# SACANet: scene-aware class attention network for semantic segmentation of remote sensing images

*Xiaowen Ma*[1], *Rui Che*[1], *Tingfeng Hong*[1], *Mengting Ma*[1], *Ziyan Zhao*[1], *Tian Feng*[1,2*] and *Wei Zhang*[1]

[1]Zhejiang University    [2]Alibaba-Zhejiang University Joint Research Institute of Frontier Technologies

*Abstract*—Spatial attention mechanism has been widely used in semantic segmentation of remote sensing images given its capability to model long-range dependencies. Many methods adopting spatial attention mechanism aggregate contextual information using direct relationships between pixels within an image, while ignoring the scene awareness of pixels (i.e., being aware of the global context of the scene where the pixels are located and perceiving their relative positions). Given the observation that scene awareness benefits context modeling with spatial correlations of ground objects, we design a scene-aware attention module based on a refined spatial attention mechanism embedding scene awareness. Besides, we present a local-global class attention mechanism to address the problem that general attention mechanism introduces excessive background noises while hardly considering the large intra-class variance in remote sensing images. In this paper, we integrate both scene-aware and class attentions to propose a scene-aware class attention network (SACANet) for semantic segmentation of remote sensing images. Experimental results on three datasets show that SACANet outperforms other state-of-the-art methods and validate its effectiveness. Code is available at https://github.com/xwmaxwma/rssegmentation.

*Index Terms*—Semantic Segmentation, Scene Awareness, Class Attention

## I. INTRODUCTION

Semantic segmentation aims to predict the semantic class (or label) of each pixel, and is one of the fundamental and challenging tasks in remote sensing image analysis. By providing cues on semantic and localization information for the ground objects of interest, semantic segmentation has been regarded as a vital role in the areas of road extraction [1], urban planning [2], environmental detection [3], etc. In recent years, convolutional neural networks (CNNs) have facilitated the development of semantic segmentation because of their strength in feature extraction. However, it is limited by the local receptive fields and short-range contextual information due to the fixed geometric structure. Consequently, context modeling, including spatial context modeling and relational context modeling, becomes a noticeable option to capture long-range dependencies.

Spatial context modeling methods, such as PSPNet [4] and DeepLabv3+ [5], integrate spatial context information using

spatial pyramid pooling and atrous convolution, respectively. These methods focus on capturing homogeneous context dependencies while often ignoring categorical differences, probably resulting in unreliable contexts if confusing categories are in the scene.

Relational context modeling methods adopt the attention mechanism [6]–[9], which calculates pixel-level similarity in an image for weighted aggregation of heterogeneous contextual information, and have achieved remarkable results in semantic segmentation tasks. However, these methods concentrate on the relationships between pixels while disregarding their awareness of the scene (i.e., global contextual information and position prior), leaving the spatial correlations of ground objects underexplored in remote sensing images.

In this paper, we firstly refine the spatial attention mechanism and propose a scene-aware attention (SAA) module, which contributes effectively to semantic segmentation of remote sensing images. Besides, remote sensing images are characterized by complex backgrounds and large intra-class variance, whereas the conventional attention mechanism overintroduces the background noises due to dense affinity operations and can hardly deal with intra-class variance. In this regard, we introduce the local-global class attention, which associates pixels with the global class representations using the local class representations as intermediate aware elements, achieving efficient and accurate class-level context modeling.

Our contributions are as follows.

- We improve the attention mechanism by embedding the scene awareness of pixels to exploit the spatial correlations of ground objects, which is for the first time to integrate the attention mechanism and scene awareness into a unified module for semantic segmentation of remote sensing images.

- We introduce the local-global class attention mechanism associating pixels with global class representations using local class representations, which achieves efficient and accurate class-level context modeling and tackles complex backgrounds and large intra-class variance in remote sensing images.

