

# SCORE: Scene Context Matters in Open-Vocabulary Remote Sensing Instance Segmentation

Shiqi Huang<sup>1</sup> Shuting He<sup>2</sup> Huaiyuan Qin<sup>3</sup> Bihan Wen<sup>1\*</sup>

<sup>1</sup>Nanyang Technological University

<sup>2</sup>MoE Key Laboratory of Interdisciplinary Research of Computation and Economics,  
Shanghai University of Finance and Economics

<sup>3</sup>Institute for Infocomm Research (I<sup>2</sup>R), A\*STAR, Singapore

## Abstract

Most existing remote sensing instance segmentation approaches are designed for close-vocabulary prediction, limiting their ability to recognize novel categories or generalize across datasets. This restricts their applicability in diverse Earth observation scenarios. To address this, we introduce open-vocabulary (OV) learning for remote sensing instance segmentation. While current OV segmentation models perform well on natural image datasets, their direct application to remote sensing faces challenges such as diverse landscapes, seasonal variations, and the presence of small or ambiguous objects in aerial imagery. To overcome these challenges, we propose **SCORE** (Scene Context matters in Open-vocabulary REmote sensing instance segmentation), a framework that integrates multi-granularity scene context, i.e., regional context and global context, to enhance both visual and textual representations. Specifically, we introduce Region-Aware Integration, which refines class embeddings with regional context to improve object distinguishability. Additionally, we propose Global Context Adaptation, which enriches naive text embeddings with remote sensing global context, creating a more adaptable and expressive linguistic latent space for the classifier. We establish new benchmarks for OV remote sensing instance segmentation across diverse datasets. Experimental results demonstrate that, our proposed method achieves SOTA performance, which provides a robust solution for large-scale, real-world geospatial analysis. Our code is available at <https://github.com/HuangShiqi128/SCORE>.

## 1. Introduction

Instance segmentation is a fundamental task in remote sensing, aiming to localize aerial objects with pixel-wise instance masks and category labels [2, 3, 33, 44, 45, 57, 62,

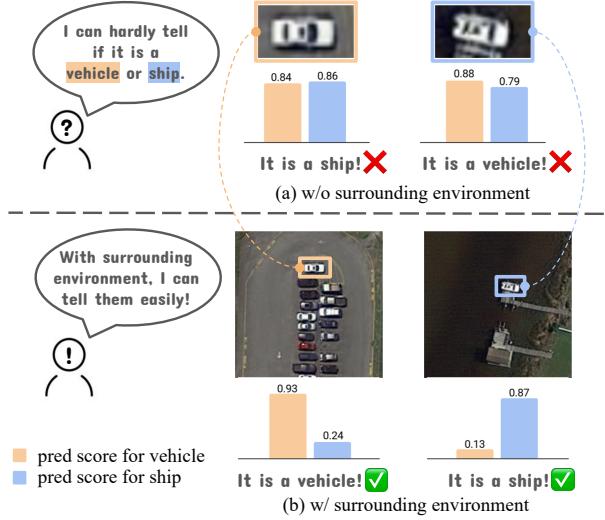


Figure 1. The example illustrates how human perception & model prediction differ (a) without and (b) with the surrounding environment. Prediction scores are derived from FC-CLIP [63].

66]. It plays a vital role in diverse applications, including environmental monitoring [16, 43], urban development [32, 51, 53, 59], and agricultural planning [69, 70]. With the vast amount of remote sensing data collected from satellites, drones, and aerial surveys, instance segmentation serves as a key tool for large-scale geospatial analysis.

Most existing approaches rely on pixel-wise annotated datasets to train instance segmentation models [44, 45, 64]. However, manually labeling aerial imagery is both time-consuming and challenging due to the small, densely packed, and visually ambiguous nature of instance objects in remote sensing [17, 58]. Furthermore, these models are constrained by the training data, limiting their ability to recognize novel categories or generalize across different domains. This restricts their applicability in diverse Earth observation tasks. Open-vocabulary (OV) learning [52, 55,

\*Corresponding author

[65] offers a promising solution by enabling models to predict novel classes without requiring exhaustive annotations or retraining. This capability is particularly valuable for remote sensing, where rapid adaptation is essential for monitoring evolving environmental conditions.

Although numerous OV segmentation models have been developed for natural images [12, 22, 54, 63], directly applying the models to remote sensing is less effective. This is due to the unique challenges posed by diverse landscapes, seasonal variations, and different imaging conditions that must be accounted for aerial and satellite imagery [1, 61]. Additionally, aerial images, typically captured from a bird’s-eye view at large scales, often contain small or ambiguous objects that are difficult to recognize [17, 58]. For example, as shown in Figure 1(a), differentiating between vehicles and ships can be challenging due to their shared elongated shapes and similar appearances.

On the other hand, in remote sensing, objects are often closely correlated with their surrounding environment, providing useful prior knowledge [27, 28, 36]. For instance, ships are typically found near coastal areas, cars are associated with roads or parking lots, airplanes are near airports, and agricultural fields are common in rural areas. This regional scene context can help differentiate instance-level objects. As illustrated in Figure 1(b), given the presence of water and harbors, the object on the left is likely a ship, whereas the object on the right, situated in a parking lot, is a car. Motivated by this, we propose leveraging **regional context**, derived from the surrounding environment, to enhance object representations. With the rise of remote sensing-specific vision-language models, *i.e.*, remote sensing CLIPs [30, 38, 50, 68], these models encode valuable domain knowledge, enabling the extraction of meaningful regional scene context. By utilizing regional context as a guiding cue, our approach aims to improve object recognition in remote sensing imagery, particularly in the challenging open-vocabulary settings.

While incorporating prior knowledge enhances recognition by associating instances with their regional context, existing OV segmentation models, which rely on frozen classifiers trained on natural images, lack domain-specific adaptability for remote sensing. Their fixed text embeddings struggle to capture the significant intra-class variance, resolution differences, and environmental diversity present in remote sensing datasets [9, 18, 67]. To address this challenge, we propose enhancing the text embeddings by incorporating domain-specific **global context**, a subset of scene context extracted from remote sensing CLIP. This adaptation enriches the classifier’s linguistic latent space with domain-relevant visual cues, ultimately improving class predictions for aerial objects.

Our main contributions can be summarized as follows:

- We put forward open-vocabulary remote sensing instance

segmentation task and develop a comprehensive framework **SCORE** to improve OV instance segmentation performance across diverse remote sensing datasets. We establish new experimental benchmarks and achieve SOTA performance on the proposed task.

- We design Region-Aware Integration, leveraging the correlation between objects and their surrounding environment in aerial images. By incorporating regional context retrieved from domain-specific CLIP, our method enhances class embeddings, improving object classification in open-vocabulary remote sensing instance segmentation.
- To bridge the gap between general-domain and remote sensing-specific classifiers, we introduce Global Context Adaptation which injects aerial global context into CLIP text embeddings. This adaptation enhances richness of text representation, making them more expressive and suited for classifying remote sensing objects.

