

# Learning to Aggregate Multi-Scale Context for Instance Segmentation in Remote Sensing Images

Ye Liu<sup>1</sup>, Huifang Li<sup>1</sup>, *Member, IEEE*, Chao Hu<sup>1</sup>, Shuang Luo<sup>1</sup>, Yan Luo<sup>1</sup>,  
and Chang Wen Chen<sup>2</sup>, *Fellow, IEEE*

**Abstract**—The task of instance segmentation in remote sensing images, aiming at performing per-pixel labeling of objects at instance level, is of great importance for various civil applications. Despite previous successes, most existing instance segmentation methods designed for natural images encounter sharp performance degradations when they are directly applied to top-view remote sensing images. Through careful analysis, we observe that the challenges mainly come from the lack of discriminative object features due to severe scale variations, low contrasts, and clustered distributions. In order to address these problems, a novel context aggregation network (CATNet) is proposed to improve the feature extraction process. The proposed model exploits three lightweight plug-and-play modules, namely dense feature pyramid network (DenseFPN), spatial context pyramid (SCP), and hierarchical region of interest extractor (HROIE), to aggregate global visual context at feature, spatial, and instance domains, respectively. DenseFPN is a multi-scale feature propagation module that establishes more flexible information flows by adopting inter-level residual connections, cross-level dense connections, and feature re-weighting strategy. Leveraging the attention mechanism, SCP further augments the features by aggregating global spatial context into local regions. For each instance, HROIE adaptively generates RoI features for different downstream tasks. Extensive evaluations of the proposed scheme on iSAID, DIOR, NWPU VHR-10, and HRSID datasets demonstrate that the proposed approach outperforms state-of-the-arts under similar computational costs. Source code and pre-trained models are available at <https://github.com/yeliudev/CATNet>.

**Index Terms**—Instance Segmentation, Object Detection, Feature Pyramid Networks, Global Context Aggregation

## I. INTRODUCTION

RECENT advances in satellites and remote sensing techniques have generated a large variety of high-resolution remote sensing images, bringing great challenges to manual

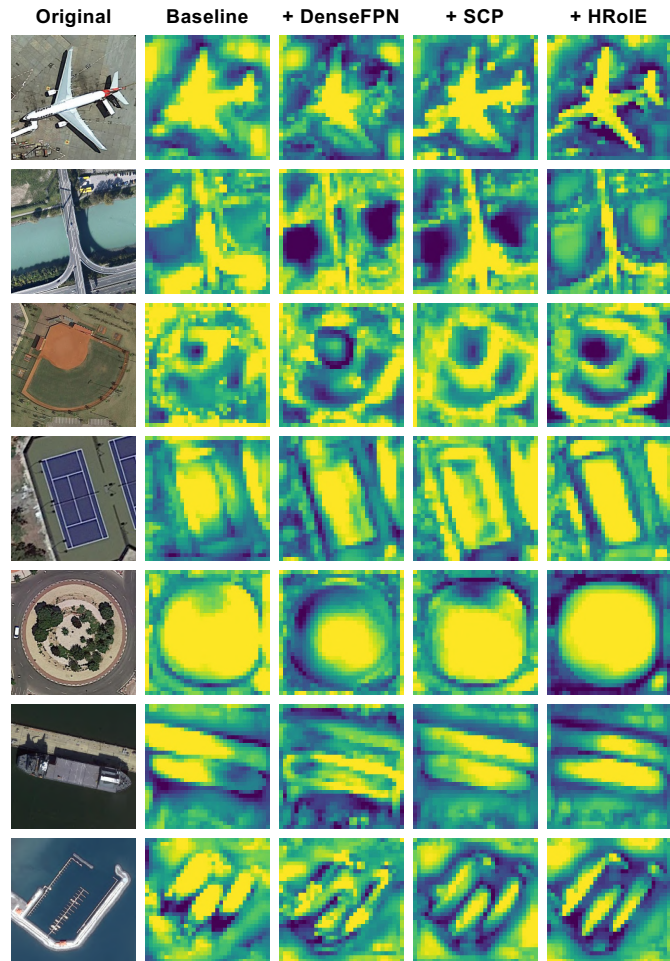


Fig. 1. Illustration of global context aggregation. Processed by the proposed modules, object features are enhanced by aggregating global visual context, as can be seen from the more discriminative feature maps.

manipulation and processing. Therefore, automatic analysis and understanding of these images are becoming increasingly essential for various civil applications [1]–[6]. As a fundamental yet challenging task in computer vision, instance segmentation, which is a combination of object detection and semantic segmentation aiming at predicting binary masks of objects at instance level, has been widely used to extract fine-grained object information from both optical remote sensing images and synthetic aperture radar (SAR) images. It has attracted considerable attention in recent years.

Most existing works on object detection and instance seg-

Manuscript received March 22, 2022; revised June 19, 2023 and August 24, 2023; accepted November 20, 2023. This work was supported by the National Natural Science Foundation of China under Grant 42371366, the Hubei Provincial Key Research and Development Program under Grant 2023BAB066, and the Fundamental Research Funds for the Central Universities under Grant 2042023kfy04. (*Corresponding author: Huifang Li*)

Ye Liu is with the Department of Computing, The Hong Kong Polytechnic University, Hong Kong SAR, China (e-mail: coco.ye.liu@connect.polyu.hk).

Huifang Li and Chao Hu are with the School of Resource and Environmental Sciences, Wuhan University, Wuhan 430079, China (e-mail: huifangli@whu.edu.cn; chaohu@whu.edu.cn).

Shuang Luo is with Changjiang Spatial Information Technology Engineering Co., Ltd., Wuhan 430074, China (e-mail: sluo@whu.edu.cn).

Yan Luo is with the Department of Computing, The Hong Kong Polytechnic University, Hong Kong SAR, China (e-mail: silver.luo@connect.polyu.hk).

Chang Wen Chen is with the Department of Computing, The Hong Kong Polytechnic University, Hong Kong SAR, China, and also with Peng Cheng Laboratory, Shenzhen 518055, China (e-mail: changwen.chen@polyu.edu.hk).

Digital Object Identifier: 10.1109/TNNLS.2023.3336563



Fig. 2. The challenges of performing instance segmentation in remote sensing images. Here, scale variation, arbitrary orientation, and clustered distribution lead to complicated object patterns while low contrast and cluttered background bring interfering information from the background.

mentation [4], [7]–[12] have been successful in conventional front-view scenes. However, when directly applied to remote sensing images, a large number of these methods inevitably encounter performance degradations [13]–[17]. Compared with natural images, remote sensing images are typically viewed from the top, capturing large areas with limited object discrepancies. We analyze the peculiarities of these scenes and divide the challenges into five distinct aspects, *i.e.*, scale variation, arbitrary orientation, clustered distribution, low contrast, and cluttered background, as illustrated in Fig. 2. The first three aspects lead to complicated object patterns while the remaining two bring interfering information from the background. These phenomena are scarce in natural scenes, thus only a few works have considered these aspects. We argue that all the challenges above are due to the lack of discriminative object features in remote sensing images. That is, visual appearances of individual objects in remote sensing images are not informative enough for directly adopting existing schemes to perform instance segmentation. Such a view is also supported by previous works [18]–[22].

A natural question will be: *How to enhance the inadequate features to achieve better instance segmentation results in remote sensing images?* Considering that, in a general instance segmentation pipeline [9]–[11], object representations are directly cropped from either of the feature maps from the backbone or neck, containing only local features with irreversible information loss. In this work, we mitigate this problem by introducing CATNet, a novel framework for global context aggregation. The key idea is that context information of images, coming from different feature pyramid levels, spatial positions, or receptive fields, shall provide extra prior for segmenting indistinguishable objects. Note that existing works [23], [24] only regard context as spatial correlations. We expand and explicitly disentangle the concept of context into three domains, *i.e.*, *feature*, *space*, and *instance*. The implication is that when detecting and segmenting objects, this

model may augment the visual information by: 1) balancing the heterogeneous features, 2) fusing information from the background or other correlative objects, and 3) adaptively refining intermediate representations for each instance and task. These three different domains are capable of modeling global visual context from coarse to fine at different granularities, capturing more discriminative object information.

The proposed framework intends to leverage three plug-and-play modules and construct the aforementioned context aggregation pipeline. Fig. 1 shows a glimpse of object features processed by these modules. In the feature domain, we argue that, in the feature pyramid built up by a backbone, flexible information flows may reduce information confounding and handle multi-scale features more effectively. It is based on this analysis that a dense feature pyramid network (DenseFPN) is proposed to enable adaptive feature propagation. This module has a pyramid structure with stackable basic blocks consisting of top-down and bottom-up paths. We adopt inter-level residual connections [25], cross-level dense connections [26], and feature re-weighting strategy to enable the module to learn its optimal feature propagation manner. In the spatial domain, long-range spatial dependencies in remote sensing images bring more complementary information to blurry objects than in natural scenes. So the spatial context pyramid (SCP) is adopted to capture global spatial context in each feature pyramid level. This module learns to aggregate features from the whole feature map, and combines them into each pixel using adaptive weights. Such a strategy guarantees that only useful global information is fused into local regions, without decreasing the discrepancies among objects. As for the instance domain, we argue that object representations should be adaptively refined for each instance and downstream task. For example, performing object classification needs an overall view, while segmentation requires more zoomed details. The demand for different sizes of receptive fields also varies among instances. Hence, we introduce hierarchical region of interest extractor (HRoIE) to generate RoI features per instance and task. After cropping the instance feature maps from all levels, this module starts from the highest or lowest scale, and fuses the features level-by-level in a hierarchical manner. Pixel-wise attention mechanism is exploited to combine neighboring feature maps. These modules are lightweight while having the flexibility for scalable model design. Overall, the main contributions of this paper are three-fold:

- The expansion and explicit disentangling of the concept of context into *feature*, *spatial*, and *instance* domains have resulted in superior performance in both optical remote sensing image and SAR image segmentation. To the best of our knowledge, this is the first work that considers global visual context beyond spatial dependencies.
- The proposed CATNet is capable of utilizing DenseFPN, SCP, and HRoIE to learn and aggregate the global visual context from different domains for object detection and instance segmentation in remote sensing images.
- The proposed scheme has been tested on a wide variety of datasets, including iSAID, DIOR, NWPU VHR-10, and HRSID, and new state-of-the-art performance has been

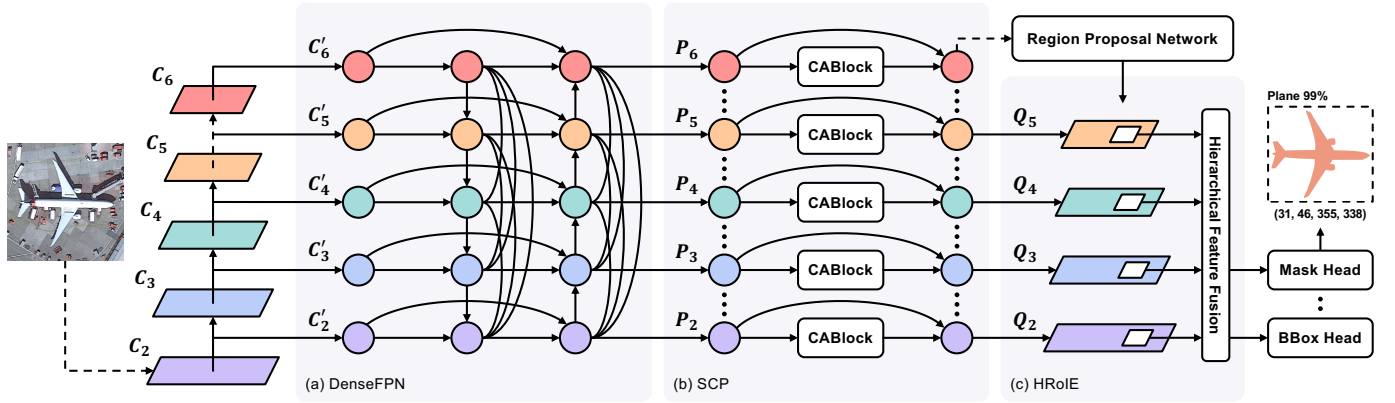


Fig. 3. Overall architecture of the proposed framework. The process of global context aggregation is realized by three modules, namely a) dense feature pyramid network, b) spatial context pyramid, and c) hierarchical region of interest extractor. These modules are designed to aggregate global context information from different feature pyramid levels, spatial positions, and receptive fields at feature, spatial, and instance domains, respectively.

obtained with similar computational costs.