- We propose a scene-aware class attention network (SACANet) combining the class attention with scene awareness for semantic segmentation of remote sensing images. Experimental results show that SACANet achieves the state-of-the-art performance on three benchmark datasets, while reaching a decent trade-off between
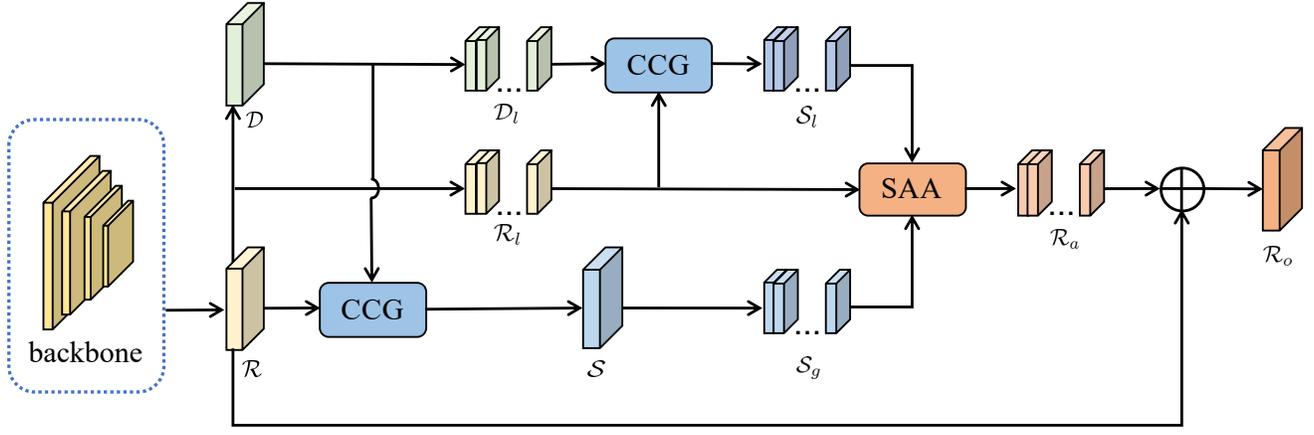
Fig. 1. Architecture of the SACANet. Given the extracted feature representations $\mathcal{R}$ from the backbone and the pre-classified representations $\mathcal{D}$, $\mathcal{R}_l$ and $\mathcal{D}_l$ are obtained after spatial dimensional splitting. Global class center $\mathcal{S}_g$ and local class center $\mathcal{S}_l$ are generated by the class center generation (CCG) module. Then, $\mathcal{R}_l$, $\mathcal{S}_l$, and $\mathcal{S}_g$ are input to the SAA module to obtain the enhanced representations $\mathcal{R}_a$. $\mathcal{R}_a$ is recovered to its original spatial dimension and concatenated with $\mathcal{R}$ to obtain the output representations $\mathcal{R}_o$.

efficiency and accuracy.

## II. METHOD

The proposed SACANet, whose architecture is depicted by Fig. 1, comprises three major components: feature extraction, class center generation and scene-aware attention. To begin with, the HRNetv2-w32 [10] pretrained on ImageNet is deployed as the backbone to extract the feature representations $\mathcal{R}$ from an input image, followed by $\mathcal{R}$ being pre-classified to obtain $\mathcal{D}$. The class center generation (CCG) module then takes as input both $\mathcal{R}$ and $\mathcal{D}$ to achieve the global class center $\mathcal{S}$, which is further cropped to output $\mathcal{S}_g$. Likewise, $\mathcal{R}_l$ and $\mathcal{D}_l$, which are obtained by cropping $\mathcal{R}$ and $\mathcal{D}$, are processed by the CCG module to achieve the local class center $\mathcal{S}_l$. Furthermore, the scene-aware attention (SSA) module takes $\mathcal{R}_l$, $\mathcal{S}_l$, $\mathcal{S}_g$ as inputs to obtain the enhanced feature representations $\mathcal{R}_a$. After original spatial dimensions are recovered, $\mathcal{R}_a$ and $\mathcal{R}$ are concatenated to achieve the output feature representations $\mathcal{R}_o$. The final segmentation map is output after quadruple upsampling.

To illustrate our design of scene-aware attention and local-global class attention, we describe the general form of the attention mechanism as follows: Given feature representations $X^Q, X^K, X^V \in \mathbb{R}^{H \times W \times \hat{C}}$, where $H$, $W$ and $\hat{C}$ denote height, width and dimension of the feature representations respectively. As shown in Fig. 2, the attention mechanism applies three different $1 \times 1$ convolutions $W^Q, W^K, W^V \in \mathbb{R}^{\hat{C} \times C}$ to obtain $q, k, v \in \mathbb{R}^{H \times W \times C}$ as follows,

$$q = X^Q W^Q, k = X^K W^K, v = X^V W^V. \tag{1}$$

Each output element $Z_i$ is a weighted sum of input elements $\{v_j\}$ as follows,