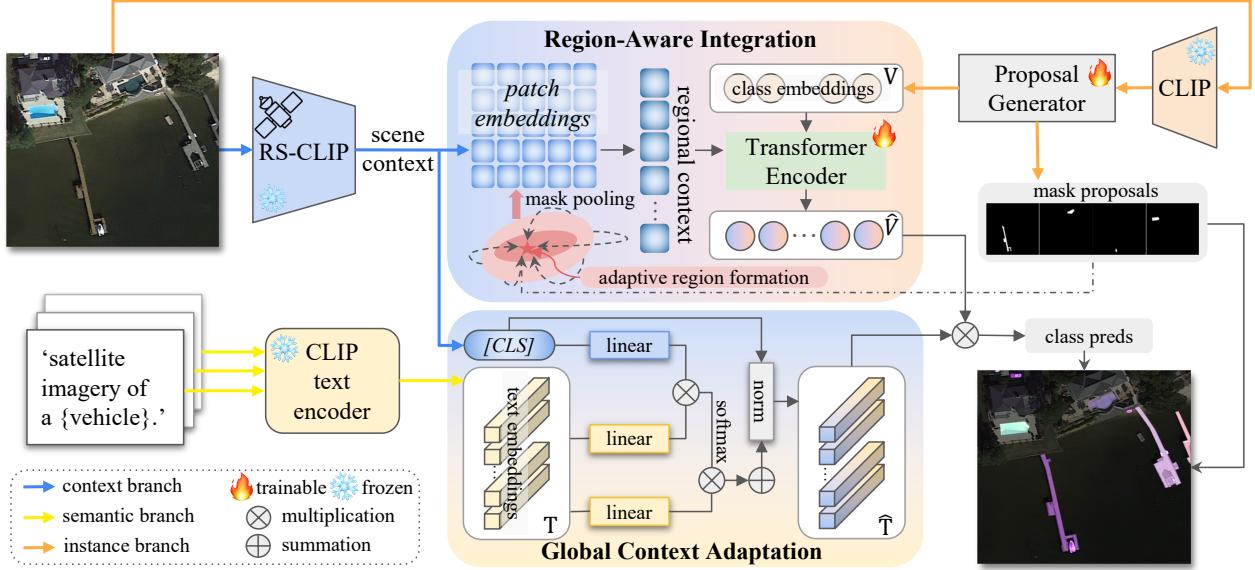
## 2. Related Work

### 2.1. Remote Sensing Instance Segmentation

With the advent of deep learning, segmentation frameworks for remote sensing images have made notable progress. For remote sensing instance segmentation, methods have been developed to tackle challenges such as scale variations, foreground-background confusion, and class ambiguity [17, 31, 33, 44–46, 62, 66]. More recently, the integration of foundation models into remote sensing segmentation has been explored. For example, RSPromoter [4] enhances the instance segmentation capabilities of SAM [23] on remote sensing images using prompt learning. SAMRS [48] expands remote sensing segmentation datasets by leveraging SAM. ZoRI [18] introduces zero-shot learning for remote sensing instance segmentation. Despite these advancements, the ability to make cross-dataset prediction, *i.e.*, open vocabulary learning, for remote sensing instance segmentation has not been studied.

### 2.2. Open-Vocabulary Segmentation

With the emergence of large-scale pre-trained multimodal foundation models like CLIP [40], numerous methods have been developed for OV segmentation in natural images by leveraging their strong generalization capabilities [5, 10, 12, 14, 22, 24, 29, 39, 54–56, 63]. MaskCLIP [12] adapts CLIP visual encoder with relative mask attention for universal OV segmentation. ODISE [54] leverages the semantically differentiated representation pretrained in text-image diffusion models to generalize to novel classes. FreeSeg [39] builds a unified model that captures task-aware and category-sensitive concepts for more robust segmentation. OPSNet [5] dynamically modulates query embeddings with CLIP features to enhance adaptation to novel categories. FC-CLIP [63] incorporates a frozen CNN-



**Figure 2. Overview of SCORE training framework.** It consists of three main branches: context branch (blue) extracts multi-granularity scene context from RS CLIP, semantic branch (yellow) encodes text embeddings with a frozen CLIP text encoder, and instance branch (orange) generates class embeddings and instance proposals. These branches interact through Region-Aware Integration (RAI) and Global Context Adaptation (GCA) to derive scene-context-enhanced class embeddings  $\hat{V}$  and text embeddings  $\hat{T}$ , which are used for class prediction.  $[CLS]$  represents the domain-specific [CLS] token, *i.e.*, global context.

based CLIP image encoder as the backbone, preserving cross-modal alignment while avoiding redundant feature extraction. MAFT+ [22] further optimizes CLIP’s vision-text representations to improve OV segmentation results.

While significant progress has been made in OV segmentation for natural images, there is also a growing interest in extending it to the field of remote sensing. Recent works have explored OV remote sensing semantic segmentation [1, 61]. OVRS [1] introduces orientation-adaptive semantics to address the varying orientations of aerial objects. GSNet [61] fuses features from CLIP backbone with remote sensing backbone to derive domain-specific representation. However, existing methods are limited to semantic segmentation, there remains a critical need for an OV instance segmentation framework to address the challenges in remote sensing.

### 2.3. Vision-Language Models

Vision-language models (VLMs) are designed to learn a shared representation space that bridges visual and textual modalities through contrastive learning. By pretraining on large-scale web-based datasets [21, 40], these models align image and text embeddings to capture meaningful semantic relationships. This cross-modal alignment enables VLMs to generalize effectively across a wide range of downstream tasks, including image retrieval, dense prediction, and cross-modal reasoning [15, 34, 37, 42, 49].

Inspired by the success of VLMs in natural image applications, researchers have adapted these models to the remote sensing domain. Recent efforts [30, 38, 50, 68]

have explored applying vision-language contrastive learning for aerial and satellite imagery, adapting general VLMs to remote sensing context with the constructed specialized cross-modal datasets. Therefore, these models encode domain-specific knowledge, enhancing their ability to extract and utilize remote sensing features effectively.

## 3. Method

### 3.1. Task Formation

Open-vocabulary remote sensing instance segmentation aims to segment an image  $I \in \mathbb{R}^{H \times W \times 3}$  into a set of pixel-wise instance masks  $m_i$  with associated aerial categories  $c_i$ :

$$\{y_i\}_{i=1}^K = \{(m_i, c_i)\}_{i=1}^K. \quad (1)$$

During training, the model learns from a predefined set of categories  $C_{\text{train}}$ , while during inference, it is evaluated on a separate set of categories  $C_{\text{test}}$ , where  $C_{\text{train}} \neq C_{\text{test}}$ . The category names of  $C_{\text{test}}$  is available during testing.

### 3.2. Framework Overview

The overall framework of our proposed **SCORE** is illustrated in Figure 2. The pipeline consists of three main parts: the context branch, semantic branch, and instance branch. In the context branch (blue arrows), we extract multi-granularity scene context with a remote-sensing (RS) CLIP model. The semantic branch (yellow arrows) represents the classifier, which consists of text embeddings encoded by a

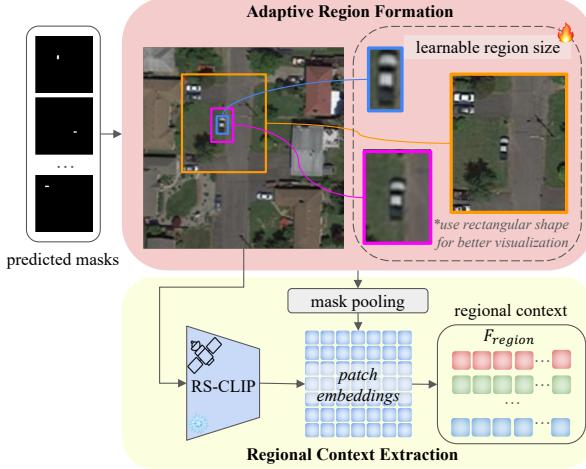


Figure 3. Illustration of Adaptive Region Formation and Regional Context Extraction in Region-Aware Integration.

frozen CLIP text encoder. The instance branch (orange arrows) is responsible for instance proposal generation. A frozen general CLIP image encoder extracts backbone features, which are then fed into the proposal generator to produce class embeddings and mask proposals. These three branches interact through the Region-Aware Integration (RAI) and Global Context Adaptation (GCA) modules. In RAI, regional context is integrated into class embedding, enabling a richer understanding of surrounding environment. In GCA, domain-specific [CLS] token, *i.e.*, global context, is injected into text embeddings, enhancing the classifier’s adaptability for aerial objects. Finally the class prediction is derived from region-aware class embeddings  $\hat{\mathbf{V}}$  and the domain adapted classifier  $\hat{\mathbf{T}}$ .

### 3.3. Scene Context Extraction

To have a domain-specific knowledge about the scene, we utilize remote-sensing vision-language model, *i.e.*, RemoteCLIP [30], to extract scene context from an input image  $I \in \mathbb{R}^{H \times W \times 3}$ . The ViT-based image encoder of RemoteCLIP, denoted as  $\text{CLIP}_{\text{RS}}^*$  (\* indicates frozen), extracts features as follow:

$$\{\mathbf{F}_{\text{CLS}}^i, \mathbf{F}_{\text{HW}}^i\} = \text{CLIP}_{\text{RS}}^*(I), \quad (2)$$

where  $i$  indicates the  $i$ -th transformer layer,  $\mathbf{F}_{\text{CLS}}^i \in \mathbb{R}^{1 \times C}$  is the [CLS] token representing global image features, and  $\mathbf{F}_{\text{HW}}^i \in \mathbb{R}^{\frac{H}{14} \times \frac{W}{14} \times C}$  contains patch embeddings encoding spatially dense image features. The final [CLS] token  $\mathbf{F}_{\text{CLS}}^{\text{final}}$  encapsulates the global context of the image, and the final patch embeddings  $\mathbf{F}_{\text{HW}}^{\text{final}}$  provide rich contextual visual representations for the input image.

### 3.4. Region-Aware Integration

Due to the significant scale variations and resolution disparities in aerial images, it is challenging to accurately recognize small or ambiguous aerial objects. However, remote sensing images contain valuable prior knowledge, as objects often correlate with their surrounding environments [27]. To exploit this correlation, we propose Region-Aware Integration, which dynamically incorporates region-level scene context into class embeddings.