The rest of this paper is organized as follows. Related works and comparisons are discussed in Section II. Detailed formulations of DenseFPN, SCP, and HRoIE are introduced in Section III. Section IV presents extensive experimental results and in-depth analysis on public datasets. Finally, concluding remarks are summarized in Section V.

## II. RELATED WORK

### A. Instance Segmentation in Remote Sensing Images

Instance Segmentation is a challenging and broadly studied problem in computer vision. Similar to object detection [7], [8], the majority of instance segmentation approaches can be divided into two schemes, namely one-stage methods and two-stage methods. As a straightforward design, one-stage methods [12], [27] adopt the bottom-up strategy that performs semantic segmentation at image level, and further separates individual objects using clustering or metric learning. These methods often possess considerable efficiencies, but are largely restricted by their localization accuracy. Compared to this paradigm, two-stage methods [9]–[11], [28] separate the segmentation pipeline into two phases, *i.e.*, region proposal generation and task-specific post-processing, resulting in a top-down style. Benefiting from two-time bounding box regression, these methods usually achieve better results on object localization and mask prediction. Some recent works [18]–[21], [29]–[34] try to tackle the problem of scene classification and object detection in remote sensing images, but they do not pay particular attention to instance segmentation. Our proposed context aggregation strategy can be integrated into both one-stage and two-stage methods, while HRoIE is not used in one-stage methods since it's not necessary to crop the feature maps. Further experimental results demonstrate that our modules can steadily boost performances.

### B. Multi-Scale Feature Propagation

Scale variation is a long-standing challenge in most visual recognition tasks, and it is a common solution to leverage a pyramid structure to represent the visual features under

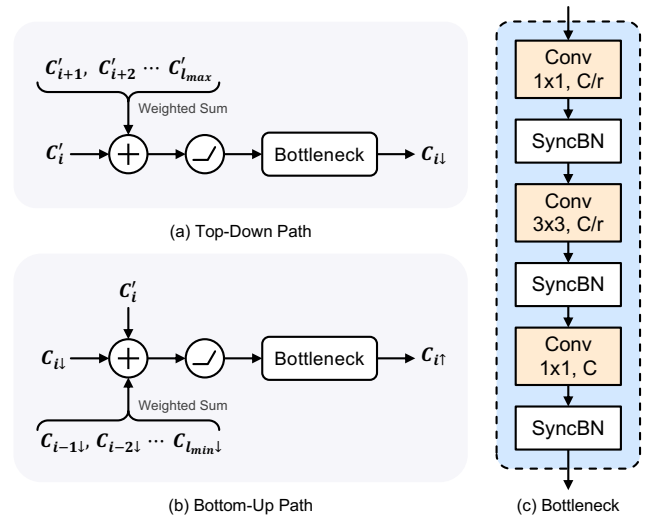


Fig. 4. Detailed procedure of multi-scale feature propagation in dense feature pyramid network.  $\oplus$  and  $\odot$  denote element-wise addition and ReLU activation, respectively. The cross-level feature maps are adaptively combined using a feature re-weighting strategy.

different resolutions. In the area of dense prediction, FPN [35] is the first work that builds up a feature pyramid, and propagates the information among different levels. Such a design rapidly became a standard for most instance segmentation models. However, FPN only propagates the features using a top-down path, which is suboptimal for multi-level feature fusion. Thus PAFPN [28] was proposed to incorporate an extra bottom-up propagation path, and FPG [36] further introduced a multi-pathway feature pyramid that can better capture cross-level information. MHN [37] tackles the semantic gap problem by leveraging semantic feature maps. A similar multi-branch structure is also used in TridentNet [38] and NETNet [39] for scale-aware training and generating scale-aware features. From the perspective of basic operators, scale-equalizing pyramid convolution (SEPC) [40] tends to fuse the feature maps using 3D convolutions. More recently, recursive feature pyramid (RFP) [41] re-uses the backbone to capture deeper semantics via feedback connections from FPN. In order to reduce the computational costs, NAS-FPN [42]

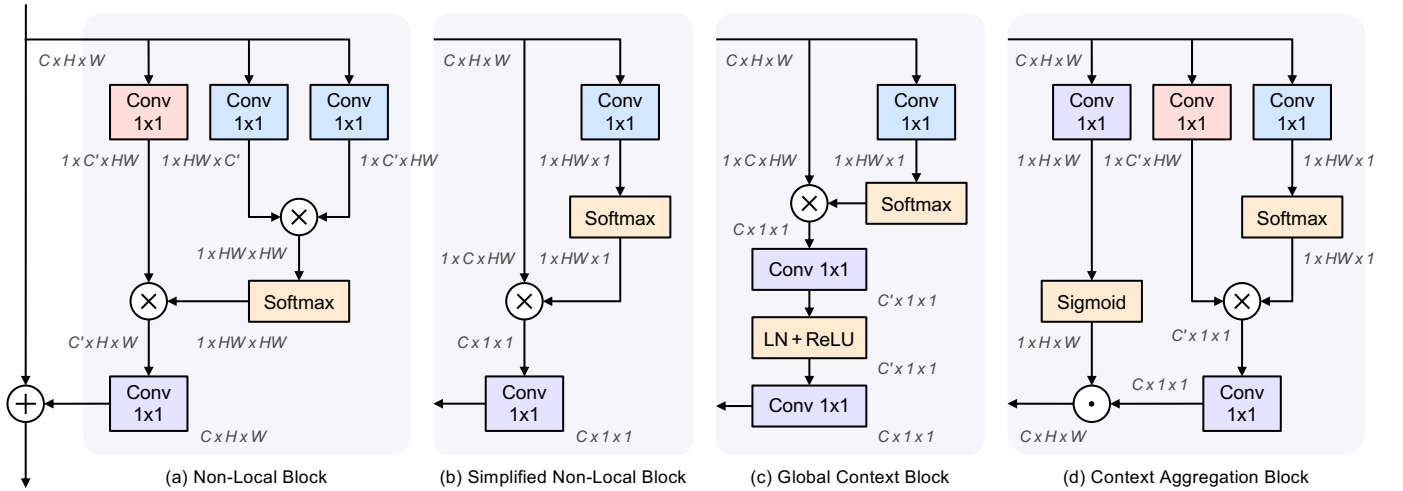


Fig. 5. Designs of spatial context modules. Permutations of feature maps are represented as their dimensions, e.g.,  $C \times H \times W$  indicates a matrix with number of channels  $C$ , height  $H$ , and width  $W$ .  $\otimes$ ,  $\odot$ , and  $\oplus$  denote batched matrix multiplication, broadcast hadamard product, and broadcast element-wise addition, respectively. Convolution layers used for attention map generation, feature mapping, and context refinement are annotated as blue, red, and purple.

is obtained by introducing neural architecture search (NAS) [43] to optimize the architecture of FPN, and BiFPN [44] is constructed by carefully designing the feature fusion blocks. We argue that compared with hard-coded feature aggregation paths, flexible information flows may reduce confounding and handle multi-scale features more effectively, so that DenseFPN is proposed to learn the optimal feature aggregation strategy during training.

### C. Global Context Modeling

One of the most representative properties of convolutional neural networks (CNNs) is local dependency modeling. The receptive fields can only be enlarged by stacking multiple convolution layers. Although feature pyramids have been introduced to capture multi-scale information, long-range spatial dependency modeling has also been proven to be effective for dense prediction tasks [23], [45]. As a pathbreaking work, non-local neural network (NLNet) [23] shows that global spatial context can be aggregated by computing pixel-level pairwise correlations, but it suffers from the problem of high computational cost. Therefore, some extensions of NLNet tackle this problem by simplifying the pairwise correlation computation. For example, criss-cross attention (CCNet) [46] aggregates global context by fusing the information along axes twice. Xu *et al.* observed that the per-pixel attention maps in NLNet are almost the same for different positions, so that GCNet [47] is proposed to generate a single attention map for a feature map. GANet [48] acts as a unified attention module for global context modeling. Although promising results have been achieved by these works, all these methods merely consider context as long-range spatial correlations, ignoring the dependencies in feature and instance domains.

## III. THE PROPOSED APPROACH

In this section, we introduce our approach on global context aggregation. As shown in Fig. 3, the entire framework can be divided into three sub-modules, namely DenseFPN, SCP, and

HRoIE. These modules aim to aggregate global context information in order from feature, spatial, and instance domains, respectively.

### A. Overview

Given an image  $x$  and a set of object categories of interest  $S = \{1, \dots, N\}$ , the task of instance segmentation aims to detect and segment all the objects in  $x$ , where they belong to whichever the pre-defined categories. The output of instance segmentation would be a collection of tuples  $\mathcal{T} = \{\langle b, m, s \rangle\}$ , where  $b \in \mathbb{R}^4$  denotes the bounding box of the object,  $m$  represents a binary mask in which  $m_{i,j} \in \{0, 1\}$  indicates whether the pixel  $(i, j)$  belongs to this object, and  $s \in S$  is a one-hot vector describing the object category. Note that a single object may be presented by separate masks.

We adopt Mask R-CNN [9], a common two-stage instance segmentation framework, as our baseline. The whole pipeline is constructed by extracting visual features, generating region proposals, and performing bounding box regression, object classification as well as mask prediction on each proposal. A heterogeneous feature pyramid is first built by extracting visual features from each stage of the backbone. In order to make the features more discriminative, we exploit DenseFPN and SCP to propagate object information among different levels and regions. After enhancing the feature pyramid, task-specific RoI features are generated by HRoIE for each proposal. Details of these modules are introduced in the following sections.

### B. Dense Feature Pyramid Network

Multi-scale feature propagation aims to aggregate visual features from different backbone stages, that is given an input feature pyramid  $C = \{C_{l_1}, C_{l_2}, \dots\}$ , where  $C_i$  denotes the feature map from the stage  $i$ , the goal is to propagate the features among different levels to produce an enhanced feature pyramid  $P = \{P_{l_1}, P_{l_2}, \dots\}$ , in which the features are more informative for downstream tasks. Formally, the resolution of feature map  $C_i$  or  $P_i$  is  $1/2^i$  of the input image.

