

RS-TinyNet: Stage-wise Feature Fusion Network for Detecting Tiny Objects in Remote Sensing Images

Xiaozheng Jiang, Wei Zhang, Xuerui Mao

Abstract—Detecting tiny objects in remote sensing (RS) imagery has been a long-standing challenge due to their extremely limited spatial information, weak feature representations, and dense distributions across complex backgrounds. Despite numerous efforts devoted, mainstream detectors still underperform in such scenarios. To bridge this gap, we introduce RS-TinyNet, a multi-stage feature fusion and enhancement model explicitly tailored for RS tiny object detection in various RS scenarios. RS-TinyNet comes with two novel designs: tiny object saliency modeling and feature integrity reconstruction. Guided by these principles, we design three step-wise feature enhancement modules. Among them, the multi-dimensional collaborative attention (MDCA) module employs multi-dimensional attention to enhance the saliency of tiny objects. Additionally, the auxiliary reversible branch (ARB) and a progressive fusion detection head (PFDH) module are introduced to preserve information flow and fuse multi-level features to bridge semantic gaps and retain structural detail. Comprehensive experiments on public RS dataset AI-TOD show that our RS-TinyNet surpasses existing state-of-the-art (SOTA) detectors by 4.0% AP and 6.5% AP₇₅. Evaluations on DIOR benchmark dataset further validate its superior detection performance in diverse RS scenarios. These results demonstrate that the proposed multi-stage feature fusion strategy offers an effective and practical solution for tiny object detection in complex RS environments.

Index Terms—Tiny Object Detection, remote sensing, multi-attention, feature-enhanced.

I. INTRODUCTION

TINY object detection from remote sensing (RS) images aims to accurately detect objects with tiny sizes and extremely low signal-to-noise ratio [1], [2]. It has been a focus of study in the RS field since it has broad real-world applications, including suspicious object surveillance, military reconnaissance, intelligent transportation, and precision agriculture [3], [4]. As shown in Fig. 1, tiny objects are widely present in RS images, which are mostly captured by satellites, featuring unique imaging conditions and extensive geographic coverage. Compared to small objects in natural scene images, those in RS data tend to be more blurred, densely packed, with fewer pixels, and lack clear structural details. Despite the advancements in deep learning for general object detection [5]–[8], the distinctive feature of tiny objects in RS imagery

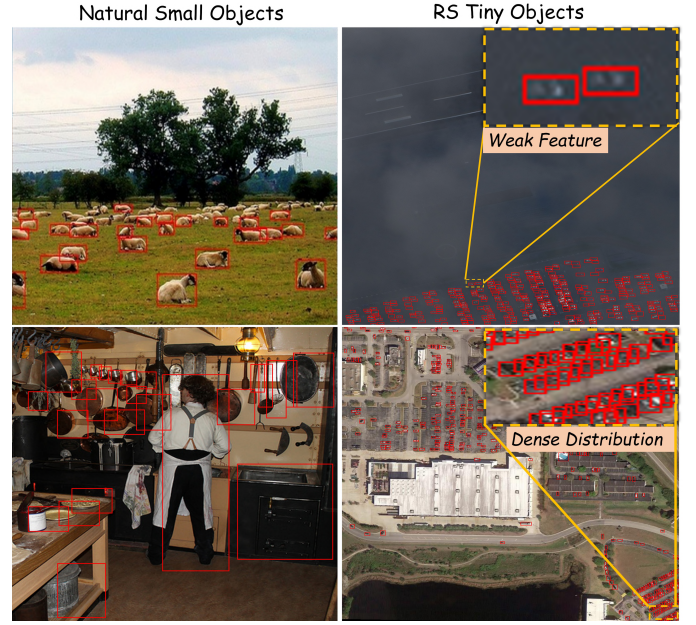


Fig. 1. Comparison of small and tiny objects in natural scene images (left) and RS imagery (right). Tiny objects in RS data exhibit extremely fewer pixels, higher object density, and quite weaker structural signals due to satellite-based imaging and wide-area coverage, posing greater challenges for detection.

continue to pose challenges to the generalization and precision of current popular models.

Most of the existing detection methods are designed for normal-sized objects in natural sensing and RS imagery. For example, SSD model [9], YOLO series [10]–[12], R-CNN [13], Fast R-CNN [14], and RS-specific models like S²ANet [15], Popeye [8], etc., which have achieved good performance in general object detection. Nevertheless, those mainstream detectors encounter difficulties in accurately localizing and classifying tiny objects that may appear less than 16×16 pixels. Notably, significant work has been focused on tiny object detection. Super-resolution-based models [16]–[18] incorporated super-resolution modules into the object detection pipeline, enhancing the visual details of tiny objects, whereas those modules significantly increase computational costs and severely degrade the applicability of the detector. In addition, Chen et al. [19] introduced a cross-stage multi-head attention module to enhance the quality of feature representations. Xu et al. [20] refined the label assignment process to mitigate the challenge of inaccurate positive and negative sample assignment for tiny objects in the training phase. However, these research efforts ignore the long-range and spatial context

Xiaozheng Jiang and Wei Zhang are co-first authors.

Corresponding author: Xuerui Mao.

Xiaozheng Jiang is with the School of Interdisciplinary Science, Beijing Institute of Technology, Beijing 100081, China. (e-mail: jxz15290303632@163.com).

Wei Zhang is with the School of Interdisciplinary Science, Beijing Institute of Technology, Beijing 100081, China. (e-mail: w.w.zhanger@gmail.com).

Xuerui Mao is with the School of Interdisciplinary Science, Beijing Institute of Technology, Beijing 100081, China, and also with Beijing Institute of Technology (Zhuhai), Zhuhai, 519088, China (e-mail: maouxuerui@sina.com).

information refinement for improving the feature representations for tiny objects. Moreover, most of them often sacrifice feature integrity and discriminative power in exchange for efficiency. Consequently, the insufficient exploitation of informative features continues to constrain the detection accuracy of tiny object detectors.

To address the above issues and better capture the characteristics of tiny objects, this paper proposes a novel RS tiny object detection framework, termed RS-TinyNet, which significantly improves detection accuracy through a multi-stage feature fusion and enhancement strategy. The design of the proposed model centers on two core principles: improving the saliency of tiny objects and preserving the integrity of feature representations. Specifically, a multi-dimensional collaborative attention (MDCA) mechanism is designed for efficient feature refinement through a two-branch co-design that establishes dynamic connections between channel-space dimensions and global-local features. Additionally, to address the challenges of feature degradation in RS tiny object detection, we introduce an auxiliary reversible branch (ARB) and a progressive fusion detection head (PFDH) to preserve feature integrity across the network [21], [22]. The ARB module alleviates the information fading problem in deep networks and enhances the integrity and reliability of the feature gradient flow [21]. Meanwhile, the PFDH progressively fuses features at different levels to minimize semantic gaps and reduce information loss [22]. Collectively, these innovations significantly improve the model's capacity to detect tiny objects in complex RS scenarios by enhancing attention-guided representation refinement and maintaining the integrity of semantic information.

Extensive experiments on RS tiny object dataset AI-TOD demonstrate that RS-TinyNet not only significantly outperforms existing mainstream algorithms in terms of detection accuracy, but also offers a new perspective for the tiny object detection field. Specifically, our method achieves a notable performance gain over the state-of-the-art (SOTA) detectors, with the overall AP improved by 4.0%, along with respective boosts of 2.1%, 4.2%, 3.7%, and 1.8% in AP_{vt} , AP_t , AP_s , and AP_m . These improvements highlight the model's ability to enhance feature representation and preserve integrity, especially for extremely tiny objects with limited pixels. The consistent gains across all scales further validate the robustness of RS-TinyNet in diverse RS scenarios.

The main contributions of this work are as follows.