$$Z_i = \sum_{j=1}^{H \times W} \alpha_{ij} v_j, \tag{2}$$

where $\alpha_{ij}$ denotes the weight from the softmax function on $e_{ij}$, which is obtained by a scaled dot-product attention as follows,

$$e_{ij} = \frac{q_i k_j^T}{\sqrt{C}}. \tag{3}$$

### A. Scene-Aware Attention

Ground objects in remote sensing images have intrinsic spatial correlations that are frequently observable. For example, vehicles are usually found to stay on a road; buildings are densely distributed on both sides of a road. Therefore, it is supposed to facilitate the modeling of corresponding patterns by embedding the scene awareness of pixels (i.e., considering the global context of pixels as well as their relative position in the attention).

**Contextual Information Embedding.** In remote sensing images, pair-wise relationships between ground objects may vary in different scenes. For example, a road that usually coexists with buildings in an urban area may be surrounded by cropland in a rural area. It suggests that embedding contextual information can benefit the modeling of pixel-level relationships. Inspired by [11], we propose a context matrix and reformulate the previous equation as follows,

$$e_{ij} = \frac{(q_i c) k_j^T}{\sqrt{C}}, \tag{4}$$

and the context matrix $c$ is computed as follows,

$$c = diag(\sigma(W_1(W_0(AvgPool(Q))) + W_1(W_0(MaxPool(Q))))) \tag{5}$$

where $\sigma$ denotes the sigmoid function, $W_0 \in \mathbb{R}^{(C/\epsilon) \times C}$, $W_1 \in \mathbb{R}^{C \times (C/\epsilon)}$, and $\epsilon$ is the reduction ratio. We create a diagonal matrix for the vector to connect the summarized contextual information with the input features. Finally, the context matrix $c$ contextualizes the input features $q$ so that the attention can be adjusted by the given context.
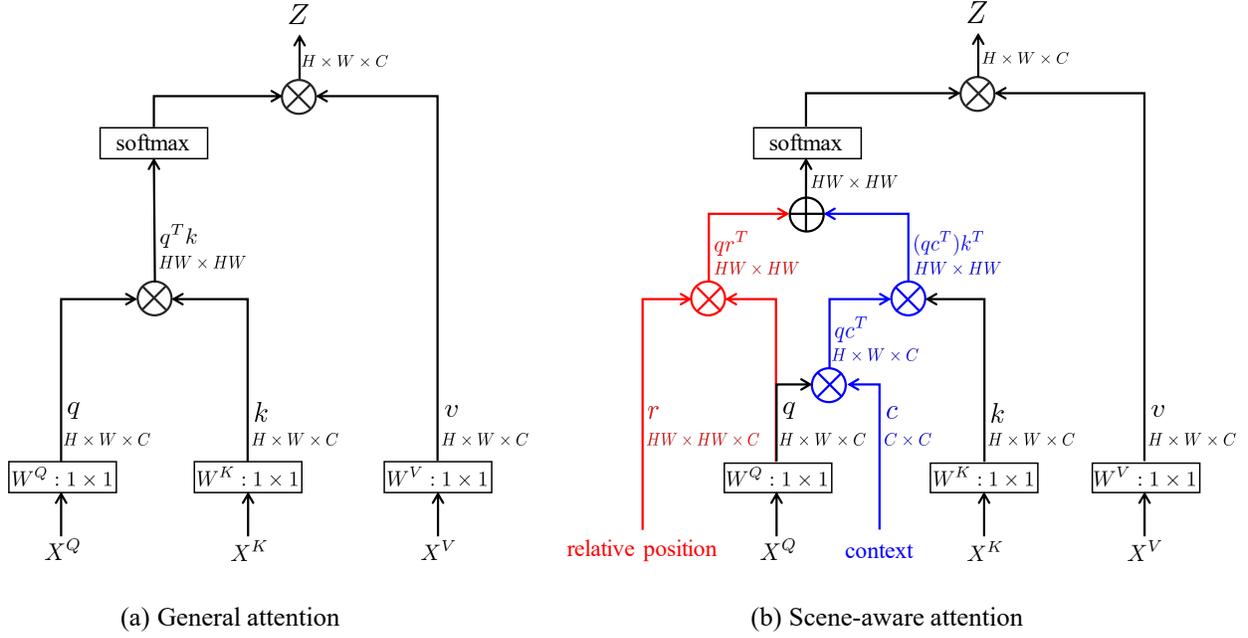
(a) General attention

(b) Scene-aware attention

Fig. 2. Illustration of (a) the general attention (GA) and (b) the scene-aware attention (SAA). The red and blue parts in the SAA are newly added compared to the GA, representing the relative position and context to embed the pixels' scene awareness in the attention.