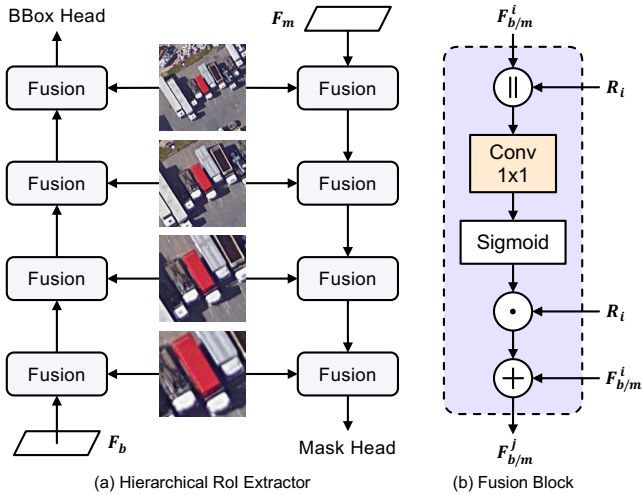


Fig. 6. Detailed structure of hierarchical region of interest extractor.  $\oplus$ ,  $\odot$ , and  $\otimes$  denote channel-wise concatenation, hadamard product, and element-wise addition, respectively. The RoI features are generated by fusing multi-scale feature maps along top-down or bottom-up paths in a hierarchical manner.

The basic architecture of DenseFPN is shown in Fig. 3 (a), where each node represents a feature map and the lines stand for information flows. This module takes  $C_2 \sim C_5$  as inputs and firstly down-samples them to 256 channels using  $1 \times 1$  convolutions, producing  $C'_2 \sim C'_5$ . An extra  $3 \times 3$  convolution with *stride* = 2 is applied to  $C'_5$  to generate  $C'_6$ . So that  $C'_2 \sim C'_6$  are with the same number of channels but different resolutions. Subsequently, these features are passed through several stacked basic blocks for feature-level context aggregation. In each block, the input feature pyramid is processed by a top-down and a bottom-up aggregation path, in which inter-level residual connections [25], cross-level dense connections [26], and feature re-weighting strategy are adopted.

Fig. 4 illustrates the detailed feature propagation strategy in basic blocks. In the top-down path, output features  $C_{i\downarrow}$  of each feature pyramid level are generated by fusing the features from the current level and all upper levels, then performing a parameterized transform upon the fused features.

$$C_{i\downarrow} = \text{Bottleneck}\left(C'_i + \sum_{j=i+1}^{l_{max}} [\text{Resize}(C'_j) \cdot w_{i\downarrow}^j]\right) \quad (1)$$

Here,  $\text{Bottleneck}(\cdot)$  denotes a ReLU activation layer followed by a  $3 \times 3$  bottleneck [25] without activations. We observe that adopting only one nonlinearity before the bottleneck structure brings better performance.  $\text{Resize}(\cdot)$  represents a max pooling layer, and  $w_{i\downarrow}^j$  is a learnable re-weighting term for aggregating features from level  $j$  to level  $i$ . The weights  $w_{i\downarrow}$  are vectors with lengths correspond to their levels, the values are normalized from raw values using softmax by

$$w_{i\downarrow}^j = \frac{\exp(v_{i\downarrow}^j)}{\sum_{k=1}^{N_i} \exp(v_{i\downarrow}^k)} \quad (2)$$

where  $v_{i\downarrow}$  denotes the raw weight vector and  $j$  is the index of each element. Using the normalization above can stabilize the learning process. Similarly to the top-down path, bottom-up

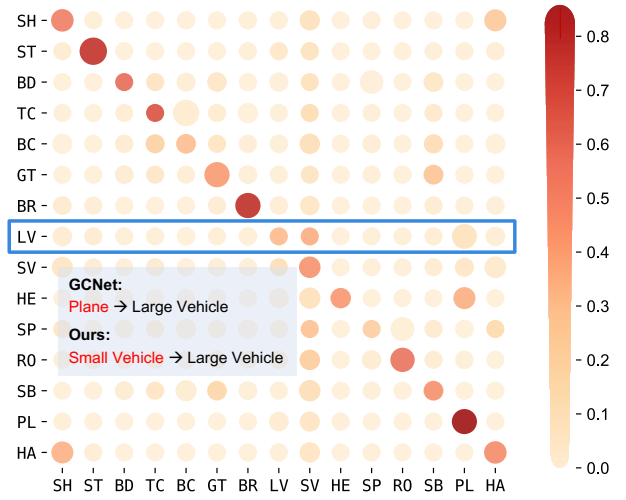


Fig. 7. Visualization of feature aggregation weights in iSAID dataset. X-axis refers to the categories that used to construct the global context, while the y-axis denotes the categories that receive features. The weights in GCNet and SCP are represented as sizes and depths of colors, respectively.

features  $C_{2\uparrow} \sim C_{6\uparrow}$  are computed by

$$C_{i\uparrow} = \text{Bottleneck}\left(C'_i + C_{i\downarrow} + \sum_{j=l_{min}}^{i-1} [\text{Resize}(C_{j\downarrow}) \cdot w_{i\uparrow}^j]\right) \quad (3)$$

where  $\text{Resize}(\cdot)$  represents a bilinear interpolation layer and other notations are consistent with Eq. 1. We adopt residual connections to preserve the original features and prevent gradient vanishing. Leveraging the flexible architecture and feature re-weighting strategy, DenseFPN has the capability to optimize the information flow of context aggregation in feature domain during training.

### C. Spatial Context Pyramid

After aggregating feature maps across different levels, the feature pyramid remains to contain spatially local information, thus we introduce spatial context pyramid (SCP) to further augment the features by learning global spatial context within each level. Former attempts in this area [23], [45]–[48] normally integrate several spatial or channel attention blocks into the backbone to enable global receptive fields. Some architectures of these blocks are presented in Fig. 5. Among these methods, non-local neural network (NLNet) [23] is a general solution that computes pixel-wise correlations via embedded gaussian for every spatial position, while squeeze-and-excitation network (SENet) [45] tackles the problem from the perspective of channel attentions. To combine the spatial and channel modulation abilities, global context network (GCNet) [47] is proposed to learn a single attention map for an image. However, we observe that in remote sensing images with objects only covering small areas, this design may bring too much useless background information to objects. To tackle this problem, we propose adding an extra path on top of this structure to learn the informativeness of each pixel. Our core idea is that if the features of a pixel are informative enough,

TABLE I  
OBJECT DETECTION AND INSTANCE SEGMENTATION RESULTS ON ISAID TEST SET.

Method	Backbone	AP <sub>b</sub>	AP <sub>b</sub> <sup>50</sup>	AP <sub>b</sub> <sup>75</sup>	AP <sub>b</sub> <sup>S</sup>	AP <sub>b</sub> <sup>M</sup>	AP <sub>b</sub> <sup>L</sup>	AP <sub>m</sub>	AP <sub>m</sub> <sup>50</sup>	AP <sub>m</sub> <sup>75</sup>	AP <sub>m</sub> <sup>S</sup>	AP <sub>m</sub> <sup>M</sup>	AP <sub>m</sub> <sup>L</sup>
Mask R-CNN [9]	ResNet-50	42.2	63.1	47.8	44.7	49.9	18.9	36.7	59.7	39.7	39.5	45.1	11.5
	ResNet-101	43.9	65.5	49.7	46.4	53.4	18.8	38.3	61.7	40.9	41.0	48.4	13.4
CenterMask [49]	ResNet-50	44.6	65.6	49.7	47.4	48.7	21.4	37.4	61.4	39.4	40.2	43.6	15.0
	ResNet-101	44.4	65.6	49.2	47.5	48.2	17.3	37.2	61.0	38.9	40.3	42.7	11.9
BlendMask [50]	ResNet-50	45.1	66.1	50.1	47.5	50.6	21.9	38.2	62.4	40.0	40.6	46.9	18.0
	ResNet-101	45.2	65.7	50.1	47.9	51.6	22.9	38.3	62.5	40.0	<b>48.9</b>	47.4	17.9
DB-BlendMask [51]	ResNet-50	46.0	66.8	50.8	48.8	50.7	21.2	39.2	63.5	41.4	42.0	46.5	16.3
	ResNet-101	46.0	66.8	50.8	48.8	52.3	22.8	38.9	63.1	41.0	41.6	48.1	17.8
Cascade R-CNN [52]	ResNeXt-101	45.3	65.8	51.1	48.3	53.7	22.8	37.0	60.7	39.0	39.8	46.6	15.1
HTC [10]	ResNet-152	46.3	66.4	51.8	48.5	56.4	25.3	38.3	61.6	40.9	40.6	50.0	17.4
ISDNet [53]	ResNet-152	46.8	67.1	52.4	49.0	57.4	27.3	38.7	62.2	42.2	41.1	51.0	19.8
MS R-CNN [54]	ResNet-101	–	–	–	–	–	–	37.0	57.8	40.5	39.7	46.0	14.3
SCNet [55]	ResNet-101	–	–	–	–	–	–	38.1	60.4	41.2	40.9	46.9	12.6
CPISNet [22]	AFEN-4GF	–	–	–	–	–	–	39.1	62.2	42.5	41.8	49.6	17.6
<b>CATNet</b>	ResNet-50	48.0	66.1	53.9	51.3	57.6	22.6	39.9	62.8	43.5	43.1	50.2	13.4
<b>CATNet + Aug.</b>	ResNet-50	<b>50.7</b>	<b>69.4</b>	<b>57.1</b>	<b>54.3</b>	<b>59.9</b>	<b>33.8</b>	<b>41.9</b>	<b>66.0</b>	<b>45.5</b>	45.2	<b>52.8</b>	<b>25.7</b>

AP<sub>b</sub> - Bounding box AP, AP<sub>m</sub> - Mask AP, Aug. - Multi-Scale Training

TABLE II  
CLASS-WISE INSTANCE SEGMENTATION RESULTS ON ISAID TEST SET.