- We propose RS-TinyNet, an effective framework tailored for tiny object detection in RS imagery. The design of RS-TinyNet is guided by two newly introduced principles: tiny object saliency modeling and feature integrity reconstruction. These principles drive a multi-stage feature enhancement strategy that improves the model's ability to extract discriminative representations while preserving critical structural information, thereby substantially boosting detection accuracy.
- We propose a stage-wise multi-level feature fusion and enhancement strategy to jointly improve tiny object saliency and preserve feature integrity across the network. A MDCA module is designed, based on the principle of tiny object saliency modeling, to synergistically integrate

channel and spatial information, as well as global and local contexts, thereby enabling the network to efficiently extract discriminative features of tiny objects. Moreover, driven by the goal of maintaining feature integrity, the ARB and PFDH modules are introduced to alleviate the loss of information during transmission and fusion of different layers of features, and therefore significantly improve the tiny object detection performance.

- To the best of our knowledge, RS-TinyNet is the first model that achieves the most significant performance improvement on multiple RS tiny object detection benchmarks. Compared with SOTA models, the AP metrics increased by more than 4.0% on the AI-TOD dataset. Moreover, RS-TinyNet performs well on other RS datasets such as DIOR, highlighting its robustness and practical applicability in diverse RS scenarios.

II. RELATED WORK

A. General Object Detection

Object detection is a fundamental task in computer vision, aiming to identify and localize all instances of objects within an image. General object detection algorithms in natural scenarios have witnessed significant progresses in recent years, with a range of classical approaches emerging. The R-CNN [13] series perform feature classification by extracting candidate regions, which is highly accurate but slow. Based on them, Fast R-CNN [14] and Faster R-CNN [5] significantly improve detection efficiency and accuracy through the implementation of ROI Pooling and Region Proposal Network (RPN). The YOLO series [10]–[12], [23]–[26] transform the detection problem into a regression problem. From YOLOv1 [10] to YOLOv12 [26], the model structure and training strategy are continuously optimised to achieve a balance between real-time performance and accuracy. The SSD [9] improves the model's adaptability to objects of different sizes by making predictions on multi-scale feature maps. Subsequent derivatives such as DSSD [27], RefineDet [28], etc., add decoding paths, feature resampling and other mechanisms to enhance the feature expression capability. On natural image datasets such as COCO [29], PASCALVOC [30], etc., these methods have generally demonstrated excellent performance and have become mainstream frameworks in the field of object detection.

In RS object detection, researchers have proposed specific detection methods for its characteristics such as variable viewpoint, wide coverage area, and dense objects. S²ANet [15] introduces an alignment convolution (AlignConv) to adaptively align features with rotated anchors, and a rotation-sensitive detection head to mitigate the inconsistency between classification and localization. The method significantly improves detection accuracy for densely distributed and arbitrarily oriented objects in aerial imagery. Chalavadi et al. [31] proposed a multi-scale object detection network (mSODANet) for aerial images that employs hierarchical dilated convolutions to effectively capture context at multiple scales. The network integrates a bidirectional feature aggregation module (BFAM) to fuse dense multiscale contextual information, enabling enhanced detection of objects with varying sizes and improving

the localization of small and sparse objects in complex aerial scenes. Gao et al. [32] proposed an attention-free global multiscale fusion network for RS object detection, which eliminates the traditional attention module and instead utilizes a global context modeling strategy to fuse multiscale features and overcome the challenges of complex backgrounds. This design reduces the computational overhead and maintains high detection accuracy, thus enabling efficient detection of densely distributed small objects in high-resolution RS images. Zhang et al. [8] proposed Popeye for multi-source ship detection in RS imagery. By fusing multi-source image representations and incorporating language-guided information, Popeye effectively improves the robustness and generalization ability of ship detection models under complex backgrounds and varying imaging conditions. Although these object detection methods in natural and RS images have seen significant success, they are primarily designed for normal-sized objects and struggle with detecting tiny objects effectively.

B. Tiny Object Detection

To address the difficulties associated with tiny object detection, researchers have explored various improvement strategies, mainly focusing on the directions of super-resolution-based methods, optimizing the training mechanism, and multiscale feature fusion. Rabbi et al. [16] introduced an end-to-end small-object detection framework that integrates an edge-enhanced GAN with a detection network, where a super-resolution module improves the visual details of tiny objects and improves the detection performance of remotely sensed images, but significantly increases the computational complexity. Chen et al. [19] introduced a RS tiny ship detection method based on degradation reconstruction enhancement, which combines image reconstruction with cross-stage multi-attention mechanism to effectively improve the feature quality and feature recognition of tiny objects. Wang et al. [33] proposed a tiny object detection framework with a normalized Wasserstein distance loss, which provides a more stable and scale-sensitive optimization objective for tiny object localization. Ge et al. [2] proposed a cross-attention based feature fusion enhancement network (CAF²ENet) to improve tiny object detection accuracy by refining the up-sampling results of deep features. Meanwhile, a regression-driven refocusing learning strategy is designed to improve the model's capacity to acquire high-quality tiny object detection frames. MARNet [34] improves multiscale feature fusion by integrating a global attention mechanism that captures channel context from deep feature maps and directs the enhancement of shallow features to highlight information in tiny object regions. The multi-branch feature pyramid network (MB-FPN) [35] presents a tiny object detection method that utilizes a global-local attention module, which integrates different levels of semantic information through a multi-scale feature fusion strategy, using global attention to capture contextual relationships and local attention to enhance the detailed representation of the tiny object region, which greatly enhances the ability to accurately identify tiny objects in complex contexts.

Despite various improvements, existing normal-sized/tiny object detectors still struggle with feature degradation and

under-utilization of spatial context in tiny objects. The lack of distinguishing information in tiny objects highlights the need for approaches that preserve feature integrity and enhance spatially-aware representation. To this end, we design RS-TinyNet to explicitly address these two key challenges.

III. METHODOLOGY

This section introduces the proposed RS-TinyNet framework. The model architecture is first summarized in Section III-A. Subsequently, three major improvements in RS-TinyNet are described: the MDCA module for tiny object saliency modeling (Section III-B) and the ARB and PFDH modules for feature integrity reconstruction (Section III-C).

A. Overview

In order to cope with the challenges of blurred features, fewer pixels, and missing detail structure in RS tiny object detection, this article constructs a multi-level feature fusion detection framework, RS-TinyNet. The overall framework of RS-TinyNet is shown in the top part of Fig. 2. The overall structure is based on YOLOv11 [25], which is enhanced by a multidimensional structure to achieve better detection performance. The key improvement of RS-TinyNet is reflected in three aspects. Firstly, the MDCA module fuses channel, spatial, local, and global information to capture the salient features of tiny objects. Then, the ARB module is proposed to reconstruct the feature connectivity of feedforward networks, and the PFDH module to hierarchically and progressively fuse different layers of features, so as to alleviate the information loss during feature transfer and fusion and improve the overall detection accuracy. Concurrently, the proposed step-wise feature fusion and enhancement strategy is implemented across the backbone, neck, and detection head of the network. At each stage, complementary mechanisms are employed to progressively refine multi-level feature representations, enhance salient cues, and preserve information integrity, thereby enabling robust tiny object detection under complex RS conditions. The designed step-wise feature enhancement core components are elaborated in detail as follows.

B. Tiny Object Saliency Modeling

Tiny objects in RS imagery exhibit ambiguous structures and significant scale variations, rendering single-scale or channel-only attention mechanisms insufficient for discriminative feature extraction. While conventional attention approaches emphasize channel dependencies, they often neglect precise spatial localization, which is crucial for tiny object detection. To address this, we design the core module MDCA of RS-TinyNet, which aims to break the limitation of separate channel and space modeling, and realize the dynamic aggregation of multi-dimensional features through the dual collaborative structure, enabling more effective representation of salient object cues and leading to improved detection performance. The detailed structure of MDCA is presented in the bottom left of Fig. 2, comprising two parallel branches.