**Position Prior Embedding.** In remote sensing images, ground objects are spatially distributed following specific intrinsic patterns. In particular, certain combinations or concurrences usually occur to the objects in close proximity, and the pixels in an object's vicinity may demonstrate high correlation. These observations suggest that a pixel's awareness of the scene relies on its sensitivity to relative positions, which are embedded as,

$$e_{ij} = \frac{(q_i c) k_j^T + q_i r_{ij}^T}{\sqrt{C}}. \tag{6}$$

Unlike previous work in the field of natural language processing [12], we extend the words in a one-dimensional sequence to the pixels in a two-dimensional plane, considering their relative positions as a combined effect along both horizontal and vertical directions. Specifically, the encoding of relative position $r_{ij}$ is defined as follows,

$$r_{ij} = P_{I^x(i,j), I^y(i,j)}, \tag{7}$$

where $P \in \mathbb{R}^{(2\xi+1) \times (2\xi+1) \times C}$ is a bucket storing a set of indexed trainable vectors, $I^x(i,j) = g(x_i - x_j)$ and $I^y(i,j) = g(y_i - y_j)$ represent subscripts for horizontal and vertical directions, forming two-dimensional indices for $P$, and $g$ denotes an index function as follows,

$$g(x) = max(-\xi, min(x, \xi)), \tag{8}$$

where $\xi$ refers to the maximum pixel-level distance. Actually, $g(x)$ maps the distance to an integer in finite set, largely reducing the number of parameters and computation cost needed for high resolution remote sensing images.
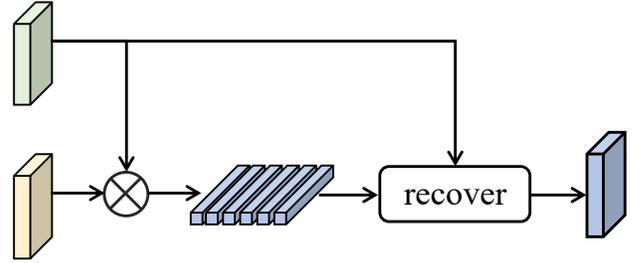


Fig. 3. Architecture of the CCG.

### B. Local-Global Class Attention

Self-attention mechanism [6], [17], [22] has so far dominated the attention-based semantic segmentation methods, whose inputs $X^Q$, $X^K$ and $X^V$ to the attention module are all set to the given feature representations $\mathcal{R} \in \mathbb{R}^{H \times W \times \hat{C}}$.

Considering that remote sensing images are characterized by complex backgrounds and large intra-class variance, these methods can result in massive background noise towards poor performances due to dense affinity operations. Several class attention methods attempt to resolve this problem using global class center for class-level context modeling; However, they are yet to consider intra-class variance and the case that pixels may be semantically distant from the global class center impacting on the class context modeling. Therefore, we present the local-global class attention to improve the performance of class-level context modeling. In particular, pixels are indirectly associated with global class representations by introducing local class representations.