Method	Backbone	AP <sub>b</sub>	AP <sub>m</sub>	SH	ST	BD	TC	BC	GT	BR	LV	SV	HE	SP	RO	SB	PL	HA
Mask R-CNN [9]	ResNet-50	42.2	36.7	46.8	35.4	48.0	75.7	48.2	28.1	17.4	30.4	15.6	13.4	38.1	43.2	30.9	42.9	31.0
	ResNet-101	43.9	38.3	47.0	32.0	53.1	76.7	48.2	30.4	18.6	31.1	16.1	14.5	39.0	46.4	35.7	44.1	31.1
CenterMask [49]	ResNet-50	44.6	37.4	46.9	37.7	54.0	76.7	50.5	23.6	16.7	32.5	16.1	10.0	38.1	41.5	35.9	42.2	29.9
	ResNet-101	44.4	37.2	46.4	36.9	54.2	75.8	50.1	25.7	16.5	32.2	15.6	10.5	37.2	42.2	34.6	41.2	29.9
BlendMask [50]	ResNet-50	45.1	38.2	46.9	37.4	55.5	76.4	50.9	21.3	15.9	33.2	15.5	11.0	38.4	45.3	38.6	43.8	33.5
	ResNet-101	45.2	38.3	47.1	38.0	56.4	77.3	51.2	22.4	16.7	33.4	16.2	11.8	38.2	45.5	36.0	42.6	34.0
DB-BlendMask [51]	ResNet-50	46.0	39.2	48.6	38.2	56.3	78.3	52.2	23.7	17.6	33.6	<b>16.4</b>	12.4	38.1	44.7	36.9	44.7	<b>37.4</b>
	ResNet-101	46.0	38.9	48.6	37.3	54.8	77.5	51.2	23.8	17.8	<b>34.1</b>	16.3	12.2	38.4	44.4	36.3	45.2	36.1
Cascade R-CNN [52]	ResNeXt-101	45.3	37.0	46.7	36.7	54.7	76.4	52.7	20.6	18.1	30.5	14.1	11.4	39.3	45.8	37.3	41.8	28.4
HTC [10]	ResNet-152	46.3	38.3	47.3	37.9	53.5	76.4	51.0	32.5	19.0	31.2	16.3	11.4	38.4	45.6	39.6	43.1	30.7
ISDNet [53]	ResNet-152	46.8	38.7	47.8	38.7	54.0	76.3	52.4	32.8	19.3	31.8	16.6	12.0	39.0	45.0	40.1	43.3	31.0
MS R-CNN [54]	ResNet-101	–	37.0	46.6	33.9	54.2	76.1	49.9	29.7	17.7	30.0	14.0	11.8	37.9	44.4	35.7	43.8	30.1
SCNet [55]	ResNet-101	–	38.1	48.0	35.9	56.6	77.0	51.5	30.2	18.7	31.7	14.2	9.7	39.2	46.7	36.3	45.1	31.1
CPISNet [22]	AFEN-4GF	–	39.1	<b>49.5</b>	35.8	54.3	77.6	52.9	31.9	19.9	32.9	14.9	13.2	<b>40.6</b>	45.2	39.3	46.2	32.7
<b>CATNet</b>	ResNet-50	48.0	39.9	49.4	36.7	55.9	77.7	55.9	29.7	20.0	31.5	15.1	14.8	40.1	43.5	39.0	<b>47.0</b>	35.1
<b>CATNet + Aug.</b>	ResNet-50	<b>50.7</b>	<b>41.9</b>	48.9	<b>39.4</b>	<b>60.9</b>	<b>79.2</b>	<b>56.7</b>	<b>33.8</b>	<b>21.2</b>	31.6	15.6	<b>17.8</b>	39.7	<b>48.6</b>	<b>45.8</b>	46.6	36.5

AP<sub>b</sub> - Bounding box AP, AP<sub>m</sub> - Mask AP, Aug. - Multi-Scale Training

there's not much need to aggregate features from other spatial positions. Such a soft re-weighting strategy can effectively fuse global features while reducing information confounding.

The architecture of SCP is shown in Fig. 3 (b). This module also has a pyramid structure, thus can be easily inserted after the backbone or neck. Each layer consists of a context aggregation block (CABlock) with a residual connection. The detailed design of this block is presented in Fig. 5 (d). In each block, pixel-wise spatial context is aggregated by

$$Q_i^j = P_i^j + a_i^j \cdot \sum_{j=1}^{N_i} \left[ \frac{\exp(w_k P_i^j)}{\sum_{m=1}^{N_i} \exp(w_k P_i^m)} \cdot w_v P_i^j \right] \quad (4)$$

where  $P_i$  and  $Q_i$  denote the input and output feature maps of level  $i$  in the feature pyramid, each contains  $N_i$  pixels.  $j, m \in \{1, N_i\}$  indicate the indices of each pixel.  $w_k$  and  $w_v$  are linear projection matrices for projecting the feature maps.

In practice, we use  $1 \times 1$  convolutions to perform the mapping. The formula above simplifies the widely used self-attention mechanism [67] by replacing the matrix multiplication between query and key with a linear projection, largely reducing the parameters and computational costs. Beyond GCNet, we apply  $a_i$ , a re-weighting matrix with the same shape as  $P_i$  and  $Q_i$ , to balance the extent of aggregating global spatial context for each pixel. This matrix can also be generated as simply as a linear projection from  $P_i$  with sigmoid activation.

$$a_i^j = \frac{1}{1 + \exp(-w_a P_i^j)} \quad (5)$$

Similarly,  $j \in \{1, N_i\}$  is the matrix index. Here, the output of sigmoid function shall be regarded as the ratio of information that should be aggregated from global context. We conducted extensive visualizations and experiments on the effectiveness



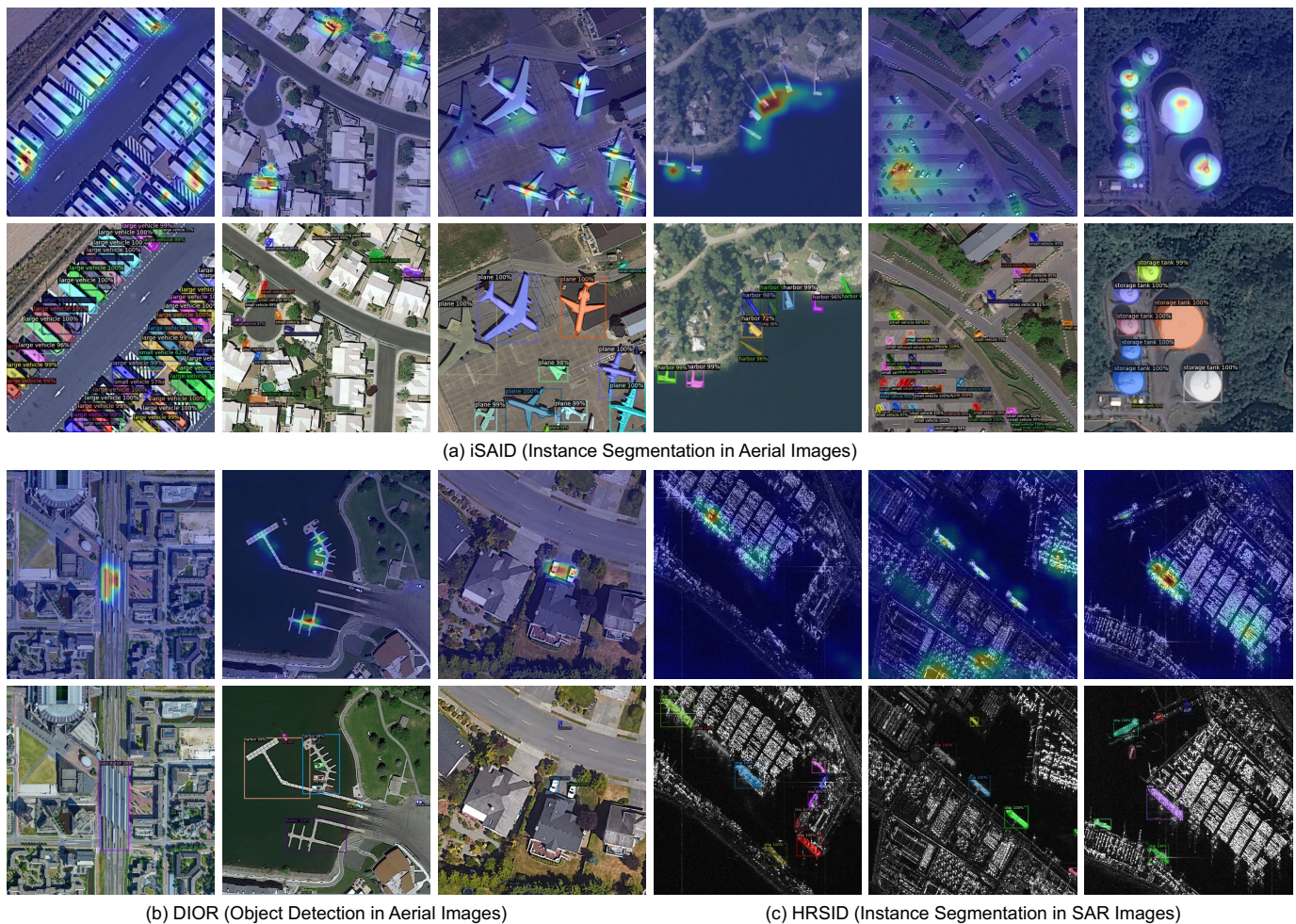


Fig. 8. Qualitative results on multiple datasets. We visualize the spatial context attention maps in SCP (odd rows) and the final outputs (even rows) on iSAID, DIOR, and HRSID datasets. The results demonstrate that our model has the ability to perform object detection and instance segmentation accurately in both optical remote sensing images and SAR images.

[28], [68] also proved that simply computing the sum of RoI features cropped from all layers achieves slightly better performance, indicating that leveraging multi-scale features might alleviate the information loss.

In this work, we further deal with this problem by proposing hierarchical region of interest extractor (HRoIE) to perform task-specific hierarchical feature fusion for each instance. This module is inserted after SCP as displayed in Fig. 3 (c). Our hypothesis is that humans can easily perform object detection and segmentation because they focalize their attention on objects in a hierarchical manner. For example, when a person tries to classify an object, he or she would look at the object itself at first. If the object’s appearance is not discriminative, the person would gradually look at the surrounding things to gain better information. On the opposite, when segmenting an object at pixel level, a human would look at the whole object to have a comprehensive understanding of its shape, and then zoom in to obtain detailed boundary information for accurate segmentation. We implement the idea above by cropping the features of proposals  $R_i$  from all feature pyramid levels in  $Q_i$  using RoIAlign [9], and utilize several attention blocks to fuse the features hierarchically and adaptively per instance and

task. As shown in Fig. 6, for each task, the RoI features are initialized from an empty matrix and combined with features from different levels in a hierarchical manner by

$$F_{b/m}^j = F_{b/m}^i + R_i^j \cdot \text{Sigmoid}([F_{b/m}^i \parallel R_i^j] \cdot w_i) \quad (6)$$

Here,  $R_i$  denotes the cropped features at level  $i$ ,  $F_{b/m}^i$  and  $F_{b/m}^j$  represent the aggregated RoI features at different levels,  $w_i$  is the linear projection weight, and  $\parallel$  means matrix concatenation at channel dimension. The procedure above computes the pixel-wise attention weights for feature aggregation, thus RoI features can be generated adaptively per instance and task. In practice, we adopt a bottom-up path for object detection head and a top-down path for mask prediction head.

#### IV. EXPERIMENTS

In this section, we extensively evaluate the proposed method on iSAID, DIOR, NWPU VHR-10, and HRSID datasets, in which the first three datasets are based on optical remote sensing images, while the last one is for synthetic aperture radar (SAR) images. Note that DIOR only provides bounding box annotations, thus we only evaluate the object detection performance on this dataset.



TABLE VI  
CLASS-WISE INSTANCE SEGMENTATION RESULTS ON ISAID VAL SET. ALL THE MODELS USE RESNET-50 AS BACKBONE.