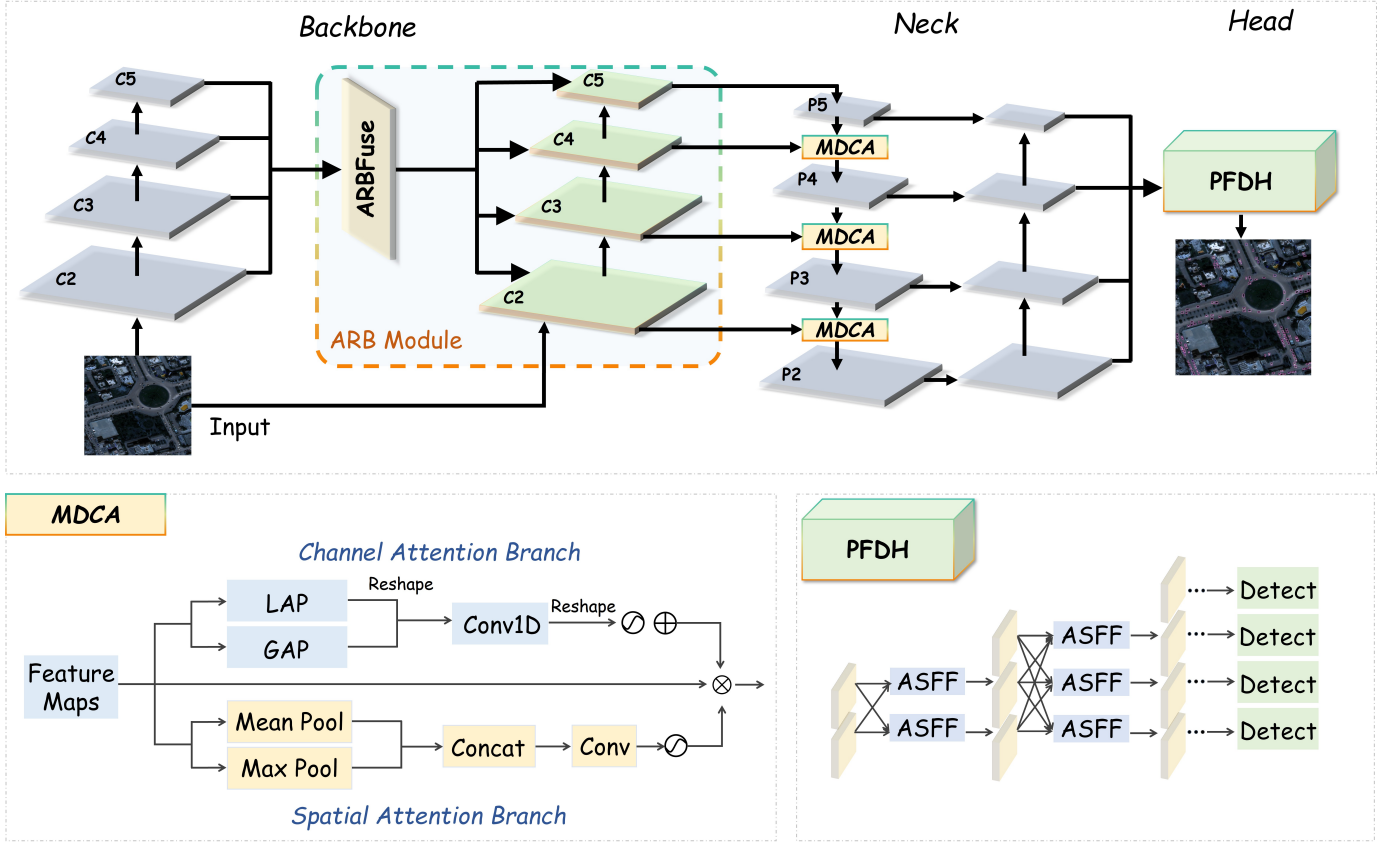


Fig. 2. The overall architecture of RS-TinyNet, which is built upon YOLOv11 and enhances RS tiny object detection through stage-wise multi-level feature fusion. The proposed framework integrates three key modules: ARB, MDCA, and PFDH to improve feature integrity and tiny object saliency.

1) *Channel Attention Branch*: In this branch, we fuse local and global contextual information, extract features using local average pooling and global average pooling for the input feature map $X \in \mathbb{R}^{B \times C \times H \times W}$, respectively, and then model inter-channel dependencies by lightweight 1D convolution. Finally, the local and global features are reshaped and weighted for fusion. This fusion of visual features from different scales facilitates the capture of subtle differences between tiny objects and complex backgrounds in RS images, leading to a deeper comprehension of the images. This process can be mathematically described as

$$z_{local} = \frac{1}{s^2} \sum_{i=1}^s \sum_{j=1}^s \text{LAP}_{s \times s}(X)[:, :, i, j], \quad (1)$$

$$z_{global} = \text{GAP}_{H \times W}(X), \quad (2)$$

$$a_{local} = \sigma(\text{Conv1D}(z_{local})), \quad (3)$$

$$a_{global} = \sigma(\text{Conv1D}(z_{global})), \quad (4)$$

$$A_{channel} = \text{Reshape}((1 - \lambda) \cdot a_{global} + \lambda \cdot a_{local}), \quad (5)$$

where z_{local} and z_{global} denote local and global contextual feature vectors, s is the local pooling window size, i and j are spatial position indices within the pooling window, H and W represent the height and width of the input feature map, σ is the Sigmoid activation function, a_{local} and a_{global} refer to the local and global channel attention, respectively, λ is the

local channel attention weight, and $A_{channel}$ denotes the fused channel attention.

2) *Spatial Attention Branch*: This branch focuses on the spatial distribution information of object, and fuses the average pooling and the maximum pooling on the channel dimension for spatial attention modeling. By extracting the statistical features of the channel dimension from the original feature map, it encodes the key regions on the spatial dimension, and then realizes the information reweighting on the spatial level. Unlike the traditional global modeling strategy, spatial attention branching can significantly highlight the location region where the tiny object is located and weaken the background interference while keeping the computation cost controllable, thus making up for the lack of perception of spatial structure by channel attention. Specifically, the construction of spatial attention branching can be written as

$$F_{avg} = \text{Mean}(X, \text{dim} = 1), F_{max} = \text{Max}(X, \text{dim} = 1), \quad (6)$$

$$F_{spatial} = \text{Concat}[F_{avg}, F_{max}], \quad (7)$$

$$A_{spatial} = \sigma(\text{Conv}_{3 \times 3}(F_{spatial})), \quad (8)$$

where F_{avg} and F_{max} denote the average and maximum pooling characteristics of the channel dimension, respectively, $F_{spatial}$ is the spliced spatial feature, and $A_{spatial}$ refers to the spatial attention.

These two branches work together by sharing input and modeling in parallel to jointly construct a channel-spatial bidirectional enhanced feature representation mechanism, which offers improved assistance in distinguishing tiny objects in RS images. The overall output of MDCA is consistent with the shape of the inputs, and it has the flexibility of plug-and-play, which is suitable for the enhancement of the features of the main trunk at all scales, and it greatly enhances the model's ability of localizing tiny objects in complex backgrounds. With the introduction of MDCA, the network can not only capture local structural changes, but also take into account the global layout information, effectively alleviating the common feature omission and ambiguous expression problems in multi-scale object detection.

C. Feature Integrity Reconstruction

Deep neural networks face the "information bottleneck" problem, where feature information decays during transmission and fusion. This can be formally expressed as

$$I(\mathcal{D}; \mathcal{D}) \geq I(\mathcal{D}; \mathcal{F}_\alpha(\mathcal{D})) \geq I(\mathcal{D}; \mathcal{K}_\beta(\mathcal{F}_\alpha(\mathcal{D}))), \quad (9)$$

where \mathcal{D} represents the original input, \mathcal{F}_α and \mathcal{K}_β are shallow and deep feature transformation functions, respectively, and $I(\cdot; \cdot)$ denotes mutual information. This decay leads to feature loss and gradient bias, particularly detrimental for tiny object detection, where features are already scarce. To mitigate this, we introduce two complementary modules. The ARB module focuses on preserving information and stabilizing gradients throughout the feature extraction process, while the PFDH module zooms in on alleviating information loss due to semantic gaps during multi-scale feature integration. Together, they enhance feature integrity and detection robustness for tiny objects in complex RS scenarios. The details of these two modules are shown below.