As shown in Fig. 3, for the feature representations $\mathcal{R} \in$

| Method | LoveDA | | | | | | | | Vaihingen | | | Potsdam | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Back. | Buil. | Road | Water | Barren | Forest | Agri. | mIoU | AF | mIoU | OA | AF | mIoU | OA |
| PSPNet [4] | 44.4 | 52.1 | 53.5 | 76.5 | 9.7 | 44.1 | 57.9 | 48.3 | 86.47 | 76.78 | 89.36 | 89.98 | 81.99 | 90.14 |
| DeepLabv3+ [5] | 43.0 | 50.9 | 52.0 | 74.4 | 10.4 | 44.2 | 58.5 | 47.6 | 86.77 | 77.13 | 89.12 | 90.86 | 84.24 | 89.18 |
| DANet [6] | 44.8 | 55.5 | 53.0 | 75.5 | 17.6 | 45.1 | 60.1 | 50.2 | 86.88 | 77.32 | 89.47 | 89.60 | 81.40 | 89.73 |
| Semantic FPN [13] | 42.9 | 51.5 | 53.4 | 74.7 | 11.2 | 44.6 | 58.7 | 48.2 | 87.58 | 77.94 | 89.86 | 91.53 | 84.57 | 90.16 |
| FarSeg [14] | 43.1 | 51.5 | 53.9 | 76.6 | 9.8 | 43.3 | 58.9 | 48.2 | 87.88 | 79.14 | 89.57 | 91.21 | 84.36 | 89.87 |
| OCRNet [7] | 44.2 | 55.1 | 53.5 | 74.3 | 18.5 | 43.0 | 60.5 | 49.9 | 89.22 | 81.71 | 90.47 | 92.25 | 86.14 | 90.03 |
| LANet [15] | 40.0 | 50.6 | 51.1 | 78.0 | 13.0 | 43.2 | 56.9 | 47.6 | 88.09 | 79.28 | 89.83 | 91.95 | 85.15 | 90.84 |
| ISNet [8] | 44.4 | 57.4 | 58.0 | 77.5 | **21.8** | 43.9 | 60.6 | 51.9 | 90.19 | 82.36 | 90.52 | 92.67 | 86.58 | 91.27 |
| Segmenter [16] | 38.0 | 50.7 | 48.7 | 77.4 | 13.3 | 43.5 | 58.2 | 47.1 | 88.23 | 79.44 | 89.93 | 92.27 | 86.48 | 91.04 |
| SwinUperNet [17] | 43.3 | 54.3 | 54.3 | 78.7 | 14.9 | 45.3 | 59.6 | 50.0 | 89.9 | 81.8 | 91.0 | 92.24 | 86.37 | 90.98 |
| MANet [18] | 38.7 | 51.7 | 42.6 | 72.0 | 15.3 | 42.1 | 57.7 | 45.7 | 90.41 | 82.71 | 90.96 | 92.90 | 86.95 | 91.32 |
| FLANet [19] | 44.6 | 51.8 | 53.0 | 74.1 | 15.8 | 45.8 | 57.6 | 49.0 | 87.44 | 78.08 | 89.60 | 93.12 | 87.50 | 91.87 |
| ConvNeXt [20] | 46.9 | 53.5 | 56.8 | 76.1 | 15.9 | **47.5** | 61.8 | 51.2 | 90.50 | 82.87 | 91.36 | 93.03 | 87.17 | 91.66 |
| PoolFormer [21] | 45.8 | 57.1 | 53.3 | **80.7** | 19.8 | 45.6 | 64.5 | 52.4 | 89.59 | 81.35 | 90.30 | 92.62 | 86.45 | 91.12 |
| SACANet (Ours) | **47.6** | **59.1** | **58.4** | 80.5 | 17.8 | 46.7 | **67.1** | **53.9** | **91.68** | **84.49** | **92.10** | **93.64** | **87.89** | **92.28** |

$\mathbb{R}^{\hat{C} \times H \times W}$, a pre-classification is deployed (i.e., two consecutive $1 \times 1$ convolution layers) to obtain the corresponding distribution $\mathcal{D} \in \mathbb{R}^{K \times H \times W}$, where $K$ is the number of classes. The global class center $S$ is defined as follows,

$$\mathcal{S} = \psi(\mathcal{D}^{K \times (H \times W)} \otimes \mathcal{R}^{(H \times W) \times \hat{C}}), \qquad (9)$$

where $\mathcal{S}$ denotes a $H \times W \times \hat{C}$ matrix, $\psi$ represents a function to place class centers in the original feature map according to the pre-classification generated mask. Then, $\mathcal{R}$ and $\mathcal{D}$ are split along the spatial dimension to reach $\mathcal{R}_l$ and $\mathcal{D}_l$, followed by calculating the local class representations $\mathcal{S}_l$ as follows,

$$\mathcal{S}_l = \psi(\mathcal{D}_l^{(N_h \times N_w) \times K \times (h \times w)} \otimes \mathcal{R}_l^{(N_h \times N_w) \times (h \times w) \times C}), \quad (10)$$

where $h$ and $w$ are the height and width of the selected local patch, $N_h = \frac{H}{h}$, and $N_w = \frac{W}{w}$. Similarly, $\mathcal{S}$ is split along the spatial dimension to obtain $\mathcal{S}_g \in \mathbb{R}^{(N_h \times N_w) \times (h \times w) \times \hat{C}}$. Hence, the inputs to the attention module are as follows:

$$X^Q = \mathcal{R}_l, \quad X^K = \mathcal{S}_l, \quad X^V = \mathcal{S}_g. \qquad (11)$$

The design of local-global class attention combines scene awareness with local-global class attention. In addition, the slicing operation provides a noticeable decrease in the number of parameters and computation, enabling the lightweight of the method.

