model are also the highest compared to other state-of-the-art models. The results show that the proposed method is very efficient in terms of the detection ability of very tiny and tiny objects. In terms of the detection ability of small and medium objects, Cascade R-CNN achieves the best detection accuracy, at 25.5 and 26.6, respectively. Since Cascade R-CNN introduced a multi-stage extension

of the R-CNN, where detector stages deeper into the cascade are sequentially more selective against close false positives to avoid the problems of overfitting at training and quality mismatch at inference, it enhances the detection performance of large objects in massive background. However, the proposed method achieves significant improvements on very tiny and tiny objects compared to Cascade R-CNN. In addition, most of the detectors in the benchmark can just obtain performance less than 3% for $AP_{vt}$, which actually cannot be applied in the real scenario

applications. The visualization of detection results of the proposed is shown in Figure 7. As can be seen from Figure 7, the proposed method can effectively detect very tiny and tiny objects in massive background.

*Table 2: Detection Results of Different Methods on the AI-TOD Dataset*

| Method | Backbone | AP | $AP_{vt}$ | $AP_t$ | $AP_s$ | $AP_m$ |
|---|---|---|---|---|---|---|
| Faster R-CNN | ResNet-50-FPN | 11.4 | 0.0 | 8.3 | 23.1 | 24.5 |
| Cascade R-CNN | ResNet-50-FPN | 13.8 | 0.0 | 10.6 | **25.5** | **26.6** |
| YOLOv3 | DarkNet-53 | 4.5 | 2.1 | 4.6 | 5.9 | 6.2 |
| RetinaNet | ResNet-50-FPN | 4.7 | 2.0 | 5.4 | 6.3 | 7.6 |
| SSD-512 | VGG-16 | 7.0 | 1.0 | 4.7 | 11.5 | 13.5 |
| FCOS | ResNet-50-FPN | 9.8 | 1.4 | 8.0 | 15.1 | 17.4 |
| Grid R-CNN | ResNet-50-FPN | 12.2 | 0.2 | 10.3 | 22.6 | 23.3 |
| **Proposed method** | **ResNest** | **15.2** | **5.5** | **16.4** | 24.8 | 24.1 |

## 4. CONCLUSIONS

Tiny object detection in aerial images remains a very challenging problem since tiny objects contain a small number of pixels and are easily confused with massive background. An efficient method for detecting tiny objects in massive background is proposed in this paper. In the proposed framework, ResNest architecture is first adopted as the backbone to extract features from input images. ResNest is a simple architecture which combines the channel-wise attention strategy with multipath network layout. ResNest captured cross-channel feature correlations, while preserving independent representation in the meta structure. As a result, ResNest achieves better speed-accuracy trade-offs than state-of-the-art CNN models without incurring excessive computational costs. In addition, this paper proposes an attention network to solve problems of occlusion, noise, and blurring and effectively enhance the representations of tiny objects in complex backgrounds based on multi-dimensional attention network and inception module. According to the numerical results, the detection performance of the proposed approach is significantly better than that of state-of-the-art methods in terms of the detection ability of very tiny and tiny objects.

**REFERENCES:**

[1] Simonyan, Karen, and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition." *arXiv preprint arXiv:1409.1556* (2014).

[2] He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep residual learning for image recognition." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770-778. 2016.

[3] Uijlings, Jasper RR, Koen EA Van De Sande, Theo Gevers, and Arnold WM Smeulders. "Selective search for object recognition." *International journal of computer vision* 104, no. 2 (2013): 154-171.

[4] Ren, Shaoqing, Kaiming He, Ross Girshick, and Jian Sun. "Faster R-CNN: towards real-time object detection with region proposal networks." *IEEE transactions on pattern analysis and machine intelligence* 39, no. 6 (2016): 1137-1149.

[5] Li, Zeming, Chao Peng, Gang Yu, Xiangyu Zhang, Yangdong Deng, and Jian Sun. "Light-head r-cnn: In defense of two-stage object detector." *arXiv preprint arXiv:1711.07264* (2017).