**Adaptive Region Formation.** Given  $N$  class embeddings  $\mathbf{V} = [v_1, v_2, \dots, v_N] \in \mathbb{R}^{N \times C}$  and their corresponding predicted mask proposals  $\mathbf{M} = [m_1, m_2, \dots, m_N] \in \mathbb{R}^{B \times N \times H \times W}$ , we suggest that the predicted masks  $M$  serve as spatial references for each object instance. To get the surrounding region, we design an learnable dilation mechanism that adaptively expands predicted masks. We introduce a learnable dilation factor  $\delta$  that adjusts the expansion dynamically. We define the expanded mask as:

$$\mathbf{M}' = \max_{\mathbf{x} \in \mathcal{N}(M, k)} \mathbf{M}(\mathbf{x}), \quad (3)$$

where  $\mathcal{N}(M, k)$  denotes the local neighborhood defined by a learnable kernel size  $k$ , computed as:

$$k = 3 + \text{clamp}(\delta, 0, 10), \quad (4)$$

where  $\delta$  is a learnable parameter initialized to 1 and optimized during training. The expansion is performed via max-pooling with kernel size  $k$ , stride 1, and padding  $k//2$  to ensure spatial alignment.

**Regional Context Extraction.** Once the expanded masks are obtained, we leverage the patch embeddings  $\mathbf{F}_{\text{HW}}^{\text{final}}$ , extracted from RemoteCLIP, to encode scene-level semantics. These embeddings capture dense visual features across the input image. To extract the regional context for each query object, with the corresponding expanded mask  $m'_i$ , we pool all the features in  $\mathbf{F}_{\text{HW}}^{\text{final}}$  that are within the expanded mask to compute a mask pooled image feature as follows:

$$F_{\text{region}} = \sum_{\mathbf{x} \in M'} \omega_{\mathbf{x}} \cdot \mathbf{F}_{\text{HW}}^{\text{final}}, \quad (5)$$

with  $\omega_{\mathbf{x}}$  denotes the normalized weight within the mask  $m'_i$ . The adaptive region formation and context extraction process is illustrated in Figure 3.

**Regional Context Integration.** With the extracted regional context  $F_{\text{region}}$ , we then integrate it into the class embeddings  $\mathbf{V}$  through  $l$  sequential Transformer Layers:

$$\mathbf{V}_{i+1} = \text{TransLayer}_i(\mathbf{V}_i, \lambda \cdot \mathbf{F}_{\text{region}}), \quad (6)$$

where  $i = 1, 2, \dots, l$  and  $\lambda$  is a temperature coefficient that controls the contribution of the regional context. Finally, we get region-aware class embeddings  $\hat{\mathbf{V}}$ .

This process seamlessly inject the regional scene context into class embeddings, utilizing domain-specific knowledge pretrained in remote sensing-specific CLIP. By leveraging the correlation between objects and their surrounding environments, the model gains a richer contextual understanding, improving its ability to distinguish aerial objects. Additionally, integrating scene information helps recalibrate the feature space of class embeddings trained on seen categories  $\mathbf{C}_{\text{train}}$ , aligning them with a more generalized latent space. This adjustment mitigates overfitting to training categories, enabling better generalization to novel, unseen objects in open-vocabulary remote sensing instance segmentation.

### 3.5. Global Context Adaptation

Associating predicted instances with their regional context improves their visual distinguishability, yet the performance is still limited by the frozen classifier built on CLIP text embeddings. Although the classifier, derived from general-domain CLIP, exhibits strong generalization abilities, its fixed text embeddings lack the specificity needed to accurately distinguish remote sensing objects. This limitation arises from factors such as high intra-class variation, resolution differences, and complex environmental conditions. To address this, we complement Region-Aware Integration, which enhances the visual modality, by further refining the textual modality. Our approach bridges the semantic gap between general-domain knowledge and remote sensing-specific concepts, which improves the classifier's adaptability for open-vocabulary remote sensing instance segmentation. Specifically, we incorporate global scene context to adapt domain-specific visual priors into the text embeddings, ensuring a more robust alignment between visual and textual representations.

The classifier  $\mathbf{T}$  consists of text embeddings encoded by the frozen CLIP text encoder. Given a set of category names  $\mathbf{C}_{\text{train}}$  represented in natural language, we derive category text embeddings by place the category name into prompt templates, *e.g.*, ‘satellite imagery of .’, ‘aerial imagery of .’, and then fed into CLIP text encoder. The text embedding for each class is the average across all templates. We denote the classifier as  $\mathbf{T} \in \mathbb{R}^{M \times C}$ , where  $M$  is the number of categories,  $C$  is the dimension of text embeddings.

To equip the classifier with remote sensing knowledge, global visual context, *i.e.*,  $\mathbf{F}_{\text{CLS}}^{\text{final}}$  from RemoteCLIP is injected into the text embeddings using multi-head cross-attention with  $\mathbf{T}$ , formulated as:

$$\begin{aligned} Q &= \mathbf{w}_Q \mathbf{F}_{\text{CLS}}^{\text{final}}, K = \mathbf{w}_K \mathbf{T}, V = \mathbf{w}_V \mathbf{T}, \\ \hat{\mathbf{T}} &= \text{MHA}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V, \end{aligned} \quad (7)$$

where  $\mathbf{w}_Q, \mathbf{w}_K, \mathbf{w}_V$  are linear projection matrices that map inputs to the shared feature dimension  $d_k$ . This adaptation process ensures that the text embeddings retain their open-vocabulary generalization while being enriched with remote sensing-specific global context.

By integrating domain-specific global context into the text embeddings, our approach facilitates cross-domain interactions to produce more representative and adaptable textual features tailored for aerial images. Since these context are extracted from visual representations, the method also strengthen cross-modal interactions between vision and language, promoting better semantic alignment. As a result, the classifier becomes more robust in handling the variability of remote sensing objects, ultimately leading to improved performance across diverse datasets.

### 3.6. Open-Vocabulary Inference

During inference, the classification score is computed using an ensemble approach [54, 63] that combines both in-vocabulary and out-vocabulary classification. The in-vocabulary classification is based on the learned region-aware class embeddings and then classified with text embeddings enhanced with domain global context. For out-vocabulary classification, our framework offers flexibility in implementation. By leveraging the instance branch (built on a general CLIP backbone) and the context branch (using a remote sensing-specific CLIP backbone), dense visual features can be extracted from both general and remote sensing domains. While the use of remote sensing CLIP models appears promising, our experimental results in Section 4.4 indicate that they still fall short compared to general-domain CLIP in our task. This observation aligns with the results reported in [61]. As a result, we continue to utilize the general CLIP to get out-vocabulary classification due to its superior generalization capabilities.

## 4. Experiments

### 4.1. Datasets and Evaluation Metric

**Training Dataset.** Following the open-vocabulary benchmarks for natural images [55], we train the model on one dataset and evaluate its cross-dataset performance on other datasets. We select two datasets for training, *i.e.*, iSAID [64] and SIOR [48]. iSAID is a large scale instance segmentation dataset for remote sensing images. It contains 18732 images for training across 15 categories. SIOR is developed from aerial object detection dataset DIOR [25], with segmentation annotations generated in SAMRS [48]. It contains 11725 images with 20 aerial categories.

**Evaluation Dataset.** To evaluate the effectiveness of our method, we conduct cross-dataset evaluation on 4 aerial instance segmentation datasets, *i.e.*, NWPU-VHR-10 [7, 44], SOTA [48], FAST [48], and SIOR [48]. NWPU-VHR-10 is

Method	iSAID					SIOR				
	NWPU	SOTA	FAST	SIOR	Average	NWPU	SOTA	FAST	iSAID	Average
ODISE [54] [CVPR22]	36.40	13.91	4.65	13.68	17.16	41.77	12.3	5.88	12.68	18.16
FC-CLIP [63] [NeurIPS23]	60.67	33.62	11.88	26.79	33.24	60.69	19.84	8.67	22.24	27.86
MAFT+ [22] [ECCV24]	35.32	6.63	5.52	9.84	14.33	39.99	7.77	5.92	9.63	15.83
ZoRI [18] [AAAI25]	62.06	30.02	12.65	26.27	32.75	59.77	20.26	9.58	23.46	28.27
<b>SCORE (Ours)</b>	<b>67.59</b>	<b>42.57</b>	<b>13.67</b>	<b>30.90</b>	<b>38.68</b>	<b>69.17</b>	<b>23.68</b>	<b>10.33</b>	<b>27.15</b>	<b>32.59</b>

Table 1. **Comparison with SOTA methods on open-vocabulary remote sensing instance segmentation.** The model is trained separately on iSAID and SIOR datasets and then tested on the remaining four datasets to measure its cross-dataset generalization capabilities.