## III. EXPERIMENTAL RESULTS

### A. Datasets and Evaluation Metrics

We conduct the experiments on three publicly available datasets to evaluate our SACANet in three common metrics: average F1-score (AF), mean Intersection-over-Union (mIoU), and overall accuracy (OA).

LoveDA [23] contains 5987 fine-resolution optical remote sensing images (GSD 0.3 m) at a size of 1024 × 1024 pixels and includes 7 landcover categories, i.e., building, road, water,

| Method | Params (M) | FLOPs (G) | Memory (MB) |
|---|---|---|---|
| PPM [4] | 23.1 | 309.5 | 257 |
| ASPP [5] | 15.1 | 503.0 | 284 |
| DAB [6] | 23.9 | 392.2 | 1546 |
| OCR [7] | 10.5 | 354.0 | 202 |
| PAM+AEM [15] | 10.4 | 157.6 | 489 |
| ILCM+SLCM [8] | 11.0 | 180.6 | 638 |
| FLA [19] | 11.5 | 154.9 | 645 |
| CCG+SAA (Ours) | 2.7 | 44.4 | 76 |

barren, forest, agriculture and background. Specifically, we use 2522 images for training, 1669 images for validation and the remaining 1796 images for testing.

ISPRS Vaihingen [24] contains 33 TOP tiles and DSMs (GSD 9 cm) collected from a small village and includes 6 landcover categories, i.e., impervious surfaces, building, low vegetation, tree, car, and clutter/background. The size of the images varies from 1996 × 1995 pixels to 3816 × 2550 pixels. We use 16 images for training and the remaining 17 for testing.

ISPRS Potsdam [24] consists of 38 TOP tiles and DSMs (GSD 5 cm) collected from a historic city at a size of 6000 × 6000 pixels and includes the same six categories as the Vaihingen dataset. We use 24 images for training and the remaining 14 for testing.

### B. Implementation Details

For all experiments, the optimizer is SGD with the batch size of 4, and the initial learning rate is set to 0.01 with a poly decay strategy and a weight decay of 0.0001. Following previous work [15], [18], we randomly crop the images from three datasets to produce 512 × 512 patches, and the augmentation methods, such as random scale ([0.5, 0.75, 1.0, 1.25, 1.5]), random vertical flip, random horizontal flip and random rotate, are adopted in the training process. The number

| Model | AF | mIoU | OA |
|---|---|---|---|
| Base | 89.65 | 81.87 | 90.13 |
| Base+GA | 89.34 | 81.62 | 89.94 |
| Base+SAA | 90.84 | 83.21 | 91.57 |
| Base+CCG(g+g)+GA | 89.60 | 82.03 | 90.13 |
| Base+CCG(l+g)+GA | 90.48 | 82.76 | 91.23 |
| Base+CCG(g+g)+SAA | 90.67 | 83.47 | 91.12 |
| SACANet (Ours) | **91.68** | **84.49** | **92.10** |

of epochs on LoveDA, ISPRS Vaihingen, and ISPRS Potsdam is set to 30, 150 and 80, respectively.

### C. Comparison and Analysis

As shown in Table I, the proposed method outperforms other state-of-the-art methods. Specifically, our SACANet achieves an improvement of 1.5% in mIoU compared to PoolFormer [21] on LoveDA; On ISPRS Vaihingen and Potsdam, SACANet's improvement is about 3.1% and 1.4%, respectively. In particular, significant improvements are found to the background class, which is complex and has a large intra-class variance, as well as the road and agriculture classes, which are closely related to the corresponding scene, on LoveDA. These improvements have validated the effectiveness of SACANet on semantic segmentation of remote sensing images via embedding scene-aware and local-global attentions. Fig. 4 shows example results from PSPNet, DANet, MANet and our SACANet. The proposed method not only better preserves the integrity and regularity of semantic objects such as building and low vegetation, but also improves the segmentation performance of small objects, such as car.