[6]     Dai, Jifeng, Yi Li, Kaiming He, and Jian Sun. "R-fcn: Object detection via region-based fully convolutional networks." *arXiv preprint arXiv:1605.06409* (2016).

[7]     Cai, Zhaowei, Quanfu Fan, Rogerio S. Feris, and Nuno Vasconcelos. "A unified multi-scale deep convolutional neural network for fast object detection." In *European conference on computer vision*, pp. 354-370. Springer, Cham, 2016.

[8]     Liu, Wei, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C. Berg. "Ssd: Single shot multibox detector." In *European conference on computer vision*, pp. 21-37. Springer, Cham, 2016.

[9]     Kong, Tao, Anbang Yao, Yurong Chen, and Fuchun Sun. "Hypernet: Towards accurate region proposal generation and joint object detection." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 845-853. 2016.

[10]    Lin, Tsung-Yi, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. "Feature pyramid networks for object detection." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2117-2125. 2017.

[11]    Fu, Cheng-Yang, Wei Liu, Ananth Ranga, Ambrish Tyagi, and Alexander C. Berg. "Dssd: Deconvolutional single shot detector." *arXiv preprint arXiv:1701.06659* (2017).

[12]    Yu, Xuehui, Yuqi Gong, Nan Jiang, Qixiang Ye, and Zhenjun Han. "Scale match for tiny person detection." In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1257-1265. 2020.

[13]    Li, Yanghao, Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. "Scale-aware trident networks for object detection." In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6054-6063. 2019.

[14]    Bai, Yancheng, Yongqiang Zhang, Mingli Ding, and Bernard Ghanem. "Sod-mtgan: Small object detection via multi-task generative adversarial network." In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 206-221. 2018.

[15]    Noh, Junhyug, Wonho Bae, Wonhee Lee, Jinhwan Seo, and Gunhee Kim. "Better to follow, follow to be better: Towards precise supervision of feature super-resolution for small object detection." In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9725-9734. 2019.

[16]    NGUYEN, HOANH. "Improvement Of Detecting Small-Sized Traffic Signs Based On Deep Learning." *Journal of Theoretical and Applied Information Technology* 97, no. 19 (2019).

[17]    Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *Advances in neural information processing systems* 25 (2012): 1097-1105.

[18]    Zhang, Hang, Chongruo Wu, Zhongyue Zhang, Yi Zhu, Haibin Lin, Zhi Zhang, Yue Sun et al. "Resnest: Split-attention networks." *arXiv preprint arXiv:2004.08955* (2020).

[19]    He, Tong, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. "Bag of tricks for image classification with convolutional neural networks." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 558-567. 2019.

[20]    Yang, Xue, Jirui Yang, Junchi Yan, Yue Zhang, Tengfei Zhang, Zhi Guo, Xian Sun, and Kun Fu. "Scrdet: Towards more robust detection for small, cluttered and rotated objects." In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8232-8241. 2019.

[21]    Hu, Jie, Li Shen, and Gang Sun. "Squeeze-and-excitation networks." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7132-7141. 2018.

[22]    Wang, Jinwang, Wen Yang, Haowen Guo, Ruixiang Zhang, and Gui-Song Xia. "Tiny Object Detection in Aerial Images." In *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 3791-3798. IEEE, 2021.

[23]    Cai, Zhaowei, and Nuno Vasconcelos. "Cascade r-cnn: Delving into high quality object detection." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 6154-6162. 2018.

[24]    Redmon, Joseph, and Ali Farhadi. "Yolov3: An incremental improvement." *arXiv preprint arXiv:1804.02767* (2018).

[25]    Lin, Tsung-Yi, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. "Focal loss

for dense object detection." In *Proceedings of the IEEE international conference on computer vision*, pp. 2980-2988. 2017.

[26]   Tian, Zhi, Chunhua Shen, Hao Chen, and Tong He. "Fcos: Fully convolutional one-stage object detection." In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9627-9636. 2019.

[27]   Lu, Xin, Buyu Li, Yuxin Yue, Quanquan Li, and Junjie Yan. "Grid r-cnn." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7363-7372. 2019.