Modules		iSAID					SIOR				
RAI	GCA	NWPU	SOTA	FAST	SIOR	Average	NWPU	SOTA	FAST	iSAID	Average
✗	✗	58.59	36.44	11.56	26.43	33.25	61.02	20.85	9.14	22.44	28.36
✓	✗	66.32	39.55	12.85	28.91	36.91	68.47	21.28	9.70	25.16	31.15
✗	✓	67.21	38.14	12.37	28.96	36.67	69.07	22.05	9.79	24.61	31.43
✓	✓	<b>67.59</b>	<b>42.57</b>	<b>13.67</b>	<b>30.90</b>	<b>38.68</b>	<b>69.17</b>	<b>23.68</b>	<b>10.33</b>	<b>27.15</b>	<b>32.59</b>

Table 2. **Component analysis of SCORE.** The cross-dataset evaluation results based on models trained on the iSAID and SIOR datasets are provided. RAI denotes Region-Aware Integration and GCA denotes Global Context Adaptation.

an aerial object detection dataset with instance masks further annotated by [44]. The test set contains 731 images from 10 aerial classes. SOTA, FAST and SIOR are segmentation datasets provided in SAMRS [48], which are developed from aerial object detection datasets DOTA-V2.0 [11], FAIR1M-2.0 [47], and DIOR [25], respectively. SOTA covers 874 images with 18 object categories for testing, FAST contains 3207 images across 37 fine-grained aerial categories for testing, and SIOR comprises 11738 testing samples across 20 classes. Additional dataset details can be found in the supplementary material.

**Evaluation Metric.** We follow previous work [39] to evaluate the performance with mean Average Precision (mAP) for open vocabulary instance segmentation.

## 4.2. Implementation Details

For instance branch, we use frozen ConvNeXt-Large CLIP as the backbone [35] with the weight pretrained on LAION-2B from OpenCLIP [19]. Its generalization to high-resolution inputs is better suited to extract features for segmentation instead of using a ViT-based CLIP. The proposal generator follows Mask2Former [6] with number of object query set to 300. In semantic branch, prompt templates for remote sensing images RESISC45 [8] are employed to obtain text embeddings with the pretrained CLIP text encoder. For context branch, we employ RemoteCLIP ViT-L/14 [30] to extract multi-granularity scene context. All ablation experiments are trained for 50 epochs with batch size set to 2, on one L40S GPU. The learning rate is set to  $1.25 \times 10^{-5}$ . Input images are resized to  $512 \times 512$  during training. The model is optimized using AdamW optimizer.

## 4.3. Comparison with State-of-the-Art Methods

Table 1 presents a comparative analysis of our method **SCORE** against SOTA models training on the iSAID and SIOR datasets respectively. Our approach consistently outperforms existing methods across all benchmarks, demonstrating its effectiveness in open-vocabulary remote sensing instance segmentation. Specifically, we achieve an average improvement of 5.53% and 4.32% over the best-performing SOTA models trained on iSAID and SIOR, respectively. Our method surpasses previous methods by a large margin, especially on NWPU, SOTA, SIOR (trained on iSAID) and iSAID (trained on SIOR) datasets, with gains up to 12.55%.

## 4.4. Ablation Studies

We perform a series of ablation experiments to assess the effectiveness of our proposed modules. The component analysis presents cross-dataset evaluation results based on models trained on the iSAID and SIOR datasets. Meanwhile, the remaining ablation studies are conducted solely on experiments where iSAID serves as the training dataset.

**Component Analysis.** The results, as shown in Table 2, highlight the contributions of Region-Aware Integration (RAI) and Global Context Adaptation (GCA) modules, which are central to facilitating OV remote sensing instance segmentation. When RAI is enabled, the mAP increases nearly 8% for NWPU dataset trained in iSAID and SIOR respectively, compared to the baseline configuration. We can also observe an average performance gain of around 3 % for the remaining three datasets. This improvement underscores the effectiveness of incorporating regional context to refine class embeddings, retrieving region-aware class embedding

VLM	NWPU	SOTA	FAST	SIOR
CLIP [40]	64.03	38.74	11.42	28.11
SkyCLIP [50]	65.04	33.57	12.43	28.96
GeoRSCLIP [68]	64.72	39.33	12.56	28.13
RemoteCLIP [30]	<b>67.59</b>	<b>42.57</b>	<b>13.67</b>	<b>30.90</b>

Table 3. **Ablation Study – VLMs for Scene Context Extraction.** Comparison of various VLMs used in the context branch for multi-granularity scene context extraction. We include general CLIP and three RS CLIPs for comparison.

Choice of Context	NWPU	SOTA	FAST	SIOR
[CLS] token	67.06	40.35	12.58	29.07
Patch embeddings	66.14	37.95	12.60	28.97
Regional context	<b>67.59</b>	<b>42.57</b>	<b>13.67</b>	<b>30.90</b>

Table 4. **Ablation Study – Choice of Context in RAI.** Comparison of different context that can be used in Region-Aware Integration.

and calibrating the feature to a more generalizable space. Additionally, VCA achieves an average improvement of 3.42% on iSAID and 3.07% on SIOR. These results validate our approach of enriching the linguistic latent space with global visual context, making the classifier more adaptable to remote sensing objects. By combining both RAI and VCA, our models yields the best performance, demonstrating their complementary roles: RAI refines class embedding with regional context in the visual space, while VPA enhances text embeddings by incorporating global scene context.

**VLMs for Scene Context Extraction.** As mentioned in the related work 2.3, several remote sensing (RS) CLIPs have been developed recently. We evaluate the impact of different VLMs for scene context extraction, comparing the general CLIP [40] with three popular remote sensing-specific CLIPs, *i.e.*, SkyCLIP [50], GeoRSCLIP [68], and RemoteCLIP [30]. To ensure a fair comparison, we use the ViT-L/14 pretrained weights across all models. As shown in Table 3, RS CLIPs outperform the general CLIP on almost all test datasets. This highlights the advantages of domain-specific pretraining in capturing more accurate and specialized scene context for aerial images. Among them, RemoteCLIP achieves the best performance and is therefore employed in our framework.

**Choice of Context in RAI.** We explore various context that can be applied for RAI module in Table 4. As a starting point, we utilize the [CLS] token from RemoteCLIP, which encapsulates the global context of the input image. The performance gain over the baseline supports our intuition that integrating scene context into class prediction can help identify instances in remote sensing imagery. However, it still falls short compared to our proposed regional context. This could be attributed to the diverse concepts and semantics that can include in large-scale remote sensing imagery. Since the [CLS] token extracts high-level semantic features, it tends to focus on the dominant components within the image, introducing a global bias [26]. As a result, the

Method	NWPU	SOTA	FAST	SIOR
w/o injection	66.32	39.55	12.85	28.91
add	65.72	39.89	11.11	29.86
concat	47.53	15.51	1.74	18.23
MHA	<b>67.59</b>	<b>42.57</b>	<b>13.67</b>	<b>30.90</b>

Table 5. **Ablation Study – Adaptation Methods in GCA.** Comparison of different global visual context injection methods in Global Context Adaptation.

Context VLM	OV VLM	NWPU	SOTA	FAST	SIOR
SkyCLIP [50]	SkyCLIP	64.57	32.49	12.44	27.22
	CLIP	65.04	33.57	12.42	28.96
GeoRSCLIP [68]	GeoRSCLIP	63.67	38.76	12.55	27.15
	CLIP	64.72	39.33	12.56	28.13
RemoteCLIP [30]	RemoteCLIP	66.05	40.26	11.65	29.43
	CLIP	67.59	42.57	13.67	30.90

Table 6. **Ablation Study – Open-Vocabulary Classification.** The context VLM denotes different RS CLIPs used for scene context extraction in the context branch, while OV VLM refers to the model used for open-vocabulary classification. The results mainly investigate the impact of employing domain-specific CLIP versus general CLIP in open-vocabulary classification.

model may struggle to accurately predict less prominent objects. Additionally, intermediate patch embeddings  $F_{HW}^i$  yield the lowest performance, likely because mid-layer representations emphasize texture over semantic information, introducing unrelated noise to the process. By contrast, our regional context, which is adaptively formulated, allows the model to leverage the local context of the target object while mitigating the effects of global bias and noise.

**Adaptation Methods in GCA.** To inject remote sensing-specific global context into the text embedding, different integration methods can be used. As listed in Table 5, we experiment with straightforward approaches such as addition and concatenation, both of which do not involve trainable parameters. The results indicate that these methods fail to produce a generalizable classifier, even performing worse than the naive classifier without global context injection. We attribute this to the misalignment between the RS CLIP visual embeddings and the general CLIP textual embeddings. Simply combining these two without learnable parameters may disrupt the well-pretrained cross-modal alignment. Instead, by treating the global context as a query and allowing it to cross-attend to the text embeddings  $T$ , we facilitate a smoother interaction between the specialized and general CLIP representations.