In addition, we compare several context aggregation modules and ours in three metrics: number of parameters (Params) measured in million (M), number of floating-point operations per second (FLOPs) measured in giga (G), and the memory consumption (Memory) measured in megabytes (MB). As shown in Table II, the attention modules in our method are based on local patches, greatly reducing the required number of parameters and computation. Specifically, we require only 26% of Params, 13% of the FLOPs, and 38% of the Memory compared to the light OCR [7] module, which significantly improves the efficiency of the model.

Furthermore, we conduct an ablation study using HRNetv2-w32 as the base on ISPRS Vaihingen to investigate the impacts of SAA and CCG modules. As shown in Table 3, we observe that the semantic segmentation performance of the base is slightly degraded after adding the general attention (GA) module, but improved by introducing our SAA module. Besides, our SACANet integrating both SAA and CCG modules achieves an even higher performance. The results of the ablation study have supported that introducing scene-aware and class attentions can benefit semantic segmentation of remote sensing images.

### IV. CONCLUSION

In this paper, we present scene-aware attention that refines the spatial attention mechanism to exploit the inherent spatial correlation of ground objects in remote sensing images. Considering the complex backgrounds and large intra-class variances, we introduce local-global class attention for class-wise context modeling, which prevents dense attention from over-introducing the interference of background noises. Integrating both scene-aware and local-global class attentions, we propose SACANet that significantly improves the semantic segmentation performance on remote sensing images. Experimental results reveal our SACANet's outperformance compared to other state-of-the-art methods. Besides, the proposed method reduces the number of parameters and computation significantly, and achieves a better trade-off between accuracy and efficiency.

REFERENCES

[1] M. Maboudi, J. Amini, S. Malihi, and M. Hahn, "Integrating fuzzy object based image analysis and ant colony optimization for road extraction from remotely sensed images," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 138, pp. 151–163, 2018.

[2] Q. Zhang and K. C. Seto, "Mapping urbanization dynamics at regional and global scales using multi-temporal dmsp/ols nighttime light data," *Remote Sensing of Environment*, vol. 115, no. 9, pp. 2320–2329, 2011.

[3] Q. Yuan, H. Shen, T. Li, Z. Li, S. Li, Y. Jiang, H. Xu, W. Tan, Q. Yang, J. Wang *et al.*, "Deep learning in environmental remote sensing: Achievements and challenges," *Remote Sensing of Environment*, vol. 241, p. 111716, 2020.

[4] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2881–2890.

[5] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 801–818.

[6] J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, and H. Lu, "Dual attention network for scene segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 3146–3154.

[7] Y. Yuan, X. Chen, and J. Wang, "Object-contextual representations for semantic segmentation," in *Proceedings of the European conference on computer vision (ECCV)*. Springer, 2020, pp. 173–190.

[8] Z. Jin, B. Liu, Q. Chu, and N. Yu, "Isnet: Integrate image-level and semantic-level context for semantic segmentation," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 7189–7198.

[9] X. Ma, M. Ma, C. Hu, Z. Song, Z. Zhao, T. Feng, and W. Zhang, "Log-can: local-global class-aware network for semantic segmentation of remote sensing images," *arXiv preprint arXiv:2303.07747*, 2023.

[10] 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*, vol. 43, no. 10, pp. 3349–3364, 2020.

[11] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "Cbam: Convolutional block attention module," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 3–19.

[12] P. Shaw, J. Uszkoreit, and A. Vaswani, "Self-attention with relative position representations," *arXiv preprint arXiv:1803.02155*, 2018.

[13] A. Kirillov, R. Girshick, K. He, and P. Dollár, "Panoptic feature pyramid networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 6399–6408.