Method	AP <sub>b</sub>	AP <sub>m</sub>	SH	ST	BD	TC	BC	GT	BR	LV	SV	HE	SP	RO	SB	PL	HA
<b>Mask R-CNN</b> [9]	43.8	36.5	47.2	36.8	54.0	76.7	32.3	32.7	22.9	39.4	15.3	4.9	30.1	37.3	46.4	48.1	26.8
+ DenseFPN	45.5	37.7	<b>48.8</b>	37.7	55.6	77.5	35.8	33.4	23.8	39.6	16.3	5.7	29.9	<b>39.8</b>	46.7	49.2	28.4
+ SCP	45.9	37.8	48.5	38.2	<b>55.8</b>	<b>77.8</b>	38.6	32.5	25.0	<b>39.8</b>	16.2	5.9	28.9	39.3	46.4	49.6	27.9
+ HRoIE	<b>46.2</b>	<b>38.5</b>	48.4	<b>38.4</b>	55.4	<b>77.8</b>	<b>40.4</b>	<b>35.3</b>	<b>25.8</b>	39.5	<b>16.4</b>	<b>6.3</b>	<b>31.0</b>	38.3	<b>47.8</b>	<b>50.2</b>	<b>29.0</b>
<b>Cascade R-CNN</b> [52]	45.8	37.3	47.7	37.0	55.6	77.4	35.2	33.7	22.6	40.2	15.9	5.9	29.8	38.6	47.2	47.5	27.6
+ DenseFPN	46.8	38.2	<b>49.6</b>	38.9	55.1	<b>78.8</b>	34.4	34.4	25.1	<b>41.4</b>	<b>16.8</b>	4.8	<b>31.2</b>	41.0	45.9	49.9	28.4
+ SCP	47.4	38.4	49.5	39.0	<b>55.9</b>	77.9	35.6	34.8	25.3	40.7	16.7	5.7	31.1	41.2	<b>47.5</b>	49.7	28.8
+ HRoIE	<b>47.8</b>	<b>38.9</b>	49.1	<b>39.2</b>	55.7	78.3	<b>37.4</b>	<b>35.5</b>	<b>25.4</b>	41.2	<b>16.8</b>	<b>8.7</b>	<b>31.2</b>	<b>41.8</b>	47.1	<b>50.0</b>	<b>29.9</b>
<b>MS R-CNN</b> [54]	44.0	37.8	48.6	36.6	54.9	77.0	36.9	35.2	22.5	39.9	15.5	9.7	30.1	39.8	44.7	49.2	29.7
+ DenseFPN	45.4	38.5	49.4	37.5	55.3	77.4	39.2	35.4	24.8	<b>41.3</b>	16.3	6.9	31.4	37.6	47.1	50.6	29.5
+ SCP	45.9	38.9	49.2	<b>37.9</b>	56.2	77.8	<b>40.8</b>	35.4	24.5	40.9	16.2	8.6	31.8	40.5	46.6	50.3	29.7
+ HRoIE	<b>46.1</b>	<b>39.3</b>	<b>49.5</b>	37.8	<b>56.6</b>	<b>77.9</b>	37.5	<b>37.7</b>	<b>25.0</b>	41.2	<b>16.4</b>	<b>9.8</b>	<b>32.0</b>	<b>40.6</b>	<b>50.1</b>	<b>51.0</b>	<b>30.2</b>
<b>PointRend</b> [11]	43.4	37.4	47.8	36.6	55.6	76.8	34.3	33.3	23.0	41.2	15.8	5.0	29.9	38.0	47.4	49.8	29.8
+ DenseFPN	44.6	38.3	49.4	39.2	54.7	77.6	35.7	33.4	23.8	41.7	17.0	5.5	<b>30.9</b>	37.7	47.6	52.0	<b>31.7</b>
+ SCP	45.2	38.7	50.9	39.4	55.7	78.4	36.6	33.7	25.0	<b>43.1</b>	<b>16.8</b>	5.0	30.1	39.2	46.2	52.7	31.4
+ HRoIE	<b>45.7</b>	<b>39.3</b>	<b>51.6</b>	<b>39.7</b>	<b>55.8</b>	<b>78.6</b>	<b>37.7</b>	<b>36.0</b>	<b>25.6</b>	42.8	<b>16.8</b>	<b>5.7</b>	29.8	<b>39.9</b>	<b>47.7</b>	<b>53.1</b>	31.6

AP<sub>b</sub> - Bounding box AP, AP<sub>m</sub> - Mask AP

### A. Datasets and Evaluation Metrics

- **iSAID** [15] is a large-scale dataset for instance segmentation in aerial images. The images in iSAID are inherited from DOTA [14], which is popular for oriented object detection. It contains 15 classes of 655,451 instances in 2,806 images, with all the objects independently annotated from scratch. The spatial resolutions of images are in a large range between 800 and 13,000. Following previous works [22], [51], [53], we split them into  $800 \times 800$  patches with a stride of 200 for fair benchmarking with existing methods. When conducting detailed comparisons and ablation studies, we adopt a smaller patch size of  $512 \times 512$  with 128 stride to reduce the training cost. The abbreviations of classes are SH - ship, ST - storage tank, BD - baseball diamond, TC - tennis court, BC - basketball court, GT - ground track field, BR - bridge, LV - large vehicle, SV - small vehicle, HE - helicopter, SP - swimming pool, RO - roundabout, SB - soccer ball field, PL - plane, and HA - harbor.
- **DIOR** [16] is a complex aerial images dataset labeled by horizontal and oriented bounding boxes. It contains 23,463 images with 190,288 instances, covering 20 object classes. Object sizes in DIOR have severe inter-class and intra-class variabilities. The complexity of this dataset is also reflected in different imaging qualities, weathers, and seasons. The abbreviations of classes are AL - airplane, AR - airport, BF - baseball field, BC - basketball court, BR - bridge, CH - chimney, DA - dam, ES - expressway service area, ET - expressway toll station, GC - golf field, GT - ground track field, HA - harbor, OV - overpass, SH - ship, ST - stadium, SA - storage tank, TC - tennis court, TS - train station, VE - vehicle, and WM - wind mill.
- **NWPU VHR-10** [13] was originally designed for object detection in aerial images, and has been enriched with instance-level mask annotations [64]. It contains 10 object categories in 800 high-resolution images, among which 650 are positive and 150 are negative without any objects of interest. Since this dataset has no official train/test split available, we randomly

select 70% of the images for training while the others for testing. Our splits will be released for fair comparison with the following methods. The abbreviations of classes are AI - airplane, SH - ship, ST - storage tank, BD - baseball diamond, TC - tennis court, BC - basketball court, GT - ground track field, HA - harbor, BR - bridge, and VE - vehicle.

- **HRSID** [17] is a recently introduced dataset for ship detection and segmentation in synthetic aperture radar (SAR) images. This dataset contains a total of 5,604 high-resolution SAR images with 16,951 ship instances. All the instances in this dataset are annotated with pixel-level masks. Spatial resolutions of the images are  $0.5m$ ,  $1m$ , and  $3m$ .

We follow the standard evaluation metric that utilizes mean average precision (mAP) to measure the detection and segmentation performances. A prediction is considered a true positive (TP) when the bounding box or mask of the object has an intersection over union (IoU) with its corresponding ground truth greater than a threshold  $\theta_{IoU}$ , and the predicted class label is correct. For iSAID, NWPU VHR-10, and HRSID datasets, we compute the mean of mAPs under  $\theta_{IoU}$  ranging from 0.05 to 0.95. For DIOR dataset, only the mAPs under  $\theta_{IoU} = 0.5$  are considered following the original paper. Aside from the accuracy metrics, we also evaluate the efficiencies of the models using the number of parameters (#Params) and floating-point operations (FLOPs), measuring the storage and computational efficiencies, respectively. We consider these two metrics since they are both hardware-irrelevant and have been widely adopted in previous studies.

### B. Implementation Details

We choose Mask R-CNN [9], Faster R-CNN [7], and RetinaNet [8] with ResNet-50 [25] backbone as our baselines for different scenarios. The backbone is pre-trained on ImageNet [73] and fine-tuned with the detectors. All the parameters in the first stage are frozen after pre-training. If not specified, 5 basic blocks are included in all DenseFPN modules. In order to stabilize the training process, synchronized batch

TABLE VII  
DETAILED COMPARISONS ON iSAID VAL SET. THE BASELINE MODEL IS MASK R-CNN WITH RESNET-50-FPN AS BACKBONE AND NECK.

(a) Types of Multi-Scale Feature Propagation Modules							(b) Types of Spatial Context Modules						
Method	Chn.	Depth	AP <sub>b</sub>	AP <sub>m</sub>	#Params	FLOPs	Method	Position	Red.	AP <sub>b</sub>	AP <sub>m</sub>	#Params	FLOPs
Baseline	256	–	43.8	36.5	44.05M	114.96G	Baseline	–	–	43.8	36.5	44.05M	114.96G
PAFPN [28]	256	–	44.4	37.0	47.59M	121.31G	GABlock [48]	C <sub>3</sub> ~ C <sub>5</sub>	–	44.5	37.0	143.25M	156.52G
HRFPN [69]	256	–	44.1	36.8	44.64M	129.09G	CCBlock [46]	C <sub>3</sub> ~ C <sub>5</sub>	–	44.1	36.7	68.97M	132.43G
CARAFE [70]	256	–	44.1	36.9	49.65M	115.71G	NLBlock [23]	C <sub>3</sub> ~ C <sub>5</sub>	2	44.2	37.0	83.93M	142.92G
BFP [71]	256	–	44.2	37.0	44.31M	115.23G			4	44.0	36.8	63.99M	128.95G
AugFPN [72]	256	–	44.3	36.8	45.82M	115.03G			8	43.8	36.7	54.03M	121.97G
TridentNet [38]	256	–	42.6	32.9	35.28M	838.99G			16	43.5	36.5	49.04M	118.47G
RFP [41]	256	–	45.3	37.4	68.80M	153.88G	GCBlock [47]	C <sub>3</sub> ~ C <sub>5</sub>	2	44.5	37.2	64.02M	115.00G
BiFPN [44]	256	–	44.1	36.8	45.89M	112.49G			4	44.3	37.1	54.05M	114.99G
NAS-FPN [42]	256	1	44.2	36.9	45.89M	119.86G			8	44.1	36.9	49.06M	114.99G
		3	45.3	37.4	54.16M	155.30G			16	44.0	36.8	46.57M	114.98G
		5	45.5	37.6	62.42M	190.74G	CABlock (ours)	C <sub>3</sub> ~ C <sub>5</sub>	2	44.5	37.6	64.02M	121.99G
		7	45.6	37.7	70.69M	226.18G			4	44.4	37.4	54.05M	118.50G
FPG [36]	128	5	43.7	36.3	38.96M	92.75G			8	44.3	37.2	49.07M	116.75G
		7	44.2	36.6	41.66M	102.23G			16	44.2	37.1	46.58M	115.87G
		9	44.5	36.9	44.36M	111.70G	GABlock [48]	P <sub>2</sub> ~ P <sub>6</sub>	–	43.9	36.6	45.66M	119.20G
		11	44.6	37.1	47.05M	121.18G	CCBlock [46]	P <sub>2</sub> ~ P <sub>6</sub>	–	43.6	36.2	44.46M	116.76G
FPG [36]	256	5	44.2	36.8	60.79M	143.19G	NLBlock [23]	P <sub>2</sub> ~ P <sub>6</sub>	1	44.0	36.9	45.36M	120.71G
		7	44.8	37.1	71.55M	181.04G	GCBlock [47]	P <sub>2</sub> ~ P <sub>6</sub>	1	44.2	37.2	44.71M	114.97G
		9	45.1	37.3	82.32M	218.90G	<b>SCP (ours)</b>	P <sub>2</sub> ~ P <sub>6</sub>	1	<b>44.6</b>	<b>37.7</b>	44.71M	116.41G
		11	45.2	37.4	93.09M	256.75G	Red. - Channel Reduction Rate						
<b>DenseFPN (ours)</b>	128	1	44.0	36.5	34.87M	84.81G	(c) Orders of the Proposed Modules						
		3	44.4	36.8	35.94M	89.46G	<b>Module Order</b>	<b>AP<sub>b</sub></b>	<b>AP<sub>m</sub></b>	<b>#Params</b>	<b>FLOPs</b>		
		5	44.7	37.0	37.02M	94.10G	Backbone (w. CABlock) → DFPN	45.4	37.5	72.95M	155.67G		
		7	44.9	37.1	38.09M	98.75G	Backbone → SCP → DFPN	45.5	37.6	53.64M	150.09G		
<b>DenseFPN (ours)</b>	256	1	44.7	37.1	44.42M	111.45G	Backbone → DFPN → SCP	<b>45.9</b>	<b>37.8</b>	53.64M	150.09G		
		3	45.2	37.5	48.70M	130.05G	DFPN - DenseFPN						
		5	45.5	37.7	52.98M	148.64G	(e) Ablation Study and Efficiency Comparisons						
		7	<b>45.7</b>	<b>37.8</b>	57.26M	167.23G	<b>DFPN</b>	<b>SCP</b>	<b>HRoIE</b>	<b>Lite</b>	<b>AP<sub>b</sub></b>	<b>AP<sub>m</sub></b>	<b>#Params</b>