**Out-Vocabulary Classification.** During inference, out-vocabulary (OV) classification is used to further enhance novel class prediction. Our framework employs RS CLIP for scene context extraction in the context branch and general CLIP as the backbone in the instance branch. Thus, we provide flexible options for OV classification. As

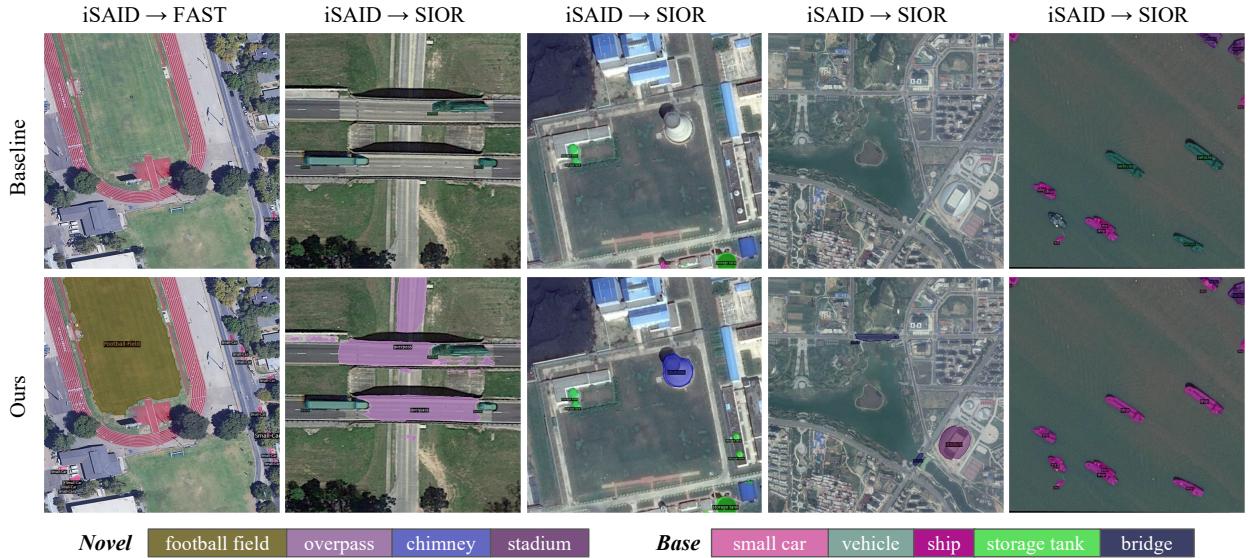


Figure 4. **Qualitative comparisons between the baseline and our model.** “iSAID → FAST” denotes training on iSAID and testing on FAST. We set jittering=False for better readability. **Novel** and **Base** indicate whether a class is absent or present in the training dataset, respectively. **SCORE** effectively segments novel classes such as *football field* (column 1), *overpass* (column 2), *chimney* (column 3), *stadium* (column 4). Furthermore, our model demonstrates strong instance segmentation capabilities in large-scale images, successfully segmenting *small car* (column 1) and *stadium* (column 4). Additionally, in the last column, we correctly identify all *ship* instances, whereas the baseline model confuses them with *vehicle*.

shown in Table 6, the context VLM refers to different RS CLIPs used for scene context extraction, while OV VLM represents the models from which our OV classification is derived. We conduct three sets of experiments based on different RS CLIPs used for scene context extraction. The trained model is then evaluated with (1) the same context VLM as the OV VLM (2) general CLIP as the OV VLM. From the results, we can see that, although RS CLIPs can produce reasonable performance, they consistently underperform general-domain CLIP in our task. This observation aligns with the results provided in [61] for OV semantic segmentation. Consequently, we adopt general CLIP to get OV classification probability, which is then combined with in-vocabulary classification to yield the final prediction. We think the performance gap primarily comes from the limited availability of image-caption pairs for cross-modal pretraining in RS CLIPs. The remote sensing pretraining datasets in these models range from 0.8M to 5M image-text pairs, which is considerably smaller than 400M web-scale data used in general CLIP. As a result, RS CLIPs acquire domain-specific knowledge but exhibit weaker generalization capabilities.

**Qualitative Results.** In Figure 4, we provide visualizations of our performance on remote sensing instance segmentation task. **SCORE** effectively segments the instances of novel classes such as *football field* (column 1), *overpass* (column 2), *chimney* (column 3) and *stadium* (column 4) across different test datasets, demonstrating

strong class prediction capabilities through the synergy of RAI and VCA modules. Furthermore, our model exhibits robustness in segmenting small instances in large-scale images. For instance, in the 1st column, most *small car* instances are segmented, whereas the baseline fails to segment any. In the 4th column, we successfully segment both novel class *stadium* and base class *bridge*, even in an extremely large-scale image captured from a high altitude. Additionally, in the last column, many instances of *ship* are misclassified as *vehicle* by the baseline model, whereas **SCORE** correctly identifies all of them as *ship*. This aligns with our motivation that surrounding context aids in aerial object recognition. More cross-dataset qualitative results are provided in the supplementary material.

## 5. Conclusion

In conclusion, we introduce the task of open-vocabulary remote sensing instance segmentation and develop **SCORE**, a comprehensive framework that integrates domain-specific scene context. By leveraging the correlation between objects and their surroundings, RAI incorporates regional context to class embeddings to gain richer contextual understanding. GCA injects domain-specific global context into text embeddings to create a more adaptive classifier for remote sensing objects. Extensive experiments validate the effectiveness of our approach, and with its state-of-the-art performance, we provide a reliable solution for large-scale real-world geospatial analysis.

**Acknowledgments.** This research is supported in part by the National Research Foundation Singapore Competitive Research Program (award number CRP29-2022-0003). The work was done at Rapid-Rich Object Search (ROSE) Lab, School of Electrical & Electronic Engineering, Nanyang Technological University. Shuting He is sponsored by Shanghai Pujiang Programme 24PJD030 and Natural Science Foundation of Shanghai 25ZR1402138.

## References

- [1] Qinglong Cao, Yuntian Chen, Chao Ma, and Xiaokang Yang. Open-vocabulary remote sensing image semantic segmentation. *arXiv preprint arXiv:2409.07683*, 2024. 2, 3
- [2] Xu Cao, Huanxin Zou, Jun Li, Xinyi Ying, and Shitian He. Obbinst: Remote sensing instance segmentation with oriented bounding box supervision. *International Journal of Applied Earth Observation and Geoinformation*, 128:103717, 2024. 1
- [3] Osmar Luiz Ferreira de Carvalho, Osmar Abilio de Carvalho Junior, Anesmar Olino de Albuquerque, Pablo Pozzobon de Bem, Cristiano Rosa Silva, Pedro Henrique Guimaraes Ferreira, Rebeca dos Santos de Moura, Roberto Arnaldo Trancoso Gomes, Renato Fontes Guimaraes, and Dibio Leandro Borges. Instance segmentation for large, multi-channel remote sensing imagery using mask-rnn and a mosaicking approach. *Remote Sensing*, 2020. 1
- [4] Keyan Chen, Chenyang Liu, Hao Chen, Haotian Zhang, Wenyuan Li, Zhengxia Zou, and Zhenwei Shi. Rsprompter: Learning to prompt for remote sensing instance segmentation based on visual foundation model. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–17, 2024. 2
- [5] Xi Chen, Shuang Li, Ser-Nam Lim, Antonio Torralba, and Hengshuang Zhao. Open-vocabulary panoptic segmentation with embedding modulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1141–1150, 2023. 2
- [6] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1290–1299, 2022. 6
- [7] Gong Cheng, Junwei Han, Peicheng Zhou, and Lei Guo. Multi-class geospatial object detection and geographic image classification based on collection of part detectors. *ISPRS Journal of Photogrammetry and Remote Sensing*, 2014. 5, 1, 2
- [8] Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Benchmark and state of the art. *Proceedings of the IEEE*, 2017. 6
- [9] Gong Cheng, Xingxing Xie, Junwei Han, Lei Guo, and Gui-Song Xia. Remote sensing image scene classification meets deep learning: Challenges, methods, benchmarks, and opportunities. *J-STARS*, 2020. 2
- [10] Seokju Cho, Heeseong Shin, Sunghwan Hong, Anurag Arnab, Paul Hongsuck Seo, and Seungryong Kim. Cat-seg: Cost aggregation for open-vocabulary semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4113–4123, 2024. 2, 1
- [11] Jian Ding, Nan Xue, Gui-Song Xia, Xiang Bai, Wen Yang, Michael Ying Yang, Serge Belongie, Jiebo Luo, Mihai Datcu, Marcello Pelillo, et al. Object detection in aerial images: A large-scale benchmark and challenges. *IEEE transactions on pattern analysis and machine intelligence*, 44(11):7778–7796, 2021. 6, 1
- [12] Zheng Ding, Jieke Wang, and Zhuowen Tu. Open-vocabulary universal image segmentation with maskclip. In *Proceedings of the 40th International Conference on Machine Learning*, pages 8090–8102, 2023. 2
- [13] Anatol Garioud, Nicolas Gonthier, Loic Landrieu, Apolline De Wit, Marion Valette, Marc Poupee, Sébastien Giordano, et al. Flair: a country-scale land cover semantic segmentation dataset from multi-source optical imagery. *Advances in Neural Information Processing Systems*, 36:16456–16482, 2023. 1, 2
- [14] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Scaling open-vocabulary image segmentation with image-level labels. In *European conference on computer vision*, pages 540–557. Springer, 2022. 2
- [15] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. *arXiv preprint arXiv:2104.13921*, 2021. 3
- [16] Zhihao Guan, Xinyu Miao, Yunjie Mu, Quan Sun, Qiaolin Ye, and Demin Gao. Forest fire segmentation from aerial imagery data using an improved instance segmentation model. *Remote Sensing*, 14(13):3159, 2022. 1
- [17] Taisei Hanyu, Kashu Yamazaki, Minh Tran, Roy A McCann, Haitao Liao, Chase Rainwater, Meredith Adkins, Jackson Cothren, and Ngan Le. Aerialformer: Multi-resolution transformer for aerial image segmentation. *Remote Sensing*, 16(16):2930, 2024. 1, 2
- [18] Shiqi Huang, Shuting He, and Bihan Wen. Zori: Towards discriminative zero-shot remote sensing instance segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3724–3732, 2025. 2, 6
- [19] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, 2021. 6
- [20] ISPRS. 2d semantic labeling potsdam dataset. <https://www.isprs.org/education/benchmarks/UrbanSemLab/2d-sem-label-potsdam.aspx>, 2013. Accessed: 2024-08-11. 1, 2
- [21] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pages 4904–4916. PMLR, 2021. 3
- [22] Siyu Jiao, Hongguang Zhu, Jiannan Huang, Yao Zhao, Yunchao Wei, and Humphrey Shi. Collaborative vision-text representation optimizing for open-vocabulary segmentation.