1) *ARB Module*: To alleviate the information bottleneck in the feature extraction network and guarantee the reliability of the gradient in each layer, we introduce the ARB module inspired by reversible residual networks. The reversible network avoids one-way loss of intermediate information by designing reversible residual blocks so that each layer's input features can be accurately reconstructed during backpropagation. It can be expressed by the following equation:

$$\text{Forward} : (y_1, y_2) = (x_1 + \mathcal{F}(x_2), x_2 + \mathcal{G}(x_1 + \mathcal{F}(x_2))), \quad (10)$$

$$\text{Inverse} : (x_1, x_2) = (y_1 - \mathcal{F}(y_2 - \mathcal{G}(y_1)), y_2 - \mathcal{G}(y_1)), \quad (11)$$

where \mathcal{F} and \mathcal{G} denote sub-functions of the invertible connected structure. However, the fully reversible structure will have degraded performance on shallow networks due to restricted mapping, while over-expanding the network size will bring huge computational overhead, which does not meet the real-time detection requirements.

The core goal of ARB is to alleviate the problem of information decay in feature transfer while avoiding an excessive increase in model complexity. As shown in the top part of Fig. 2, ARB employs an auxiliary supervision mechanism that generates the complementary gradient through classification and regression losses and fuses them with the main branch

through weighted integration. This design enables different network levels, particularly deeper layers, to access more comprehensive object information beyond specific scales or local regions. In addition, the auxiliary branch uses a small number of convolutional operations to achieve cross-layer feature reconstruction with fewer parameters and uses jump connections to retain shallow details, effectively compensating for the tiny object information that may be lost in the depth processing of the main network. Unlike the fully reversible structure, this design relaxes the requirement for the main branch to preserve complete original information, significantly reducing computational overhead while maintaining detection performance. Ultimately, with the stable gradient flow provided by the reversible connection and the efficient supervision of the auxiliary branch, the backbone network can maintain sensitivity to global features at all scales, thus effectively alleviating the information bottleneck problem and improving the accuracy and convergence speed of RS tiny object detection.

2) *PFDH Module*: While ARB tackles the preservation of information during feature extraction, the PFDH module addresses the challenge of effective multi-scale feature fusion. Conventional pyramid networks suffer from semantic inconsistencies when merging features across scales, particularly impairing tiny object detection through information conflicts and detail loss. As shown in the bottom right of Fig. 2, our PFDH module addresses these limitations by constructing feature fusion paths in a progressive manner to achieve level-by-level alignment of high-level semantics and low-level details, thereby improving detection accuracy and stability.

The core idea of PFDH is that fusion always occurs only between neighboring levels and avoids direct cross-layer connections, thus effectively mitigating the semantic gaps between different levels, and ensuring the continuous fusion of details and contextual information. It is worth noting that element-by-element summation is not an effective method during the gradual fusion of multi-level features because different objects at a certain position between levels may be contradictory. Accordingly, adaptive spatial feature fusion is applied to assign level-specific spatial weights, enhancing salient feature and alleviating contradictions among different object representations. We define the output of the l th layer of the fusion layer at the spatial location (u, v) as

$$Z_{uv}^l = \sum_{i=1}^N \omega_{uv}^{(i)} \cdot F_{uv}^{(i)}, \text{ with } \sum_{i=1}^N \omega_{uv}^{(i)} = 1, \quad (12)$$

where $F_{uv}^{(i)}$ denotes the input from the i th feature layer, and N denotes the number of fused source feature maps (usually 3). Through level-by-level fusion and adaptive spatial alignment, PFDH not only effectively preserves the semantic and spatial features of each level, but also mitigates the semantic drift and noise interference caused by cross-layer direct fusion, which significantly improves the overall modeling capability of multi-scale objects in RS images, and is especially suitable for the task of detecting far-small, overlapping, and dense tiny objects in RS images.

TABLE I
EVALUATION RESULTS OF DIFFERENT MODELS ON AI-TOD DATASET

Method	Backbone	Year	AP	AP ₅₀	AP ₇₅	AP _{vt}	AP _t	AP _s	AP _m
Faster R-CNN [5]	ResNet50	2015	11.4	27.0	8.0	0.0	8.3	23.1	24.5
Cascade R-CNN [7]	ResNet50	2018	13.8	30.8	10.5	0.0	10.6	25.5	26.6
RepPoints [36]	ResNet50	2019	9.2	23.6	5.3	2.5	9.2	12.9	14.4
FCOS [37]	ResNet50	2019	13.9	35.5	8.6	2.7	12.0	20.2	32.2
CenterNet [38]	DLA-34	2019	13.4	39.2	5.0	3.8	12.1	17.7	18.9
Dynamic R-CNN [39]	ResNet50	2020	16.9	39.4	11.6	6.3	18.2	21.2	25.6
TOOD [40]	ResNet50	2021	14.9	34.7	8.5	3.3	12.8	21.7	33.1
ATSS [41]	ResNet50	2021	12.8	30.6	8.5	1.9	11.6	19.5	29.2
DetectoRS [42]	ResNet50	2021	14.8	32.8	11.4	0.0	10.8	28.3	38.0
M-CenterNet [33]	DLA-34	2022	14.5	40.7	6.4	6.1	15.0	19.4	20.4
DetectoRS w/RFLA [43]	ResNet50	2022	24.8	55.2	18.5	9.3	24.8	30.3	38.2
FSANet [44]	Swin-T	2022	22.6	52.8	15.6	7.4	21.6	29.1	38.5
YOLOv8m [24]	-	2023	22.9	52.5	16.5	7.1	22.4	29.7	37.1
HANet [45]	ResNet50	2023	22.1	53.7	14.4	10.9	22.2	27.3	36.8
SSRDet [46]	ResNet50	2023	24.7	56.0	18.4	9.4	25.6	30.2	38.1
MENet [47]	Swin-T	2024	23.2	56.2	15.0	9.7	23.9	25.3	34.4
YOLOv11m [25]	-	2024	27.9	57.9	23.1	9.6	26.5	<u>37.8</u>	<u>44.8</u>
BRSTD [48]	-	2024	26.1	58.0	19.5	11.6	26.0	32.3	38.5
FFCA-YOLO [49]	CSPDarkNet53	2024	-	61.7	-	12.6	24.9	31.8	-
CAF ² ENet-M [2]	DarkNet-53	2024	<u>30.2</u>	<u>63.7</u>	<u>25.0</u>	<u>12.8</u>	<u>30.3</u>	36.7	41.9
DNTR [50]	ResNet50	2024	26.2	56.7	20.2	<u>12.8</u>	26.4	31.0	37.0
Zhang Y et al. [51]	-	2025	28.4	59.1	21.3	11.0	26.2	32.1	40.6
RS-TinyNet (Ours)	-	2025	34.2	65.2	31.5	14.9	34.5	41.5	46.6
<i>Compared with SOTA</i>	-	-	<i>+4.0</i>	<i>+1.5</i>	<i>+6.5</i>	<i>+2.1</i>	<i>+4.2</i>	<i>+3.7</i>	<i>+1.8</i>

Note: **Bold** values indicate the best performance achieved by our RS-TinyNet. Underlined values represent the suboptimal results. *Italic* values show the improvement margins of RS-TinyNet compared to the suboptimal results.

The incorporation of ARB and PFDH modules establishes a comprehensive framework for feature integrity preservation throughout the detection pipeline. The ARB module guarantees stable feature propagation and gradient flow through its assisted reversible architecture, whereas the PFDH module achieves precise multi-scale semantic alignment via progressive feature fusion. This synergistic combination effectively mitigates information degradation during both feature extraction and fusion phases, significantly improving detection accuracy for challenging RS tiny object scenarios.

