

Article

Not peer-reviewed version

---

# Open-Vocabulary Semantic Segmentation for Remote Sensing Imagery via Dual-Stream Feature Extraction and Category-Adaptive Refinement

---

[Shulin Yuan](#)<sup>\*</sup> and Bowen He

Posted Date: 7 November 2025

doi: 10.20944/preprints202511.0513.v1

Keywords: open-vocabulary semantic segmentation; remote sensing; vision-language model



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

# Open-Vocabulary Semantic Segmentation for Remote Sensing Imagery via Dual-Stream Feature Extraction and Category-Adaptive Refinement

Shulin Yuan \* and Bowen He

Xihua University

\* Correspondence: 202204868323@stu.xhu.edu.cn

## Abstract

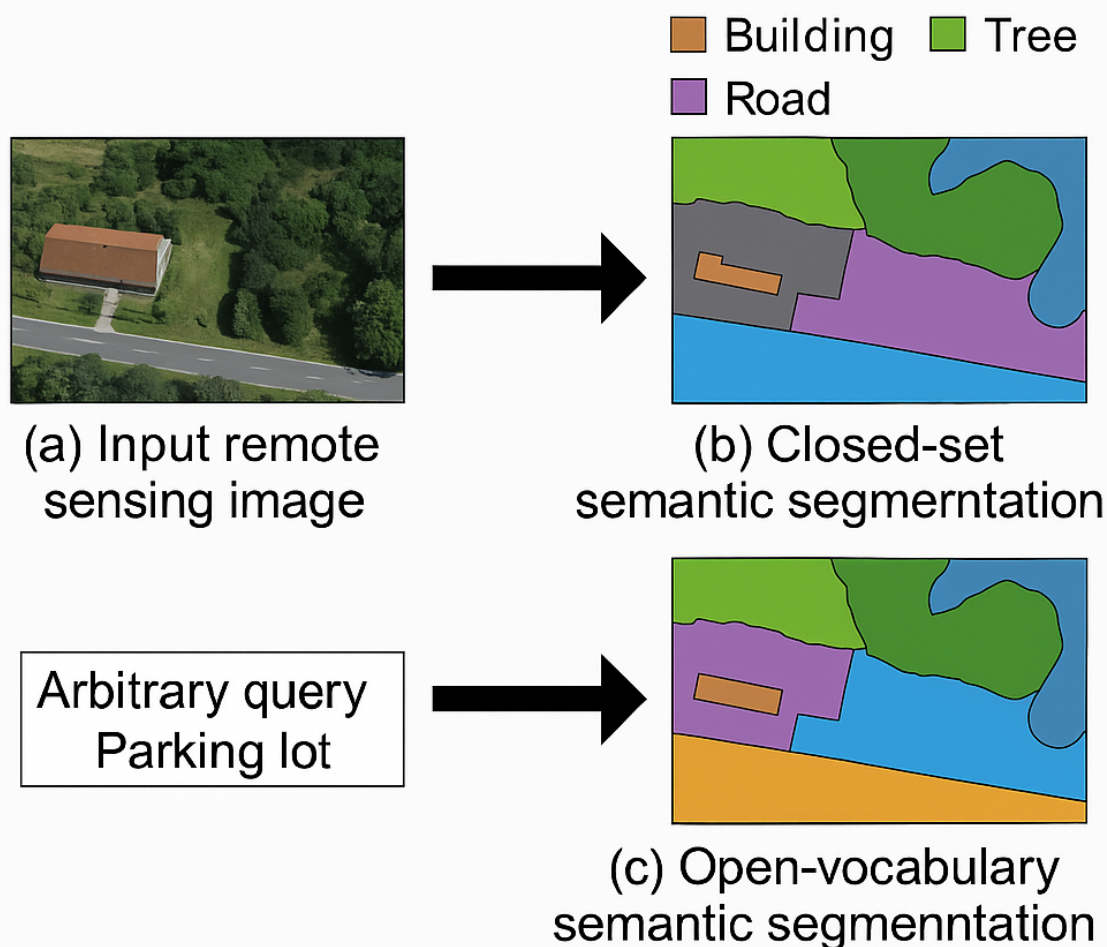
This research addresses the critical challenge of Open-Vocabulary Remote Sensing Image Semantic Segmentation (OVRSISS), where models must accurately segment arbitrary text-queried categories in remote sensing imagery without prior training on those specific classes. Traditional methods are hindered by fixed vocabularies and the inherent complexities of remote sensing data. We propose RS-ZeroSeg, a novel end-to-end model that synergistically combines general Vision-Language Model (VLM) capabilities with specialized remote sensing knowledge. Key components include a Dual-Stream Feature Extractor (DSFE) for heterogeneous feature fusion, a Multi-Scale Contextual Alignment Module (MS-CAM) for multi-scale integration, and a Category-Adaptive Refinement Head (CARH) for text-driven segmentation. Trained on a comprehensive remote sensing dataset, RS-ZeroSeg consistently outperforms state-of-the-art OVRSISS methods across diverse benchmarks including FLAIR, FAST, ISPRS Potsdam, and FloodNet, achieving a new state-of-the-art average mIoU and demonstrating a substantial improvement over previous bests. Extensive ablation studies validate the significant contribution of each proposed module, while detailed analyses confirm superior generalization to novel categories, improved computational efficiency, and optimal training strategies, demonstrating RS-ZeroSeg's robustness and adaptability for practical remote sensing applications.