- In *European Conference on Computer Vision*, pages 399–416. Springer, 2024. [2](#), [3](#), [6](#)
- [23] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. [2](#)
- [24] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. *arXiv preprint arXiv:2201.03546*, 2022. [2](#)
- [25] Ke Li, Gang Wan, Gong Cheng, Liqiu Meng, and Junwei Han. Object detection in optical remote sensing images: A survey and a new benchmark. *ISPRS journal of photogrammetry and remote sensing*, 159:296–307, 2020. [5](#), [6](#), [1](#)
- [26] Kaiyu Li, Ruixun Liu, Xiangyong Cao, Xueru Bai, Feng Zhou, Deyu Meng, and Zhi Wang. Segearth-ov: Towards training-free open-vocabulary segmentation for remote sensing images. *arXiv preprint arXiv:2410.01768*, 2024. [7](#)
- [27] Yansheng Li, Deyu Kong, Yongjun Zhang, Yihua Tan, and Ling Chen. Robust deep alignment network with remote sensing knowledge graph for zero-shot and generalized zero-shot remote sensing image scene classification. *ISPRS Journal of Photogrammetry and Remote Sensing*, 179:145–158, 2021. [2](#), [4](#)
- [28] Yuxuan Li, Qibin Hou, Zhaozhui Zheng, Ming-Ming Cheng, Jian Yang, and Xiang Li. Large selective kernel network for remote sensing object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 16794–16805, 2023. [2](#)
- [29] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7061–7070, 2023. [2](#)
- [30] Fan Liu, Delong Chen, Zhangqinyun Guan, Xiaocong Zhou, Jiale Zhu, Qiaolin Ye, Liyong Fu, and Jun Zhou. Remoteclip: A vision language foundation model for remote sensing. *IEEE Transactions on Geoscience and Remote Sensing*, 2024. [2](#), [3](#), [4](#), [6](#), [7](#)
- [31] Xinyu Liu and Xiaoguang Di. Global context parallel attention for anchor-free instance segmentation in remote sensing images. *IEEE Geoscience and Remote Sensing Letters*, 19:1–5, 2020. [2](#)
- [32] Yuanyuan Liu, Dingyuan Chen, Ailong Ma, Yanfei Zhong, Fang Fang, and Kai Xu. Multiscale u-shaped cnn building instance extraction framework with edge constraint for high-spatial-resolution remote sensing imagery. *IEEE Transactions on Geoscience and Remote Sensing*, 59(7):6106–6120, 2020. [1](#)
- [33] Ye Liu, Huifang Li, Chao Hu, Shuang Luo, Yan Luo, and Chang Wen Chen. Learning to aggregate multi-scale context for instance segmentation in remote sensing images. *IEEE Transactions on Neural Networks and Learning Systems*, 2024. [1](#), [2](#)
- [34] Zheyuan Liu, Cristian Rodriguez-Opazo, Damien Teney, and Stephen Gould. Image retrieval on real-life images with pre-trained vision-and-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2125–2134, 2021. [3](#)
- [35] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s, 2022. [6](#)
- [36] Alina Marcu and Marius Leordeanu. Dual local-global contextual pathways for recognition in aerial imagery. *arXiv preprint arXiv:1605.05462*, 2016. [2](#)
- [37] Matthias Minderer, Alexey Gritsenko, and Neil Houlsby. Scaling open-vocabulary object detection. *Advances in Neural Information Processing Systems*, 36:72983–73007, 2023. [3](#)
- [38] Chao Pang, Jiang Wu, Jiayu Li, Yi Liu, Jiaxing Sun, Weijia Li, Xingxing Weng, Shuai Wang, Litong Feng, Gui-Song Xia, et al. H2rsvlm: Towards helpful and honest remote sensing large vision language model. *arXiv e-prints*, pages arXiv–2403, 2024. [2](#), [3](#)
- [39] Jie Qin, Jie Wu, Pengxiang Yan, Ming Li, Ren Yuxi, Xuefeng Xiao, Yitong Wang, Rui Wang, Shilei Wen, Xin Pan, et al. Freeseg: Unified, universal and open-vocabulary image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19446–19455, 2023. [2](#), [6](#)
- [40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. [2](#), [3](#), [7](#)
- [41] Maryam Rahnmemoonfar, Tashnim Chowdhury, Argho Sarkar, Debraj Varshney, Masoud Yari, and Robin Roberson Murphy. Floodnet: A high resolution aerial imagery dataset for post flood scene understanding. *IEEE Access*, 9:89644–89654, 2021. [1](#), [2](#)
- [42] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou, and Jiwen Lu. Denseclip: Language-guided dense prediction with context-aware prompting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18082–18091, 2022. [3](#)
- [43] Abubakar Sani-Mohammed, Wei Yao, and Marco Heurich. Instance segmentation of standing dead trees in dense forest from aerial imagery using deep learning. *ISPRS Open Journal of Photogrammetry and Remote Sensing*, 6:100024, 2022. [1](#)
- [44] Hao Su, Shunjun Wei, Min Yan, Chen Wang, Jun Shi, and Xiaoling Zhang. Object detection and instance segmentation in remote sensing imagery based on precise mask r-cnn. In *IGARSS 2019-2019 IEEE International Geoscience and Remote Sensing Symposium*, pages 1454–1457. IEEE, 2019. [1](#), [2](#), [5](#), [6](#)
- [45] Hao Su, Shunjun Wei, Shan Liu, Jiadian Liang, Chen Wang, Jun Shi, and Xiaoling Zhang. Hq-isnet: High-quality instance segmentation for remote sensing imagery. *Remote Sensing*, 2020. [1](#)
- [46] Hao Su, Peng Huang, Jun Yin, and Xiaofeng Zhang. Faster and better instance segmentation for large scene remote sensing imagery. In *IGARSS 2022-2022 IEEE international geoscience and remote sensing symposium*, pages 2187–2190. IEEE, 2022. [2](#)