IV. EXPERIMENT

A. Implementation Details

All experiments were based on the Pytorch framework and performed on NVIDIA A800 GPUs. The model was trained for 600 epochs with a batch size of 16. The image size for both the training and testing phases is 800×800 . The optimizer used was Stochastic Gradient Descent (SGD) [52] with a momentum of 0.937 and a weight decay of 0.0005. All of our experiments were performed in the above configuration, unless otherwise stated.

B. Datasets

1) *AI-TOD*: AI-TOD [33] is a RS dataset designed for detecting tiny objects, derived from a large publicly accessible collection of aerial images. It includes 280,036 aerial images, each measuring 800×800 pixels, and contains a total of 700,621 labeled object instances across eight common categories. The average size of the objects is just 12.8 pixels, with 85.6% of the objects being smaller than 16 pixels, making it ideal for testing detectors in difficult scenarios involving tiny objects. The dataset provides 11,214 images for training and 2,804 for validation, amounting to 14,018 images in total. An additional 14,018 images are provided in the test set to evaluate the performance of the detectors.

2) *DIOR*: DIOR [53] serves as a widely-used and challenging RS dataset focused on detecting small objects in complex scenes. It contains 23,463 high-resolution aerial images, annotated with a total of 190,288 object instances spanning twenty categories. Each image in the DIOR dataset measures 800×800 pixels, with spatial resolution ranging from 0.5 to 30 meters. It includes airplane (AL), airport (AT), baseball field (BF), basketball court (BC), bridge (BG), chimney (CM), dam (DM), expressway service area (EA), expressway toll station (ES), golf course (GC), ground track

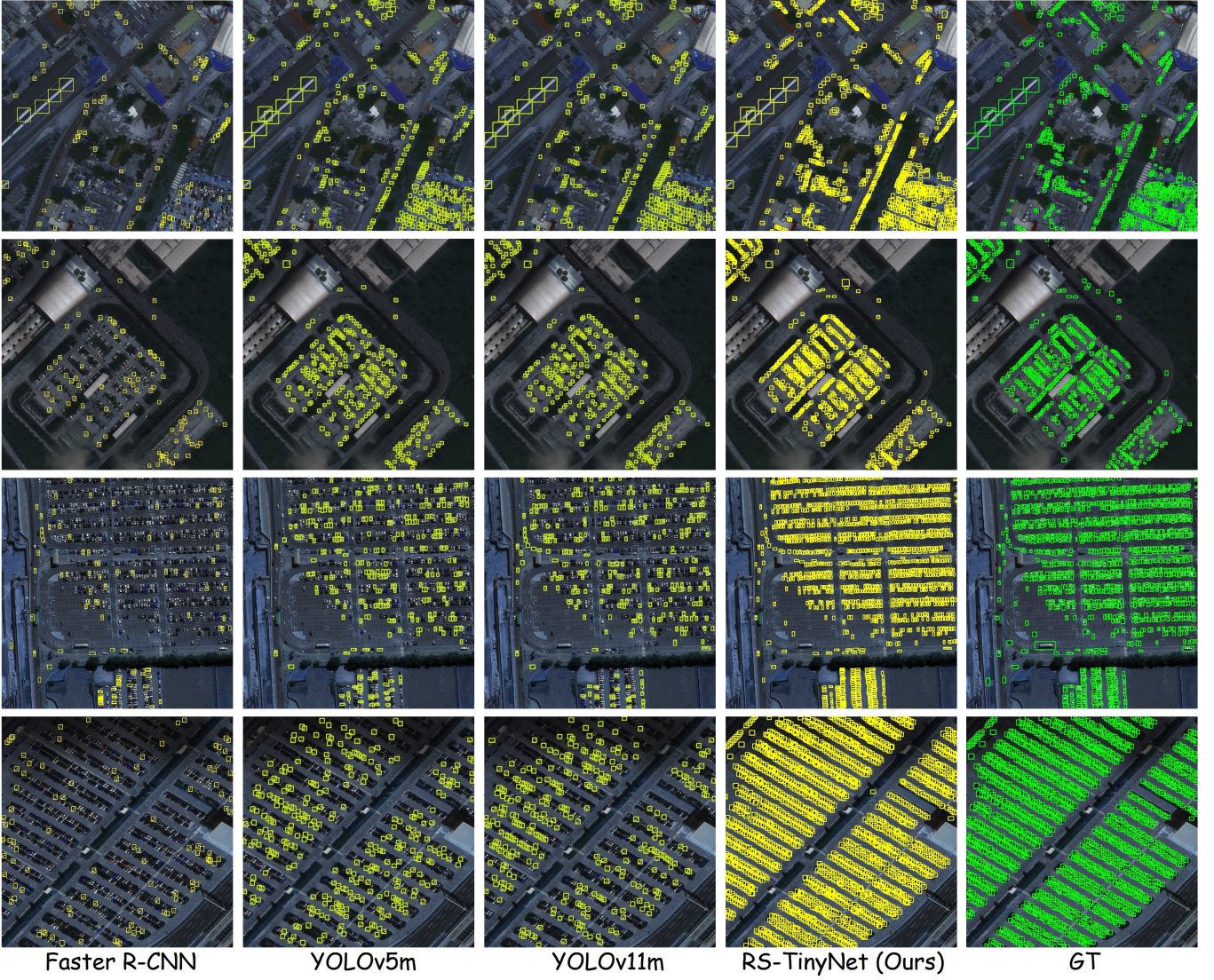


Fig. 3. Visualization of detection results on the AI-TOD dataset. While existing methods such as Faster R-CNN, YOLOv5m, and YOLOv11m tend to produce false positives and miss detections in dense tiny object scenarios, the proposed RS-TinyNet yields more precise and consistent detection outcomes.

field (GF), harbor (HB), overpass (OP), ship (SP), stadium (SD), storage tank (ST), tennis court (TC), train station (TS), vehicle (WH), windmill (WM). The training set includes 5,862 images and the validation set contains 5,863 images, totaling 11,725 images for model training. The test set consists of 11,738 images for evaluating the performance of the detector.

C. Evaluation Metrics

In our experiment, we use the average precision (AP) to evaluate the detection performance of the model. The formula is as follows

$$AP = \int_0^1 P(R) dR, \quad (13)$$

where P is the Precision and R is the Recall rate, defined as

$$P = \frac{TP}{TP + FN}, \quad (14)$$

$$R = \frac{TP}{TP + FP}, \quad (15)$$

where TP, FN, and FP mean the number of true positive, false negative, and false positive samples, respectively. To evaluate detection performance on AI-TOD, we employ the MS COCO [29] AP evaluation metrics across all detectors. In detail, AP_{50} denotes the accuracy at IoU threshold of 0.5, while AP_{75} adopts a stricter IoU threshold of 0.75. The overall AP is determined by averaging AP values at IoU thresholds from 0.5 to 0.95 in steps of 0.05. In addition, we use AP_{vt} , AP_t , AP_s , and AP_m to evaluate the robustness of model to objects of different scales. Specifically, AP_{vt} represents the AP for objects with a size less than 8^2 , AP_t for objects with sizes between $(8^2, 16^2]$, AP_s for objects with sizes between $(16^2, 32^2]$ and AP_m for objects with sizes larger than 32^2 .

To comprehensively evaluate the computational efficiency of object detectors, params is introduced as a critical indicator. Params quantifies the total learnable weights and biases in the model architecture, directly reflecting its memory footprint and storage requirements. Lower parameter counts generally imply reduced hardware resource demands.