[14] Z. Zheng, Y. Zhong, J. Wang, and A. Ma, "Foreground-aware relation network for geospatial object segmentation in high spatial resolution remote sensing imagery," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2020, pp. 4096–4105.
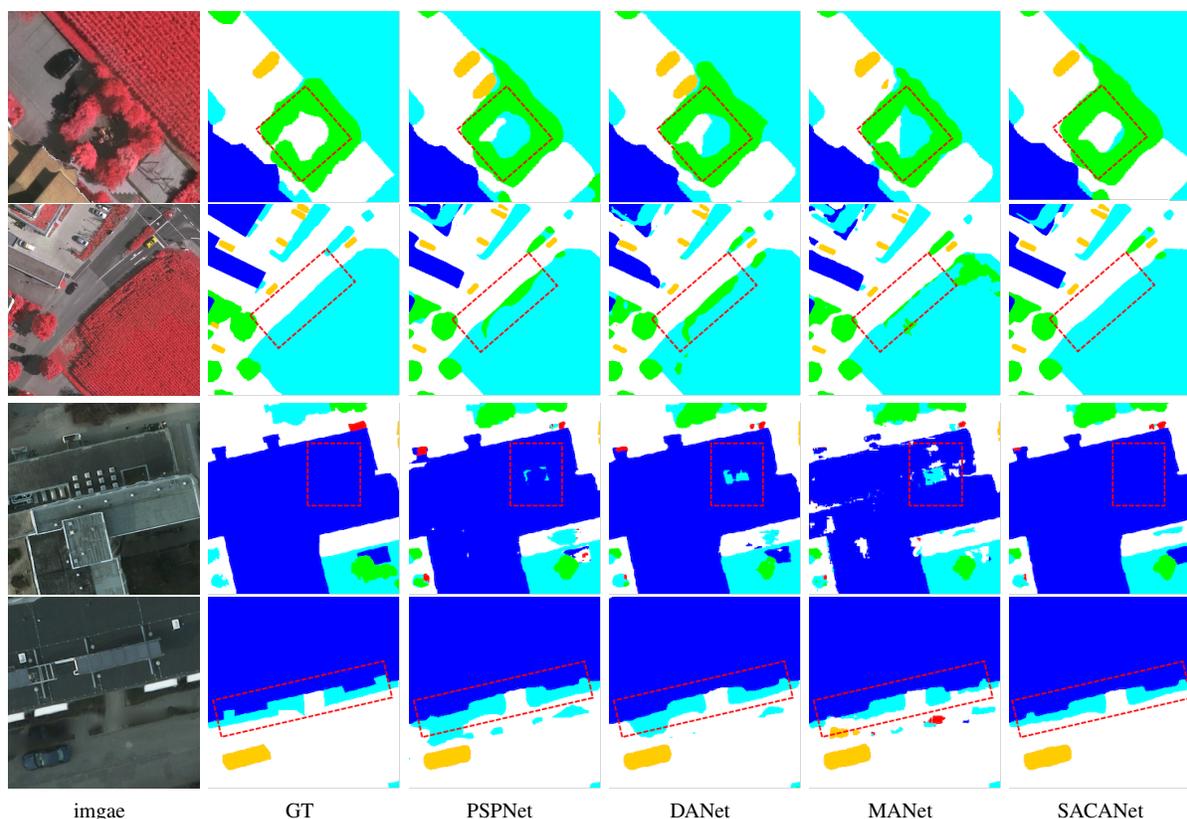
Fig. 4. Example outputs from the SACANet and other methods on the ISPRS Vaihingen test set and ISPRS Potsdam test set. Best viewed in color and zoom in.

[15] L. Ding, H. Tang, and L. Bruzzone, "Lanet: Local attention embedding to improve the semantic segmentation of remote sensing images," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 1, pp. 426–435, 2021.

[16] R. Strudel, R. Garcia, I. Laptev, and C. Schmid, "Segmenter: Transformer for semantic segmentation," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 7262–7272.

[17] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 10 012–10 022.

[18] R. Li, S. Zheng, C. Zhang, C. Duan, J. Su, L. Wang, and P. M. Atkinson, "Multiattention network for semantic segmentation of fine-resolution remote sensing images," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–13, 2022.

[19] Q. Song, J. Li, C. Li, H. Guo, and R. Huang, "Fully attentional network for semantic segmentation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 2, 2022, pp. 2280–2288.

[20] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A convnet for the 2020s," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11 976–11 986.

[21] W. Yu, M. Luo, P. Zhou, C. Si, Y. Zhou, X. Wang, J. Feng, and S. Yan, "Metaformer is actually what you need for vision," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 819–10 829.

[22] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7794–7803.

[23] J. Wang, Z. Zheng, A. Ma, X. Lu, and Y. Zhong, "Loveda: A remote sensing land-cover dataset for domain adaptive semantic segmentation," *arXiv preprint arXiv:2110.08733*, 2021.

[24] F. Rottensteiner, G. Sohn, J. Jung, M. Gerke, C. Baillard, S. Bnitez, and U. Breitkopf, "International society for photogrammetry and remote sensing, 2d semantic labeling contest," Accessed: Oct. 29, 2020., available: https://www.isprs.org/education/benchmarks.