- [47] Xian Sun, Peijin Wang, Zhiyuan Yan, Feng Xu, Ruiping Wang, Wenhui Diao, Jin Chen, Jihao Li, Yingchao Feng, Tao Xu, et al. Fair1m: A benchmark dataset for fine-grained object recognition in high-resolution remote sensing imagery. *ISPRS Journal of Photogrammetry and Remote Sensing*, 184: 116–130, 2022. [6](#), [1](#)
- [48] Di Wang, Jing Zhang, Bo Du, Minqiang Xu, Lin Liu, Dacheng Tao, and Liangpei Zhang. Samrs: Scaling-up remote sensing segmentation dataset with segment anything model. *Advances in Neural Information Processing Systems*, 36:8815–8827, 2023. [2](#), [5](#), [6](#), [1](#)
- [49] Zhaoqing Wang, Yu Lu, Qiang Li, Xunqiang Tao, Yandong Guo, Mingming Gong, and Tongliang Liu. Cris: Clip-driven referring image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11686–11695, 2022. [3](#)
- [50] Zhecheng Wang, Rajanie Prabha, Tianyuan Huang, Jiajun Wu, and Ram Rajagopal. Skyscript: A large and semantically diverse vision-language dataset for remote sensing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5805–5813, 2024. [2](#), [3](#), [7](#)
- [51] Shunjun Wei, Xiangfeng Zeng, Hao Zhang, Zichen Zhou, Jun Shi, and Xiaoling Zhang. Lfg-net: Low-level feature guided network for precise ship instance segmentation in sar images. *IEEE Transactions on Geoscience and Remote Sensing*, 60: 1–17, 2022. [1](#)
- [52] Jianzong Wu, Xiangtai Li, Shilin Xu, Haobo Yuan, Henghui Ding, Yibo Yang, Xia Li, Jiangning Zhang, Yunhai Tong, Xudong Jiang, et al. Towards open vocabulary learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(7):5092–5113, 2024. [1](#)
- [53] Tong Wu, Yuan Hu, Ling Peng, and Ruohan Chen. Improved anchor-free instance segmentation for building extraction from high-resolution remote sensing images. *Remote Sensing*, 12(18):2910, 2020. [1](#)
- [54] Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong Wang, and Shalini De Mello. Open-vocabulary panoptic segmentation with text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2955–2966, 2023. [2](#), [5](#), [6](#)
- [55] Mengde Xu, Zheng Zhang, Fangyun Wei, Yutong Lin, Yue Cao, Han Hu, and Xiang Bai. A simple baseline for open-vocabulary semantic segmentation with pre-trained vision-language model. In *European Conference on Computer Vision*, pages 736–753. Springer, 2022. [1](#), [5](#)
- [56] Mengde Xu, Zheng Zhang, Fangyun Wei, Han Hu, and Xiang Bai. Side adapter network for open-vocabulary semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2945–2954, 2023. [2](#)
- [57] Xiangkai Xu, Zhejun Feng, Changqing Cao, Mengyuan Li, Jin Wu, Zengyan Wu, Yajie Shang, and Shubing Ye. An improved swin transformer-based model for remote sensing object detection and instance segmentation. *Remote Sensing*, 2021. [1](#)
- [58] Xiwen Yao, Junwei Han, Gong Cheng, Xueming Qian, and Lei Guo. Semantic annotation of high-resolution satellite images via weakly supervised learning. *IEEE Transactions on Geoscience and Remote Sensing*, 54(6):3660–3671, 2016. [1](#), [2](#)
- [59] Muhammad Yasir, Lili Zhan, Shanwei Liu, Jianhua Wan, Md Sakaouth Hossain, Arife Tugsan Isiacik Colak, Mengge Liu, Qamar Ul Islam, Syed Raza Mehdi, and Qian Yang. Instance segmentation ship detection based on improved yolov7 using complex background sar images. *Frontiers in Marine Science*, 2023. [1](#)
- [60] Chengyang Ye, Yunzhi Zhuge, and Pingping Zhang. Towards open-vocabulary remote sensing image semantic segmentation. *arXiv preprint arXiv:2412.19492*, 2024. [1](#)
- [61] Chengyang Ye, Yunzhi Zhuge, and Pingping Zhang. Towards open-vocabulary remote sensing image semantic segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025. [2](#), [3](#), [5](#), [8](#), [1](#)
- [62] Wenhui Ye, Wei Zhang, Weimin Lei, Wencho Zhang, Xinyi Chen, and Yanwen Wang. Remote sensing image instance segmentation network with transformer and multi-scale feature representation. *Expert Systems with applications*, 234:121007, 2023. [1](#), [2](#)
- [63] Qihang Yu, Ju He, Xueqing Deng, Xiaohui Shen, and Liang-Chieh Chen. Convolutions die hard: Open-vocabulary segmentation with single frozen convolutional clip. *Advances in Neural Information Processing Systems*, 36:32215–32234, 2023. [1](#), [2](#), [5](#), [6](#)
- [64] Syed Waqas Zamir, Aditya Arora, Akshita Gupta, Salman Khan, Guolei Sun, Fahad Shahbaz Khan, Fan Zhu, Ling Shao, Gui-Song Xia, and Xiang Bai. isaid: A large-scale dataset for instance segmentation in aerial images, 2019. [1](#), [5](#), [2](#)
- [65] Hao Zhang, Feng Li, Xueyan Zou, Shilong Liu, Chunyuan Li, Jianwei Yang, and Lei Zhang. A simple framework for open-vocabulary segmentation and detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1020–1031, 2023. [2](#)
- [66] Tianyang Zhang, Xiangrong Zhang, Peng Zhu, Xu Tang, Chen Li, Licheng Jiao, and Huiyu Zhou. Semantic attention and scale complementary network for instance segmentation in remote sensing images. *IEEE Transactions on Cybernetics*, 2021. [1](#), [2](#)
- [67] Xin Zhang, Liangxiu Han, Lianghao Han, and Liang Zhu. How well do deep learning-based methods for land cover classification and object detection perform on high resolution remote sensing imagery? *Remote Sensing*, 2020. [2](#)
- [68] Zilun Zhang, Tiancheng Zhao, Yulong Guo, and Jianwei Yin. Rs5m: A large scale vision-language dataset for remote sensing vision-language foundation model. *arXiv preprint arXiv:2306.11300*, 2(3):6, 2023. [2](#), [3](#), [7](#)
- [69] Bo Zhong, Tengfei Wei, Xiaobo Luo, Bailin Du, Longfei Hu, Kai Ao, Aixia Yang, and Junjun Wu. Multi-swin mask transformer for instance segmentation of agricultural field extraction. *Remote sensing*, 15(3):549, 2023. [1](#)
- [70] Yuqian Zhu, Weitao Chen, Wenxi He, Ruizhen Wang, Xianju Li, and Lizhe Wang. Cug\_misdataset: A remote sensing instance segmentation dataset for improved wide-area high-precision mining land occupation recognition. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2024. [1](#)

# SCORE: Scene Context Matters in Open-Vocabulary Remote Sensing Instance Segmentation

## Supplementary Material

In the supplementary materials, we provide more information on the datasets used for open-vocabulary remote sensing instance segmentation benchmark and include more qualitative results along with comparisons. Moreover, we show that **SCORE** can also enhance the performance for open-vocabulary remote sensing semantic segmentation task, which further unleashes the potential of our model.

### A. Implementation Details

#### A.1. Remote Sensing Instance Segmentation

**Training Dataset.** Following the open-vocabulary benchmarks for natural images [55], we train the model on one dataset and evaluate its cross-dataset performance on other datasets. We select two datasets for training, *i.e.*, iSAID [64] and SIOR [48]. iSAID is a large scale instance segmentation dataset for remote sensing images. It contains 18732 images for training across 15 categories. SIOR is developed from aerial object detection dataset DIOR [25], with segmentation annotations generated in SAMRS [48], which contains 11725 images with 20 categories.

**Evaluation Dataset.** To evaluate the effectiveness of our method, we conduct cross-dataset evaluation on 4 aerial instance segmentation datasets, *i.e.*, NWPU-VHR-10 [7, 44], SOTA [48], FAST [48], and SIOR [48]. NWPU-VHR-10 is an aerial object detection dataset with instance masks further annotated by [44]. The test set contains 731 images from 10 aerial classes. SOTA, FAST and SIOR are segmentation datasets provided in SAMRS [48], which are developed from aerial object detection datasets DOTA-V2.0 [11], FAIR1M-2.0 [47], and DIOR [25], respectively. SOTA covers 874 images with 18 object categories for testing, FAST contains 3207 images across 37 fine-grained aerial object categories for testing, and SIOR is with 11738 testing samples. We provide the categories in each dataset in Table B.

#### A.2. Remote Sensing Semantic Segmentation

**Training Dataset.** Following the open-vocabulary benchmarks for remote sensing semantic segmentation [60], we train the model on their proposed LandDiscover50K dataset. It includes 51846 high-resolution remote sensing images annotated across 40 object categories.