TABLE II
EVALUATION RESULTS OF DIFFERENT MODELS ON DIOR DATASET

Method	AL	AT	BF	BC	BG	CM	DM	EA	ES	GC	GF	HB	OP	SP	SD	ST	TC	TS	VH	WM	mAP ₅₀
Dynamic R-CNN	56.9	65.1	62.4	85.5	33.5	76.0	45.5	54.4	47.6	72.5	69.7	36.8	53.5	70.4	45.4	57.4	78.3	47.7	37.0	73.7	58.5
TridentNet [54]	56.3	73.6	63.5	84.9	29.9	75.7	53.7	55.7	48.0	75.3	71.0	52.0	50.7	64.8	62.6	52.6	78.6	54.6	34.2	72.6	60.5
TOOD	67.6	67.9	67.1	85.8	33.1	77.9	49.0	55.7	48.7	73.2	67.6	47.9	52.8	73.7	40.4	65.7	79.9	49.7	45.9	74.2	61.2
CenterNet	67.8	68.8	75.0	83.2	44.6	67.6	44.7	71.0	67.7	66.6	72.5	41.6	56.3	72.0	54.2	70.3	85.3	48.6	51.6	79.2	64.4
Faster R-CNN	54.0	74.5	63.3	80.7	44.8	72.5	60.0	75.6	62.3	76.0	76.8	46.4	57.2	71.8	68.3	53.8	81.1	59.5	43.1	81.2	65.2
FCOS	54.8	77.2	73.6	87.3	36.8	79.2	56.2	75.7	60.7	80.2	75.6	50.2	55.2	69.7	62.5	51.1	85.7	51.2	41.2	83.5	65.4
PANet [55]	60.2	72.0	70.6	80.5	43.6	72.3	61.4	72.1	66.7	72.0	73.4	45.3	56.9	71.7	70.4	62.0	80.9	57.0	47.2	84.5	66.1
YOLOv5m [23]	87.3	61.7	73.8	89.0	42.6	77.5	55.2	63.8	63.2	66.9	78.0	58.2	58.1	87.8	54.3	79.3	89.7	50.2	54.0	79.6	68.6
DETR [56]	39.6	74.0	65.2	80.7	26.5	75.2	66.8	70.5	52.5	74.2	62.1	27.4	47.0	8.5	46.0	14.5	64.4	54.3	14.4	55.6	68.6
RAST-YOLO [57]	84.3	76.4	78.7	85.9	40.2	76.8	50.2	62.6	56.5	77.1	73.7	61.1	56.6	91.1	74.3	77.9	89.3	53.3	54.0	76.2	69.8
ASSD [58]	85.6	82.4	75.8	89.5	40.7	77.6	64.7	67.1	61.7	80.8	78.6	62.0	58.0	84.9	76.7	65.3	87.9	62.4	44.5	76.3	71.1
AGMF-Net [32]	90.9	72.8	79.3	89.7	44.7	81.4	59.3	66.0	62.7	73.8	79.2	65.0	61.7	91.7	78.6	75.8	90.7	60.0	58.0	83.1	73.2
KLDet [59]	67.5	81.1	75.5	88.2	47.4	78.6	65.0	82.9	72.7	79.9	80.8	57.1	61.2	91.2	74.6	75.2	87.4	56.9	56.0	90.2	73.5
YOLOv11m	82.5	83.2	80.3	89.3	45.0	78.5	69.9	69.3	65.4	77.6	75.6	65.6	59.9	90.8	70.5	78.2	89.5	66.1	54.7	80.3	73.6
RS-TinyNet(Ours)	81.0	83.5	77.5	89.6	48.6	82.0	64.6	71.2	72.9	80.9	74.6	65.9	63.9	92.1	63.7	78.4	90.8	62.0	58.7	84.1	74.3

D. Comparative Experiments

We compared RS-TinyNet with other excellent object detection models on the AI-TOD and DIOR datasets, including the baseline model YOLOv11m.

1) *Results on AI-TOD*: Tab. I shows the experimental results of different SOTA detectors on the AI-TOD dataset. Obviously, compared with other methods, our proposed method achieves the highest accuracy of RS tiny object detection. Firstly, under the same experimental conditions, the detection performance of our method was improved in all aspects compared to the baseline model. Specifically, the AP, AP₅₀, AP₇₅, AP_{vt}, AP_t, AP_s, and AP_m metrics are improved by 6.3%, 7.3%, 8.4%, 5.3%, 8.0%, 3.7% and 1.8%, respectively. Secondly, compared with the traditional object detection network, the detection performance of our method is greatly improved, especially for tiny objects. Additionally, compared with other tiny object detection algorithms that achieved SOTA performance on AI-TOD, our method demonstrates superior performance. Our AP₅₀ metric surpasses that of MENet, DNTR, BRSTD, FFCA-YOLO, and CAF²ENet-M by 9.0%, 8.5%, 7.2%, 3.5%, and 1.5%, respectively, accompanied by improvements in other evaluation metrics. The experimental results unequivocally indicate that our method outperforms existing SOTA approaches on the AI-TOD. Fig. 3 presents the detection outcomes obtained on the test set.

2) *Results on DIOR*: To evaluate the generalization of our model, we performed comparative experiments using the DIOR dataset, and the findings are presented in Tab. II. The comparative experiments on the DIOR dataset show that our

algorithm performs well with a good generalization ability of 74.3% mAP₅₀, which is better than other SOTA models. The proposed model achieves SOTA performance in 11 of the 20 categories, with accuracy scores of 83.5% (AT), 89.6% (BC), 48.6% (BG), 82.0% (CM), 72.9% (ES), 80.9% (GC), 65.9% (HB), 63.9% (OP), 92.1% (SP), 90.8% (TC), and 58.7% (VH). The outcomes of our model's detection on the DIOR dataset are illustrated in Fig. 4, where different color-labeled detection boxes are used to denote different classes of objects. This serves as strong evidence for the efficacy of the proposed RS-TinyNet.

E. Ablation Experiments

To verify the effectiveness of each component, we conducted a series of experiments on the test set of AI-TOD. Tab. III shows the performance and model size of our detectors under different conditions. Using YOLOv11m as the baseline model, the accuracy of the model was evaluated by adding the DCA, ARB, and PFDH modules. To ensure fair comparisons, all hyperparameters are identical in the following experiments.

1) *Effectiveness of MDCA*: To assess the impact of the suggested MDCA module, we incorporated it into the neck architecture of the baseline YOLOv11m model. As shown in Tab. III, the incorporation of the MDCA module alone yields significant performance enhancements, improving AP, AP₅₀, and AP₇₅ by 3.2%, 3.7%, and 4.2%, respectively. This improvement benefits from MDCA's ability to efficiently capture and fuse channel, spatial, local and global information, enabling the network to better distinguish objects in complex