Chn. - Feature Channels

(d) Types of Region of Interest Extractors

Method	Direction	AP <sub>b</sub>	AP <sub>m</sub>	#Params	FLOPs
Baseline	–	43.8	36.5	44.05M	114.96G
GRoIE [68]	–	44.0	36.9	47.84M	574.92G
SUM	–	43.9	36.7	44.05M	114.96G
CONCAT	–	44.0	36.7	84.35M	188.18G
ATTENTION	–	43.9	36.8	46.15M	186.97G
<b>HRoIE (ours)</b>	↓ + ↓	44.2	37.1	45.10M	151.00G
<b>HRoIE (ours)</b>	↓ + ↑	44.1	37.0	45.10M	151.00G
<b>HRoIE (ours)</b>	↑ + ↑	44.3	36.9	45.10M	151.00G
<b>HRoIE (ours)</b>	↑ + ↓	<b>44.4</b>	<b>37.2</b>	45.10M	151.00G

↓ - Top-Down, ↑ - Bottom-Up

normalizations (SyncBN) [74] are used among intermediate layers. When testing, we adopt Soft-NMS [75] to suppress the duplicate results with IoU larger than 0.5 and a maximum of 1,000 predictions would be made for each image, since most objects are heavily overlapped in remote sensing images.

We use stochastic gradient descent (SGD) optimizer with initial learning rate 0.01, momentum 0.9, and weight decay 0.0001 to learn the parameters for all models. Each training batch contains 8 images. On iSAID dataset, we follow the standard 1× training schedule that drops the learning rate by 0.1 at epoch 8 and 11, and stop training at epoch 12 for efficient parameter tuning and ablation studies on validation set. When benchmarking on test set, we train our model on

both training and validation sets under a 3× schedule, which is a common setting in existing works [15], [51]. On DIOR, NWPU VHR-10, and HRSID datasets, we adopt the 3×, 6×, and 3× training schedules, respectively. All the models are trained with PyTorch [76] on a compute node with 2 Intel Xeon Platinum 8358 (2.6GHz) Processors, 512GB RAM, and 4 NVIDIA A100 (80G) GPUs. The training of our model on iSAID, DIOR, NWPU VHR-10, and HRSID datasets costs about 27, 8, 1, and 5 hours, respectively.

### C. Instance Segmentation on iSAID Dataset

We first evaluate our approach on optical remote sensing images. Table I shows the benchmark of instance segmentation

Lite - Lite version (w. 1-layer DenseFPN@128)

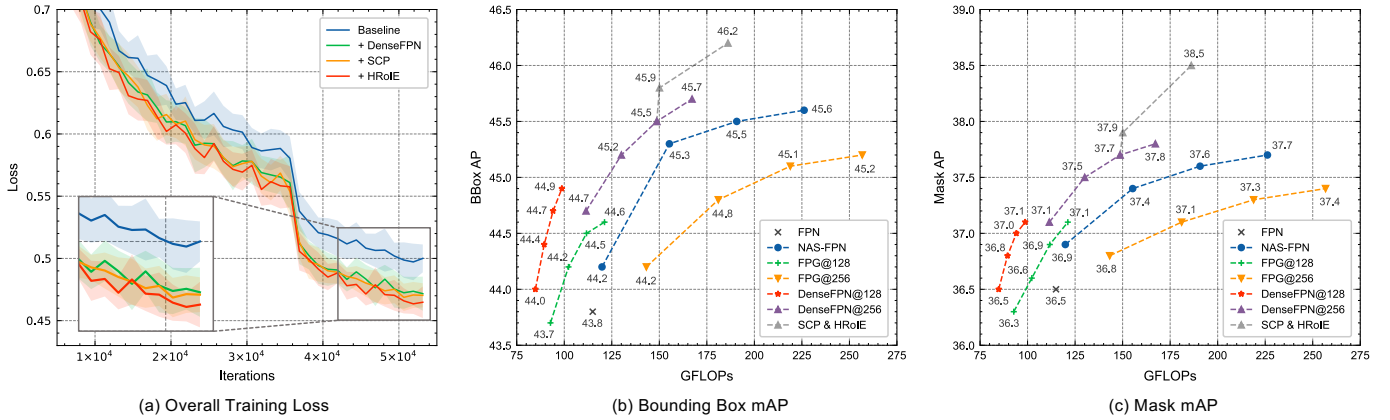


Fig. 9. Visualizations of training and test metrics on iSAID val set. a) displays the overall training losses with and without our proposed modules. b) and c) shows the performance comparison among different multi-scale propagation modules. Performances of each module are reported with 4 different depths.

methods on iSAID test set. The results are obtained by submitting the predictions onto the official evaluation server, which only accepts one submission per day, avoiding potential model tuning on the test set. Here,  $AP^{50}$  and  $AP^{75}$  indicate average precisions under  $\theta_{IoU} = 0.5$  and  $0.75$ , while  $AP^S$ ,  $AP^M$ , and  $AP^L$  represent the average precisions for small ( $10 \sim 144$  px), medium ( $144 \sim 1024$  px), and large (more than 1024 px) objects, respectively. The evaluation results of CATNet are obtained with two standard settings, with input images rescaled to  $(1400, 800)$  only and with augmentations on the short side at five scales ( $1200, 1000, 800, 600, 400$ ). Among existing methods, our proposed CATNet achieves state-of-the-art performances under both object detection and instance segmentation metrics, even with a 40% lighter backbone. When adopting multi-scale training, which is a common setting in previous works, our model can be further boosted significantly, demonstrating the scalability of our approach.

To justify the effectiveness of our method, class-wise instance segmentation results are also provided in Table II. Our model obtains better performances on most categories, including “storage tank”, “baseball diamond”, and “tennis court”. On other categories, our method is still comparable with most existing works.

#### D. Object Detection on DIOR Dataset

We then evaluate our approach on DIOR dataset. Note that this dataset only provides bounding box annotations for objects, we only validate the object detection performance on it. Such a comparison can also demonstrate the generalization ability of the modules – they could potentially benefit all the dense prediction tasks including object detection and instance segmentation. Table III presents the class-wise performance comparison with existing methods. Our methods are built upon two meta algorithms for one-stage and two-stage object detection, which are RetinaNet [8] and Faster R-CNN [7]. Note that one-stage models do not need RoI extractors, so that CATNet<sup>†</sup> only contains DenseFPN and SCP. The comparison shows that our proposed method outperforms all the previous methods, with or without data augmentations. Moreover, the performances of our model with ResNet-50 backbones are

even better than the previous state-of-the-art with ResNet-101 by a noticeable margin. This demonstrates the effectiveness and generalization ability of the proposed modules.

#### E. Instance Segmentation on NWPU VHR-10 Dataset

To validate the significance of our method, we also conduct instance segmentation experiments on an extra NWPU VHR-10 Dataset, and the performance comparisons are shown in Table IV. The first group includes one-stage methods while the second group contains two-stage methods. Following previous works, we only report the average precisions for mask predictions to focus on the instance segmentation task. The proposed CATNet can still perform better than all the previous approaches with heavier backbones.

#### F. Instance Segmentation on HRSID Dataset

Beyond optical remote sensing images, we also evaluate our model on the more challenging SAR images on HRSID dataset. The results are presented in Table V, in which the first group contains object detection models while the second group includes instance segmentation models. Practically, SAR images are considered as single-channel grayscale images, so that each input for the model is constructed by repeating the SAR images for three times along the channel dimension. The experimental results prove that compared with strong baselines for natural or optical remote sensing images, our method significantly works better on both object detection and instance segmentation.

#### G. Visualizations

We also provide some visualizations to conduct an in-depth study on the significance and effectiveness of our method.

**Context Aggregation Weights:** To demonstrate the effectiveness and significance of SCP, we visualize the context aggregation weights in GCNet [47] and SCP in Fig. 7. Each row represents the weights of all classes that be aggregated into each class. The sizes and color depths of circles denote the weights in GCNet and SCP, respectively. Larger circles or darker colors indicate higher weights. From the visualization,

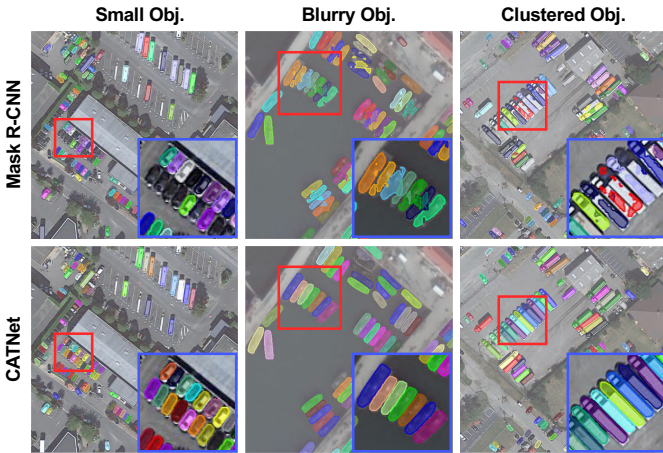


Fig. 10. Detailed comparisons on mask predictions for small, blurry, and clustered objects. The image patches are from iSAID val set.

we can observe that our method tends to aggregate global spatial context from objects with the same category. Some similar (e.g., *plane* and *helicopter*) or semantically related (e.g., *ship* and *harbor*) categories can also help each other during training. For example, SCP tends to aggregate the information of *small vehicle* to *large vehicle*, since they are semantically similar objects, while GCNet simply fuse *plane* into *large vehicle* as global context is often dominated by large objects. Compared with our method, GCNet does not focus much on similar or semantically related objects, leading to information confounding when aggregating features.