**Evaluation Dataset.** To evaluate the effectiveness of our method, we follow the evaluation settings in [60] to conduct cross-dataset evaluation on 4 remote sensing semantic

datasets, *i.e.*, FLAIR [13], FAST [48], Potsdam [20], and FloodNet [41]. Each dataset has its own bias towards different remote sensing categories. To illustrate, Potsdam [20] emphasizes the in-vocabulary performance with high category similarity to the training LandDiscover50K dataset, which contains 5472 images with 6 semantic categories. FloodNet [41] focuses more on the post-flood analysis, which contains 898 images with 9 semantic categories. FLAIR [13] is with 15700 images focusing on 12 large-scale landcover types. FAST [48] contains 3207 images, specializing in 37 fine-grained semantic classes for remote sensing. The combination of the four datasets enables a comprehensive evaluation of the open-vocabulary semantic segmentation tasks in remote sensing. We provide the categories in each dataset in Table C.

### B. Additional Experiment Results

#### B.1. Additional Results on Semantic Segmentation

The proposed **SCORE** can also be applied to diverse segmentation related tasks, *e.g.* semantic segmentation. We provide the semantic segmentation results of our method in Table A. Our approach consistently outperforms existing across three of four benchmarks, demonstrating its effectiveness in open-vocabulary remote sensing semantic segmentation. Specifically, we achieve an average improvement of 1.13% over the current SOTA model [61]. Our method surpasses previous methods by a large margin, especially on FLAIR and FAST datasets, with gains up to 9.62%.

Method	LandDiscover50K				
	FLAIR	FAST	Potsdam	FloodNet	Average
CAT-SEG [10] [CVPR24]	19.71	15.55	39.57	35.91	27.69
GSNet [61] [AAAI25]	18.35	15.21	<b>43.29</b>	37.68	28.63
<b>SCORE (Ours)</b>	<b>29.33</b>	<b>21.51</b>	26.51	<b>41.70</b>	<b>29.76</b>

Table A. Comparison with SOTA methods on open-vocabulary remote sensing semantic segmentation. The model is trained on LandDiscover50K dataset and then tested on the four evaluation benchmarks to measure its cross-dataset generalization capabilities.

#### B.2. Additional Qualitative Results

We provide additional qualitative results of our proposed method on remote sensing instance segmentation task as shown in Figure A.

Dataset	#Category	Category Name
iSAID [64]	15	ship, storage tank, baseball diamond, tennis court, basketball court, ground track field, bridge, large vehicle, small vehicle, helicopter, swimming pool, roundabout, soccer ball field, plane, harbor
SIOR [48]	20	airplane, airport, baseball field, basketball court, bridge, chimney, expressway service area, expressway toll station, dam, golffield, ground track field, harbor, overpass, ship, stadium, storage tank, tennis court, train station, vehicle, windmill
NWPU [7, 44]	10	airplane, ship, storage tank, baseball diamond, tennis court, basketball court, ground track field, harbor, bridge, vehicle
FAST [48]	37	A220, A321, A330, A350, ARJ21, baseball field, basketball court, Boeing737, Boeing747, Boeing777, Boeing787, bridge, bus, C919, cargo truck, dry cargo ship, dump truck, engineering ship, excavator, fishing boat, football field, intersection, liquid cargo ship, motorboat, other-airplane, other-ship, other-vehicle, passenger ship, roundabout, small car, tennis court, tractor, trailer, truck tractor, tugboat, van, warship
SOTA [48]	18	large vehicle, swimming pool, helicopter, bridge, plane, ship, soccer ball field, basketball court, ground track field, small vehicle, baseball diamond, tennis court, roundabout, storage tank, harbor, container crane, airport, helipad

Table B. Category Names for datasets used in our instance segmentation benchmarks.

Dataset	#Category	Category Name
LandDiscover50K [61]	40	background, bare land, grass, pavement, road, tree, water, agriculture land, buildings, forest land, barren land, urban land, large vehicle, swimming pool, helicopter, bridge, plane, ship, soccer ball field, basketball court, ground track field, small vehicle, baseball diamond, tennis court, roundabout, storage tank, harbor, container crane, airport, helipad, chimney, expressway service area, expresswalltoll station, dam, golf field, overpass, stadium, train station, vehicle, windmill
FLAIR [13]	12	building, pervious surface, impervious surface, bare soil, water, coniferous, deciduous, brushwood, vineyard, herbaceous vegetation, agricultural land, plowed land
FAST [48]	37	A220, A321, A330, A350, ARJ21, baseball field, basketball court, Boeing737, Boeing747, Boeing777, Boeing787, bridge, bus, C919, cargo truck, dry cargo ship, dump truck, engineering ship, excavator, fishing boat, football field, intersection, liquid cargo ship, motorboat, other-airplane, other-ship, other-vehicle, passenger ship, roundabout, small car, tennis court, tractor, trailer, truck tractor, tugboat, van, warship
Potsdam [20]	6	impervious surface, building, low vegetation, tree, car, clutter
FloodNet [41]	9	building-flooded, building-non-flooded, road-flooded, road-non-flooded, water, tree, vehicle, pool, grass

Table C. Category Names for datasets used in our semantic segmentation benchmarks.

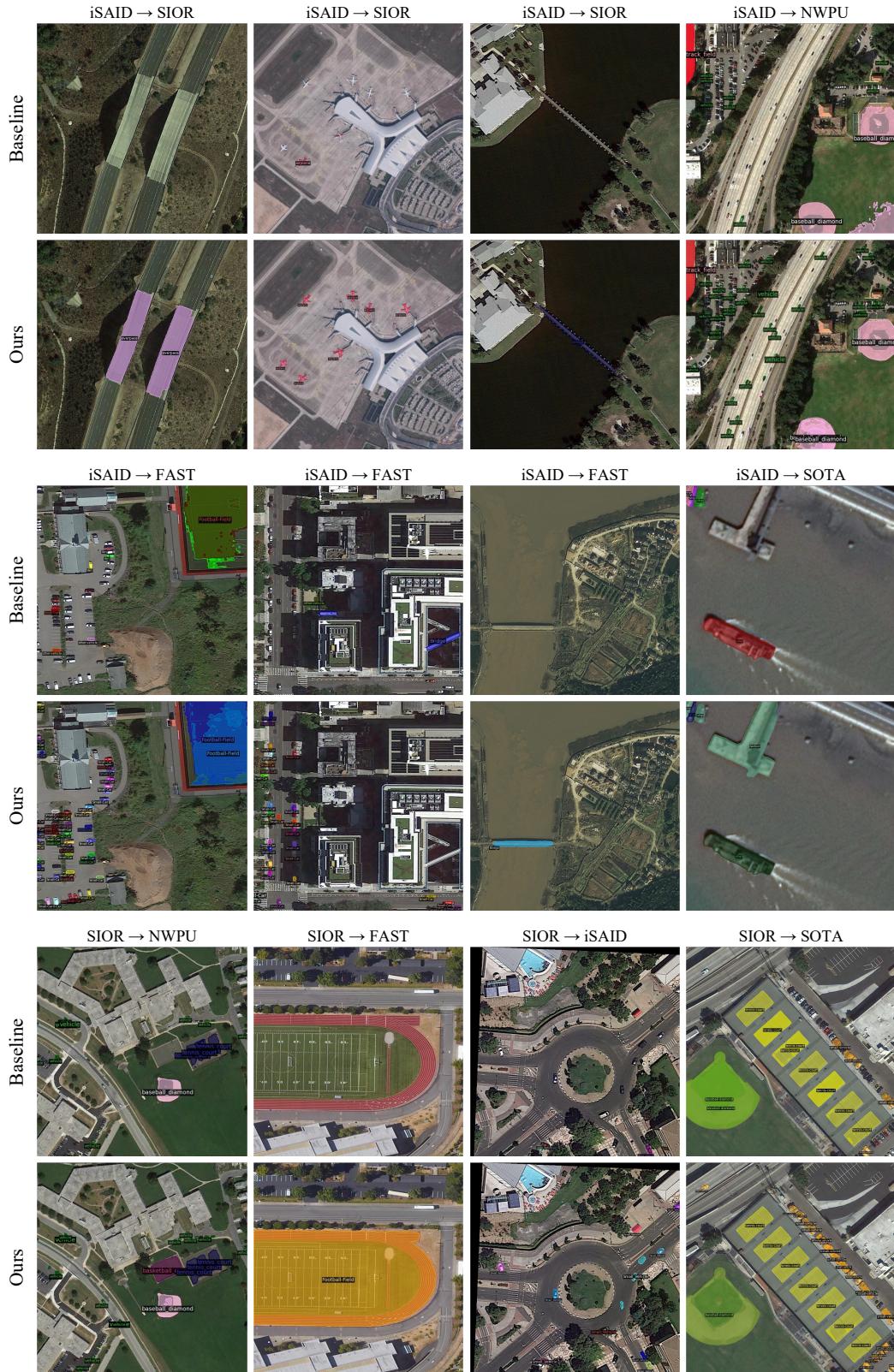


Figure A. Additional qualitative results between the baseline and our model.