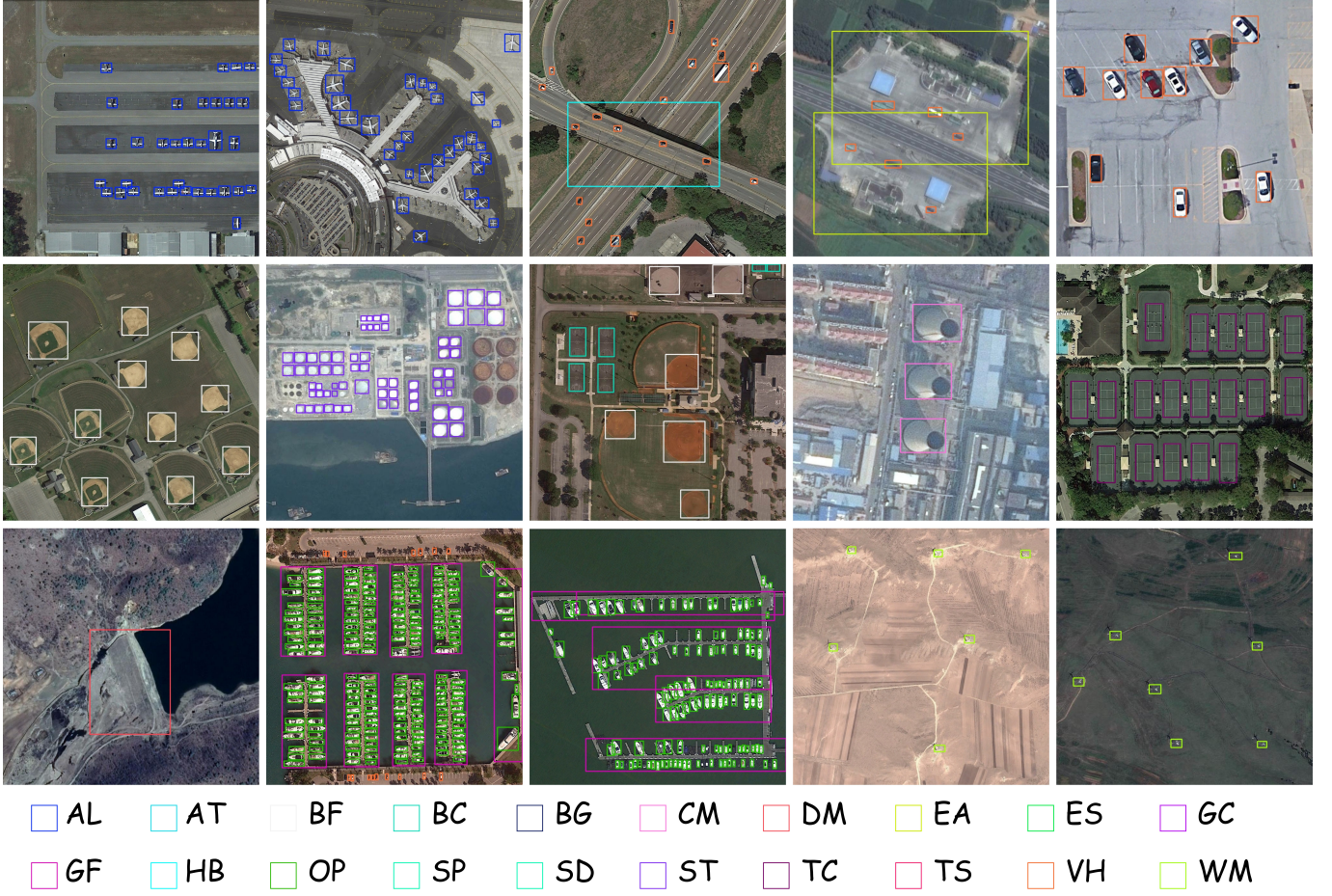


Fig. 4. Some representative detection results of RS-TinyNet on the DIOR dataset. Detection boxes for different object categories are annotated using different colors.

TABLE III
EFFECTIVENESS OF DIFFERENT COMPONENTS ON AI-TOD

Baseline	MDCA	ARB	PFDH	AP	AP ₅₀	AP ₇₅	Params
YOLOv11m				27.9	57.9	23.1	20059176
	✓			31.1	61.6	27.3	20553992
		✓		30.6	61.0	26.7	30294440
			✓	28.5	59.6	25.2	21339925
	✓	✓		32.7	63.6	28.5	30789252
		✓	✓	33.4	64.3	30.5	32159637
	✓		✓	32.8	63.7	28.6	22059206
	✓	✓	✓	34.2	65.2	31.5	32159657

RS scenes. The results confirm that MDCA can significantly improve detection sensitivity, particularly for dense and scale-variant tiny objects.

2) *Effectiveness of ARB*: We further examined the impact of the ARB module by inserting it separately into the baseline network. The results show that the introduction of ARB improves AP, AP₅₀, and AP₇₅ by 2.7%, 3.1%, and 3.6%, respectively, as shown in Tab. III. The ARB module introduces auxiliary branches with reversible transformations, which fa-

cilitates feature learning through gradient-guided feedback and effectively mitigates the loss of information of the object feature during the feedforward. This contributes to more stable representation learning, especially for tiny and fuzzy objects. Meanwhile, it can work together with the MDCA module to achieve further improvement in model detection performance. The above results show that ARB effectively enhances the capability of feature integrity reconstruction and improves the detection performance.

3) *Effectiveness of PFDH*: To evaluate the effectiveness of the proposed PFDH, we first added it independently to the baseline model and observed a significant improvement in AP, AP₅₀, and AP₇₅, as shown in Tab. III. The PFDH module enhances detection capabilities by gradually integrating cross-scale semantic features, thereby improving localization and classification, especially for tiny and multi-scale objects. Furthermore, when combined with the MDCA or ARB modules, detection performance is further improved, indicating strong compatibility and complementary advantages between these components. Notably, the combination of MDCA and PFDH better refines the prominent features of multi-scale objects, while the integration of ARB and PFDH enhances feature integrity reconstruction, preventing feature loss, especially for tiny objects. When all three modules, MDCA, ARB, and

PFDH, are applied jointly, our model achieves the best performance, confirming the cumulative and mutually reinforcing effects of each module. These results validate the overall effectiveness and rationality of our architecture design.

V. CONCLUSION

In this paper, we propose RS-TinyNet, an innovative and effective framework for tiny object detection in RS images. By integrating MDCA module for tiny object saliency modeling, ARB and PFDH modules for feature integrity reconstruction, RS-TinyNet substantially enhances the representation and discrimination of tiny objects. Experiments on benchmark datasets such as AI-TOD and DIOR demonstrate that RS-TinyNet significantly outperforms existing SOTA methods in detection precision, with particularly significant improvements for tiny objects. Further ablation studies confirm the importance of each suggested module, highlighting their synergistic effects in addressing challenges such as low pixel density, object density, and lack of clear structural details in RS scenarios. These findings underscore the exceptional robustness and practical utility of RS-TinyNet in real-world applications. Future research directions will focus on computational efficiency enhancements and framework optimization, as well as investigating its adaptability and generalization capability across a wider spectrum of RS tasks.