**Model Predictions and Attention Maps:** Fig. 8 shows some qualitative results on iSAID, DIOR, and HRSID datasets. Each image patch is visualized by its feature aggregation weights in SCP (top) and the final object detection or instance segmentation results (bottom). The feature aggregation weights indicate that our method focuses more on areas that contain objects rather than pure background, and aggregates them into regions with poor features. Furthermore, we also visualize the detailed comparisons on mask prediction for hard cases in Fig. 10. Our model can better handle small, blurry, and clustered objects, which are common in both optical and other kinds of remote sensing images. These results demonstrate that our method can effectively detect and segment objects accurately in multiple scenes.

**Training Process:** In order to study the model training process, we also visualize the overall training losses in Fig. 9 (a). All the models are trained under the  $1\times$  schedule as stated above, and the training losses of the last 8 epochs are plotted. The solid curves indicate the smoothed losses of different models, while the light-colored areas denote the standard deviations. We start from our baseline (i.e., Mask R-CNN [9]) and gradually incorporate the three modules to study the overall converge processes. The comparison indicates that all the proposed modules, especially DenseFPN, contribute to faster learning and higher performance, as the loss values decline faster and can reach lower minimums. Notably, the full version of CATNet converges about 30% faster than the baseline while reaching a 10% lower final loss.

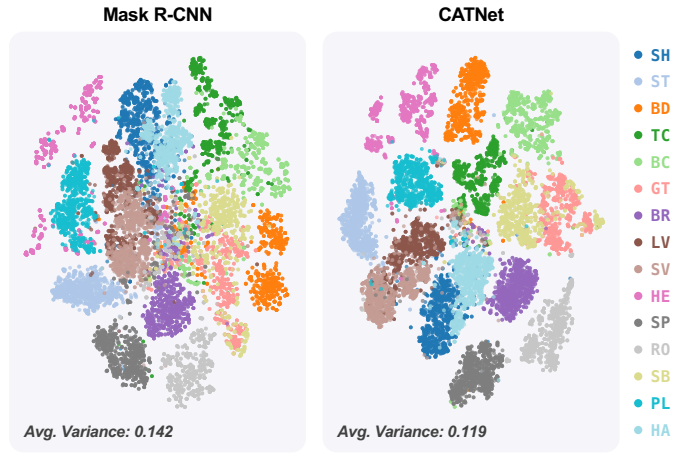


Fig. 11. Visualization of the sampled ROI features from different models using t-SNE. The averaged class-wise variances are also reported.

**Discriminativeness of Features:** To verify that the proposed method can effectively enhance the inadequate features and make them more discriminative, we conduct two studies on the hidden states. First, we visualize the ROI features for some common objects in iSAID dataset in Fig. 1. The feature maps are obtained by randomly select a channel from the ROI features. By aggregating context gradually from different domains, clearer object boundaries can be seen from the feature maps. Second, we randomly sampled 1,000 objects from each category, extracted their average pooled ROI features, and plotted them on a low dimensional manifold using t-SNE [77]. The results are shown in Fig. 11. In our baseline method, features of some categories are heavily confounded, since the top-view and low quality images cannot provide enough information for perception. By introducing multi-scale context aggregation, our method makes the features much more discriminative, as the clusters are more isolated, while features for semantically related categories (e.g., *large vehicle* & *small vehicle* and *helicopter* & *plane*) remain close. We also compared the averaged class-wise variances of features, proving that the enhanced feature clusters are more cohesive, thus more discriminative.

#### H. Detailed Comparisons and Ablation Study

To study the effectiveness of the proposed modules individually, we conduct extensive experiments to compare different settings with some representative methods. All the experiments are performed on iSAID val set using the standard  $1\times$  training recipe with  $512\times 512$  patches inputs. To verify the lightweight characteristics of our method, which means the modules merely introduce negligible extra parameters and computations to the baselines, #Params and FLOPs are reported to compare the space and time complexities of models.

**Plug-and-Play Abilities:** We first verify the plug-and-play abilities of the proposed modules by incorporating them into different instance segmentation models. Aside from Mask R-CNN, we choose Cascade R-CNN [52], Mask Scoring R-CNN [54], and PointRend [11] as base models, and sequentially

incorporate DenseFPN, SCP, and HRoIE into these frameworks. The class-wise instance segmentation performances are reported in Table VI. We observe that for both object detection and instance segmentation tasks, our methods can steadily boost performances. Among the three proposed modules, DenseFPN provides the most significant improvement on mAPs, while SCP and HRoIE also bring considerable gains with few extra parameters. We claim that the proposed modules shall be compatible with any CNN-based instance segmentation pipelines.

**Types of Multi-Scale Feature Propagation Modules:** We then study the accuracies and efficiencies of different multi-scale feature propagation modules. Table VII (a) presents a detailed comparison among different feature pyramid networks. The results are also visualized in Fig. 9 (b) and (c). Compared with existing representative methods, DenseFPN works distinctly better on both object detection and instance segmentation tasks with fewer computational costs. It is worth noting that a single-layer DenseFPN can already obtain considerable gains from the baseline with fewer or comparable computational costs. Simply stacking more basic blocks in DenseFPN can further boost the performances, indicating its capability and flexibility of model scaling.

**Types of Spatial Context Modules:** Table VII (b) shows the comparison among multiple spatial context modules. Compared with the baseline, NLNet [23] can effectively bring higher performance with a high computational cost. GCNet [47] solves the problem of computational complexity, but leads to an extra information confounding problem. With the help of re-weighting context, our proposed CABlock steadily outperforms GCNet using different channel reduction rates. In previous studies, these spatial attention modules are used as plugins of backbones to enhance their global context modeling abilities for image classification. However, we observe that such a strategy is sub-optimal for some modules, due to the large computational consumption with limited performance gains, as can be seen in the table. In this work, we propose to move these modules from the backbone to the layers after FPN, so that the input channels can be largely reduced, while the semantics in the feature maps are more consistent. The comparison of different designs of CABlocks (*i.e.*, used as backbone plugins or incorporated into SCP) shows that the architecture of SCP brings more performance gains with negligible computational costs.

**Order of DenseFPN and SCP:** Since both DenseFPN and SCP are of pyramid structures for multi-scale features, the order of these two modules could be changed. We compare the different orders of these modules in Table VII (c). The first row means using CABlocks as backbone plugins and adopt DenseFPN after backbone. This design is the de-facto standard for most existing works. However, it can only lead to sub-optimal performances. The second row is placing SCP before DenseFPN, so that spatial context is aggregated before feature context. This order still cannot perform well enough. The last row shows our final design, which places SCP after DenseFPN, achieving the best performance on both object detection and instance segmentation. We argue that this is

because the context aggregation process shall be performed in a coarse-to-fine manner, in which we need to 1) calibrate the feature semantics among feature pyramid levels, 2) aggregate spatial context in each level, and 3) refine the instance-specific context for each object. These three steps establish the correct order of DenseFPN, SCP, and HRoIE.

**Types of Region of Interest Extractors:** Table VII (d) shows the comparison between multiple RoI extractors. The baseline model only crops the RoI features from a single feature map, leading to severe information loss and achieving ordinary results. Simply computing the sum or concatenation of RoI features cropped from multiple layers can slightly boost the performance. Considering that object detection and instance segmentation tasks require different features, incorporating HRoIE for adaptive feature fusion can better generate appropriate RoI features for these tasks.

**Ablation Study and Efficiency Comparisons:** We finally conduct ablation studies to validate the effectiveness of each component. As shown in Table VII (e), all the three proposed modules can marginally bring better results on object detection and instance segmentation in remote sensing images. When collaborating with each other, the performance improvements are still stable, indicating that these modules do not interfere with each other. The best results can be achieved by combining all these modules together, enabling aggregating multi-scale context from multiple domains simultaneously. Aside from FLOPs, we also compare the inference speed with and without the proposed modules. The FPS results are computed by averaging the model inference speed on 2,000 images on a single A100 GPU. Compared with the baseline, adding one or two of the proposed modules leads to  $\sim 10\%$  efficiency drop, while the full version is  $\sim 20\%$  slower. Since all the modules are designed to be lightweight, the inference speeds are still acceptable when incorporating them with the baseline. Furthermore, we also observe that the higher latency mainly comes from the number of DenseFPN layers and channels in HRoIE. Therefore, a lite version of CATNet which leverages a 1-layer DenseFPN with 128 channels is proposed. This model has a higher inference speed than the baseline, while maintaining better performances.

## V. CONCLUSION AND FUTURE WORK

In this work, we provided an in-depth analysis on global context modeling in remote sensing images and proposed CATNet, a novel framework that leverages three lightweight plug-and-play modules, *i.e.*, dense feature pyramid network (DenseFPN), spatial context pyramid (SCP), and hierarchical region of interest extractor (HRoIE), to aggregate the global visual context in *feature*, *spatial*, and *instance* domains, respectively. It has been demonstrated that the collaboration among these modules can effectively enhance the discriminative object features for promoting both object detection and instance segmentation performances. We expect that the new understanding of global context and the design of proposed modules will benefit future research in this area. Below we discuss the limitations and future work in this direction.

The motivation of this work is to design lightweight plug-and-play modules for existing instance segmentation pipelines.

Each proposed modules aggregate global context in one domain only. Such a strategy may be sub-optimal since cross-domain context could be correlated with each other. Therefore, a more unified framework shall be designed from scratch to mitigate the limitations of existing instance segmentation models. The recently introduced transformers [67] might be a better solution to handle the global context in multiple domains. We foresee the potential of incorporating multi-scale context aggregation in transformer-based models.