REFERENCES

- [1] Luqi Gong, Haotian Chen, Yikun Chen, Tianliang Yao, Chao Li, Shuai Zhao, and Guangjie Han. Dpnet: Dynamic pooling network for tiny object detection. *arXiv preprint arXiv:2505.02797*, 2025.
- [2] Lihui Ge, Guanqun Wang, Tong Zhang, Yin Zhuang, He Chen, Hao Dong, and Liang Chen. Regression-guided refocusing learning with feature alignment for remote sensing tiny object detection. *IEEE Transactions on Geoscience and Remote Sensing*, 2024.
- [3] Xuehui Yu, Zhenjun Han, Yuqi Gong, Nan Jan, Jian Zhao, Qixiang Ye, Jie Chen, Yuan Feng, Bin Zhang, Xiaodi Wang, et al. The 1st tiny object detection challenge: Methods and results. In *Computer Vision—ECCV 2020 Workshops: Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pages 315–323. Springer, 2020.
- [4] Mahya Nikouei, Bitar Baroutian, Shahabedin Nabavi, Fateme Taraghi, Atefeh Aghaei, Ayoob Sajedi, and Mohsen Ebrahimi Moghaddam. Small object detection: A comprehensive survey on challenges, techniques and real-world applications. *arXiv preprint arXiv:2503.20516*, 2025.
- [5] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015.
- [6] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [7] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162, 2018.
- [8] Wei Zhang, Miaoxin Cai, Tong Zhang, Guoqiang Lei, Yin Zhuang, and Xuerui Mao. Popeye: A unified visual-language model for multi-source ship detection from remote sensing imagery. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2024.
- [9] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*, pages 21–37. Springer, 2016.
- [10] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [11] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7263–7271, 2017.
- [12] Joseph Redmon and Ali Farhadi. Yolo3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- [13] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Region-based convolutional networks for accurate object detection and segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 38(1):142–158, 2015.
- [14] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.
- [15] Jiaming Han, Jian Ding, Jie Li, and Gui-Song Xia. Align deep features for oriented object detection. *IEEE transactions on geoscience and remote sensing*, 60:1–11, 2021.
- [16] Jakaria Rabbi, Nilanjan Ray, Matthias Schubert, Subir Chowdhury, and Dennis Chao. Small-object detection in remote sensing images with end-to-end edge-enhanced gan and object detector network. *Remote Sensing*, 12(9):1432, 2020.
- [17] Jingqian Wu and Shibiao Xu. From point to region: Accurate and efficient hierarchical small object detection in low-resolution remote sensing images. *Remote Sensing*, 13(13):2620, 2021.
- [18] Jianxiang Li, Zili Zhang, Yan Tian, Yiping Xu, Yihong Wen, and Shicheng Wang. Target-guided feature super-resolution for vehicle detection in remote sensing images. *IEEE geoscience and remote sensing letters*, 19:1–5, 2021.
- [19] Jianqi Chen, Keyan Chen, Hao Chen, Zhengxia Zou, and Zhenwei Shi. A degraded reconstruction enhancement-based method for tiny ship detection in remote sensing images with a new large-scale dataset. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–14, 2022.
- [20] Chang Xu, Jinwang Wang, Wen Yang, Huai Yu, Lei Yu, and Gui-Song Xia. Detecting tiny objects in aerial images: A normalized wasserstein distance and a new benchmark. *ISPRS Journal of Photogrammetry and Remote Sensing*, 190:79–93, 2022.
- [21] Chien-Yao Wang, I-Hau Yeh, and Hong-Yuan Mark Liao. Yolo9: Learning what you want to learn using programmable gradient information. In *European conference on computer vision*, pages 1–21. Springer, 2024.
- [22] Guoyu Yang, Jie Lei, Zhikuan Zhu, Siyu Cheng, Zunlei Feng, and Ronghua Liang. Afpn: Asymptotic feature pyramid network for object detection. In *2023 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pages 2184–2189. IEEE, 2023.
- [23] G. Jocher. Yolo5 by ultralytics. <https://github.com/ultralytics/yolov5>, 2020.
- [24] G. Jocher. Ultralytics yolo(version 8.0.0). <https://github.com/ultralytics/ultralytics>, 2023.
- [25] G. Jocher. Ultralytics yolo(version 8.3.0). <https://github.com/ultralytics/ultralytics>, 2024.
- [26] Yunjie Tian, Qixiang Ye, and David Doermann. Yolo12: Attention-centric real-time object detectors. *arXiv preprint arXiv:2502.12524*, 2025.
- [27] Cheng-Yang Fu, Wei Liu, Ananth Ranga, Ambrish Tyagi, and Alexander C Berg. Dssd: Deconvolutional single shot detector. *arXiv preprint arXiv:1701.06659*, 2017.
- [28] Shifeng Zhang, Longyin Wen, Xiao Bian, Zhen Lei, and Stan Z Li. Single-shot refinement neural network for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4203–4212, 2018.
- [29] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision—ECCV 2014: 13th European conference, Zurich, Switzerland, September 6–12, 2014, proceedings, part v 13*, pages 740–755. Springer, 2014.
- [30] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88:303–338, 2010.
- [31] Vishnu Chalavadi, Prudviraj Jeripothula, Rajeshreddy Datla, Sobhan Babu Ch, et al. msodanet: A network for multi-scale object detection in aerial images using hierarchical dilated convolutions. *Pattern Recognition*, 126:108548, 2022.
- [32] Tao Gao, Ziqi Li, Yuanbo Wen, Ting Chen, Qianqian Niu, and Zixiang Liu. Attention-free global multiscale fusion network for remote sensing object detection. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–14, 2023.
- [33] Jinwang Wang, 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)*, pages 3791–3798. IEEE, 2021.
- [34] Lianyu Cao, Xiaolu Zhang, Zhaoshun Wang, and Guangyu Ding. Multi angle rotation object detection for remote sensing image based on

- modified feature pyramid networks. *International Journal of Remote Sensing*, 42(14):5253–5276, 2021.
- [35] Bao Liu and Jinlei Huang. Global-local attention mechanism based small object detection. In *2023 IEEE 12th Data Driven Control and Learning Systems Conference (DDCLS)*, pages 1439–1443. IEEE, 2023.
- [36] Ze Yang, Shaohui Liu, Han Hu, Liwei Wang, and Stephen Lin. Repoints: Point set representation for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9657–9666, 2019.
- [37] Zhi Tian, 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*, pages 9627–9636, 2019.
- [38] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019.
- [39] Hongkai Zhang, Hong Chang, Bingpeng Ma, Naiyan Wang, and Xilin Chen. Dynamic r-cnn: Towards high quality object detection via dynamic training. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XV 16*, pages 260–275. Springer, 2020.
- [40] Chengjian Feng, Yujie Zhong, Yu Gao, Matthew R Scott, and Weilin Huang. Tood: Task-aligned one-stage object detection. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3490–3499. IEEE Computer Society, 2021.
- [41] Shifeng Zhang, Cheng Chi, Yongqiang Yao, Zhen Lei, and Stan Z Li. Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9759–9768, 2020.
- [42] Siyuan Qiao, Liang-Chieh Chen, and Alan Yuille. Detectors: Detecting objects with recursive feature pyramid and switchable atrous convolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10213–10224, 2021.
- [43] Chang Xu, Jinwang Wang, Wen Yang, Huai Yu, Lei Yu, and Gui-Song Xia. Rfla: Gaussian receptive field based label assignment for tiny object detection. In *European conference on computer vision*, pages 526–543. Springer, 2022.
- [44] Jixiang Wu, Zongxu Pan, Bin Lei, and Yuxin Hu. Fsanet: Feature-and-spatial-aligned network for tiny object detection in remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–17, 2022.
- [45] Chengxi Han, Chen Wu, Haonan Guo, Meiqi Hu, and Hongruixuan Chen. Hanet: A hierarchical attention network for change detection with bitemporal very-high-resolution remote sensing images. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 16:3867–3878, 2023.
- [46] Lijuan Zhang, Minhui Wang, Yutong Jiang, Dongming Li, and Yue Zhou. Ssrnet: Small object detection based on feature pyramid network. *IEEE Access*, 11:96743–96752, 2023.
- [47] Tianyang Zhang, Xiangrong Zhang, Xiaoqian Zhu, Guanchun Wang, Xiao Han, Xu Tang, and Licheng Jiao. Multistage enhancement network for tiny object detection in remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–12, 2024.
- [48] Sihan Huang, Chuan Lin, Xintong Jiang, and Zhenshen Qu. Brstd: Bio-inspired remote sensing tiny object detection. *IEEE Transactions on Geoscience and Remote Sensing*, 2024.
- [49] Yin Zhang, Mu Ye, Guiyi Zhu, Yong Liu, Pengyu Guo, and Junhua Yan. Ffca-yolo for small object detection in remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–15, 2024.
- [50] Hou-I Liu, Yu-Wen Tseng, Kai-Cheng Chang, Pin-Jyun Wang, Hong-Han Shuai, and Wen-Huang Cheng. A denoising fpn with transformer r-cnn for tiny object detection. *IEEE Transactions on Geoscience and Remote Sensing*, 2024.
- [51] Yunzuo Zhang, Ting Liu, Jiawen Zhen, Yaoxing Kang, and Yu Cheng. Adaptive downsampling and scale enhanced detection head for tiny object detection in remote sensing image. *IEEE Geoscience and Remote Sensing Letters*, 22:1–5, 2025.
- [52] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [53] Ke Li, Gang Wan, Gong Cheng, Liqiu Meng, and Junwei Han. Object detection in optical remote sensing images: A survey and a new benchmark. *ISPRS journal of photogrammetry and remote sensing*, 159:296–307, 2020.
- [54] Yanghao Li, 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*, pages 6054–6063, 2019.
- [55] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8759–8768, 2018.
- [56] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020.
- [57] Xuzhao Jiang and Yonghong Wu. Remote sensing object detection based on convolution and swin transformer. *IEEE Access*, 11:38643–38656, 2023.
- [58] Tao Xu, Xian Sun, Wenhui Diao, Liangjin Zhao, Kun Fu, and Hongqi Wang. Assd: Feature aligned single-shot detection for multiscale objects in aerial imagery. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–17, 2021.
- [59] Zhuangzhuang Zhou and Yingying Zhu. Kldet: Detecting tiny objects in remote sensing images via kullback-leibler divergence. *IEEE Transactions on Geoscience and Remote Sensing*, 2024.