## REFERENCES

- [1] L. Liu, Z. Yang, G. Li, K. Wang, T. Chen, and L. Lin, "Aerial images meet crowdsourced trajectories: A new approach to robust road extraction," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 34, no. 7, pp. 3308–3322, 2022.
- [2] G. Ding, D. Yang, T. Wang, S. Wang, and Y. Zhang, "Crowd counting via unsupervised cross-domain feature adaptation," *IEEE Transactions on Multimedia*, vol. 25, pp. 4665–4678, 2022.
- [3] G. Ding, M. Cui, D. Yang, T. Wang, S. Wang, and Y. Zhang, "Object counting for remote-sensing images via adaptive density map-assisted learning," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–11, 2022.
- [4] L. Liu, B. X. Yu, J. Chang, Q. Tian, and C. W. Chen, "Prompt-matched semantic segmentation," Tech. Rep. arXiv:2208.10159, 2022.
- [5] L. Liu, M. Liu, G. Li, Z. Wu, J. Lin, and L. Lin, "Road network-guided fine-grained urban traffic flow inference," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [6] Y. Luo, Y. Liu, F.-I. Chung, Y. Liu, and C. W. Chen, "End-to-end personalized next location recommendation via contrastive user preference modeling," Tech. Rep. arXiv:2303.12507, 2023.
- [7] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2015, pp. 91–99.
- [8] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2017, pp. 2980–2988.
- [9] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2017, pp. 2961–2969.
- [10] K. Chen, J. Pang, J. Wang, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Shi, W. Ouyang *et al.*, "Hybrid task cascade for instance segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4974–4983.
- [11] A. Kirillov, Y. Wu, K. He, and R. Girshick, "Pointrend: Image segmentation as rendering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [12] X. Wang, T. Kong, C. Shen, Y. Jiang, and L. Li, "Solo: Segmenting objects by locations," Tech. Rep., 2020.
- [13] 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 Journal of Photogrammetry and Remote Sensing*, vol. 98, pp. 119–132, 2014.
- [14] G.-S. Xia, X. Bai, J. Ding, Z. Zhu, S. Belongie, J. Luo, M. Datcu, M. Pelillo, and L. Zhang, "Dota: A large-scale dataset for object detection in aerial images," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 3974–3983.
- [15] S. Waqas Zamir, A. Arora, A. Gupta, S. Khan, G. Sun, F. Shahbaz Khan, F. Zhu, L. Shao, G.-S. Xia, and X. Bai, "isaid: A large-scale dataset for instance segmentation in aerial images," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 28–37.
- [16] K. Li, G. Wan, G. Cheng, L. Meng, and J. Han, "Object detection in optical remote sensing images: A survey and a new benchmark," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 159, pp. 296–307, 2020.
- [17] S. Wei, X. Zeng, Q. Qu, M. Wang, H. Su, and J. Shi, "Hrsid: A high-resolution sar images dataset for ship detection and instance segmentation," *IEEE Access*, vol. 8, pp. 120 234–120 254, 2020.
- [18] 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," Tech. Rep. arXiv:2004.13316, 2020.
- [19] J. Zhang, C. Xie, X. Xu, Z. Shi, and B. Pan, "A contextual bidirectional enhancement method for remote sensing image object detection," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 13, pp. 4518–4531, 2020.
- [20] G. Cheng, Y. Si, H. Hong, X. Yao, and L. Guo, "Cross-scale feature fusion for object detection in optical remote sensing images," *IEEE Geoscience and Remote Sensing Letters*, 2020.
- [21] H. Lin, J. Zhou, Y. Gan, C.-M. Vong, and Q. Liu, "Novel up-scale feature aggregation for object detection in aerial images," *Neurocomputing*, vol. 411, pp. 364–374, 2020.
- [22] X. Zeng, S. Wei, J. Wei, Z. Zhou, J. Shi, X. Zhang, and F. Fan, "Cpynet: Delving into consistent proposals of instance segmentation network for high-resolution aerial images," *Remote Sensing*, vol. 13, no. 14, p. 2788, 2021.
- [23] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7794–7803.
- [24] G. Zhang, S. Lu, and W. Zhang, "Cad-net: A context-aware detection network for objects in remote sensing imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 57, no. 12, pp. 10015–10024, 2019.
- [25] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [26] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 4700–4708.
- [27] D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee, "Yolact: Real-time instance segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 9157–9166.
- [28] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia, "Path aggregation network for instance segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 8759–8768.
- [29] N. He, L. Fang, S. Li, J. Plaza, and A. Plaza, "Skip-connected covariance network for remote sensing scene classification," *IEEE transactions on neural networks and learning systems*, vol. 31, no. 5, pp. 1461–1474, 2019.
- [30] Q. Wang, W. Huang, Z. Xiong, and X. Li, "Looking closer at the scene: Multiscale representation learning for remote sensing image scene classification," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [31] Q. Zhao, S. Lyu, Y. Li, Y. Ma, and L. Chen, "Mgml: Multigranularity multilevel feature ensemble network for remote sensing scene classification," *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- [32] W. Ma, Y. Li, H. Zhu, H. Ma, L. Jiao, J. Shen, and B. Hou, "A multi-scale progressive collaborative attention network for remote sensing fusion classification," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 34, no. 8, pp. 3897–3911, 2021.
- [33] Q. Lin, J. Zhao, G. Fu, and Z. Yuan, "Crpn-sfnet: A high-performance object detector on large-scale remote sensing images," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [34] G. Wang, Y. Zhuang, H. Chen, X. Liu, T. Zhang, L. Li, S. Dong, and Q. Sang, "Fsod-net: Full-scale object detection from optical remote sensing imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–18, 2021.
- [35] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 2117–2125.
- [36] K. Chen, Y. Cao, C. C. Loy, D. Lin, and C. Feichtenhofer, "Feature pyramid grids," Tech. Rep. arXiv:2004.03580, 2020.
- [37] J. Cao, Y. Pang, S. Zhao, and X. Li, "High-level semantic networks for multi-scale object detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 10, pp. 3372–3386, 2019.
- [38] Y. Li, Y. Chen, N. Wang, and Z. Zhang, "Scale-aware trident networks for object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 6054–6063.
- [39] Y. Li, Y. Pang, J. Shen, J. Cao, and L. Shao, "Netnet: Neighbor erasing and transferring network for better single shot object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 13 349–13 358.
- [40] X. Wang, S. Zhang, Z. Yu, L. Feng, and W. Zhang, "Scale-equalizing pyramid convolution for object detection," in *Proceedings of the*

- IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 13 359–13 368.
- [41] S. Qiao, L.-C. Chen, and A. 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 (CVPR)*, 2021, pp. 10 213–10 224.
- [42] G. Ghiasi, T.-Y. Lin, and Q. V. Le, “Nas-fpn: Learning scalable feature pyramid architecture for object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 7036–7045.
- [43] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, “Learning transferable architectures for scalable image recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 8697–8710.
- [44] M. Tan, R. Pang, and Q. V. Le, “Efficientdet: Scalable and efficient object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [45] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7132–7141.
- [46] Z. Huang, X. Wang, L. Huang, C. Huang, Y. Wei, and W. Liu, “Ccnnet: Criss-cross attention for semantic segmentation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 603–612.
- [47] Y. Cao, J. Xu, S. Lin, F. Wei, and H. Hu, “Global context networks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 6, pp. 6881–6895, 2023.
- [48] X. Zhu, D. Cheng, Z. Zhang, S. Lin, and J. Dai, “An empirical study of spatial attention mechanisms in deep networks,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 6688–6697.
- [49] Y. Lee and J. Park, “Centermask: Real-time anchor-free instance segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 13 906–13 915.
- [50] H. Chen, K. Sun, Z. Tian, C. Shen, Y. Huang, and Y. Yan, “Blendmask: Top-down meets bottom-up for instance segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [51] Z. Chen, Y. Shang, A. Python, Y. Cai, and J. Yin, “Db-blendmask: Decomposed attention and balanced blendmask for instance segmentation of high-resolution remote sensing images,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–15, 2021.
- [52] Z. Cai and N. Vasconcelos, “Cascade r-cnn: Delving into high quality object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 6154–6162.
- [53] P. Garg, A. S. Chakravarthy, M. Mandal, P. Narang, V. Chamola, and M. Guizani, “Isdnet: Ai-enabled instance segmentation of aerial scenes for smart cities,” *ACM Transactions on Internet Technology*, vol. 21, no. 3, pp. 1–18, 2021.
- [54] Z. Huang, L. Huang, Y. Gong, C. Huang, and X. Wang, “Mask scoring r-cnn,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 6409–6418.
- [55] T. Vu, H. Kang, and C. D. Yoo, “Schnet: Training inference sample consistency for instance segmentation,” in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2021, pp. 2701–2709.
- [56] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 580–587.
- [57] G. Cheng, P. Zhou, and J. Han, “Learning rotation-invariant convolutional neural networks for object detection in vhr optical remote sensing images,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 54, no. 12, pp. 7405–7415, 2016.
- [58] K. Li, G. Cheng, S. Bu, and X. You, “Rotation-insensitive and context-augmented object detection in remote sensing images,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 4, pp. 2337–2348, 2017.
- [59] G. Cheng, P. Zhou, and J. Han, “Rifd-cnn: Rotation-invariant and fisher discriminative convolutional neural networks for object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 2884–2893.
- [60] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, “Ssd: Single shot multibox detector,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016, pp. 21–37.
- [61] J. Redmon and A. Farhadi, “Yolov3: An incremental improvement,” Tech. Rep. arXiv:1804.02767, 2018.
- [62] S. Wang, Y. Gong, J. Xing, L. Huang, C. Huang, and W. Hu, “Rdsnet: A new deep architecture for reciprocal object detection and instance segmentation,” in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2020.
- [63] Z. Tian, C. Shen, and H. Chen, “Conditional convolutions for instance segmentation,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 282–298.
- [64] H. Su, S. Wei, S. Liu, J. Liang, C. Wang, J. Shi, and X. Zhang, “Hq-isnet: High-quality instance segmentation for remote sensing imagery,” *Remote Sensing*, vol. 12, no. 6, p. 989, 2020.
- [65] X. Zeng, S. Wei, J. Shi, and X. Zhang, “A lightweight adaptive roi extraction network for precise aerial image instance segmentation,” *IEEE Transactions on Instrumentation and Measurement*, vol. 70, pp. 1–17, 2021.
- [66] X. Ke, X. Zhang, and T. Zhang, “Gcbnet: A global context boundary-aware network for sar ship instance segmentation,” *Remote Sensing*, vol. 14, no. 9, p. 2165, 2022.
- [67] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5998–6008.
- [68] L. Rossi, A. Karimi, and A. Prati, “A novel region of interest extraction layer for instance segmentation,” in *Proceedings of the IEEE International Conference on Pattern Recognition (ICPR)*, 2021, pp. 2203–2209.
- [69] J. Wang, K. Sun, T. Cheng, B. Jiang, C. Deng, Y. Zhao, D. Liu, Y. Mu, M. Tan, X. Wang *et al.*, “Deep high-resolution representation learning for visual recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [70] J. Wang, K. Chen, R. Xu, Z. Liu, C. C. Loy, and D. Lin, “Carafe++: Unified content-aware reassembly of features,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [71] J. Pang, K. Chen, J. Shi, H. Feng, W. Ouyang, and D. Lin, “Libra r-cnn: Towards balanced learning for object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 821–830.
- [72] C. Guo, B. Fan, Q. Zhang, S. Xiang, and C. Pan, “Augfpn: Improving multi-scale feature learning for object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 12 595–12 604.
- [73] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2012, pp. 1097–1105.
- [74] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *Proceedings of the International Conference on Machine Learning (ICML)*, 2015.
- [75] N. Bodla, B. Singh, R. Chellappa, and L. S. Davis, “Soft-nms – improving object detection with one line of code,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2017, pp. 5561–5569.
- [76] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, “Automatic differentiation in pytorch,” in *Advances in Neural Information Processing Systems (NeurIPS) Workshops*, 2017.
- [77] L. Van der Maaten and G. Hinton, “Visualizing data using t-sne,” *Journal of Machine Learning Research*, vol. 9, no. 11, pp. 2579–2605, 2008.