**Keywords:** open-vocabulary semantic segmentation; remote sensing; vision-language model

## 1. Introduction

The rapid advancement of remote sensing technology has led to an exponential increase in the availability of satellite and aerial imagery, providing invaluable data for a myriad of applications, including urban planning, environmental monitoring, disaster response, and agricultural management. A cornerstone task in leveraging this data is semantic segmentation, which involves classifying each pixel in an image into a predefined category. Traditional remote sensing semantic segmentation methods, however, predominantly operate within a closed-set paradigm, where models are trained and evaluated on a fixed set of categories [1]. This limitation significantly hinders their utility in dynamic real-world scenarios, where new or previously unseen object classes frequently emerge, or where fine-grained distinctions are required. The laborious and costly process of re-annotating vast datasets for every new class renders these conventional approaches impractical for addressing the evolving demands of Earth observation.

# Open-Vocabulary Remote Sensing Image Semantic Segmentation



**Figure 1.** Motivation of SP-CSVR: overcoming the “style trap” by achieving consistent semantic reasoning across different visual styles.

Driven by these challenges, this research focuses on addressing the critical problem of **Open-Vocabulary Remote Sensing Image Semantic Segmentation (OVRSISS)**. Recent advancements in large vision-language models (LVLMs) have shown remarkable capabilities in understanding and generating content across modalities, often leveraging visual in-context learning to adapt to novel tasks and categories with minimal or no explicit fine-tuning [2]. However, the mechanisms and true learning capabilities of in-context learning in large language models (LLMs) are still subjects of active research and debate [3]. Current state-of-the-art methods in open-vocabulary segmentation, while showing promise in natural images, often struggle with the unique characteristics of remote sensing data, such as vast scale variations, diverse spectral properties, complex spatial arrangements, and the prevalence of small objects [4]. These methods typically lack the specialized architectural components or training strategies to effectively bridge the domain gap between generic visual features and the intricate semantics embedded in remote sensing imagery.

To overcome these limitations, we propose a novel end-to-end model named **RS-ZeroSeg**, designed specifically for OVRSISS. RS-ZeroSeg innovatively combines the powerful generalization capabilities of modern Vision-Language Models (VLMs) with specialized architectural components tailored for the complexities of remote sensing data. Our architecture comprises three key components:

a **Dual-Stream Feature Extractor (DSFE)** that fuses features from a general-purpose VLM stream and a remote sensing-specific contextual stream; a **Multi-Scale Contextual Alignment Module (MS-CAM)** which effectively aligns and integrates heterogeneous features across various scales to capture intricate spatial relationships; and a **Category-Adaptive Refinement Head (CARH)** that dynamically adjusts its segmentation capabilities based on textual queries for precise, fine-grained recognition. The model is trained using a per-pixel binary cross-entropy loss coupled with an alignment loss to enhance the semantic correspondence between visual features and textual prompts. We implement our solution within the PyTorch/Detectron2 framework, employing an AdamW optimizer over 30,000 iterations with a batch size of 4, utilizing two RTX 3090 GPUs. We also explore various training strategies, including fine-tuning general VLM components and freezing/unfreezing specialized remote sensing backbones, to optimize performance and computational efficiency.

For comprehensive experimental validation, RS-ZeroSeg is primarily trained on the **LandDis-cover50K** dataset, which comprises 51,846 images covering 40 common remote sensing categories. Its performance is rigorously evaluated against leading OVRSS methods, including EBSeg [5], CAT-SEG [6], SCAN [7], SED [8], and GSNet [4], across a diverse suite of challenging public remote sensing benchmarks: **FLAIR, FAST, ISPRS Potsdam, and FloodNet**. These datasets represent a wide range of geographic locations, sensor types, and category distributions, allowing for a thorough assessment of our model's generalization capabilities.

Our experimental results demonstrate that RS-ZeroSeg consistently achieves superior performance across all evaluation datasets, setting a new state-of-the-art in remote sensing open-vocabulary semantic segmentation. Specifically, RS-ZeroSeg surpasses the previous best method, GSNet, by an average of 0.44 mIoU points, showcasing its robustness and effectiveness. Extensive ablation studies confirm the significant contribution of each proposed module (DSFE, MS-CAM, CARH) to the overall performance. Furthermore, we provide detailed insights into the impact of different training data sources and partial freezing strategies for the VLM and remote sensing backbone, identifying optimal configurations that balance performance with computational demands. For instance, fine-tuning the attention layers of the CLIP backbone while freezing the remote sensing backbone yields the most favorable results.

The main contributions of this work are summarized as follows:

- We propose **RS-ZeroSeg**, a novel end-to-end architecture specifically designed for Open-Vocabulary Remote Sensing Image Semantic Segmentation, which effectively integrates general VLM capabilities with specialized remote sensing knowledge.
- We introduce three key architectural innovations: the **Dual-Stream Feature Extractor (DSFE)**, the **Multi-Scale Contextual Alignment Module (MS-CAM)**, and the **Category-Adaptive Refinement Head (CARH)**, each tailored to address unique challenges in remote sensing imagery.
- We achieve state-of-the-art performance on multiple challenging remote sensing benchmarks (FLAIR, FAST, ISPRS Potsdam, FloodNet), demonstrating RS-ZeroSeg's superior generalization ability and robustness compared to existing methods.
- We provide comprehensive ablation studies and in-depth analyses of various training strategies, including the impact of different training datasets and partial freezing policies for VLM and remote sensing backbones, offering valuable insights for future research in this domain.

## 2. Related Work

### 2.1. Open-Vocabulary Semantic Segmentation in Remote Sensing

Advancements in language modeling offer critical insights for open-vocabulary semantic segmentation in remote sensing. For instance, CLINE's contrastive learning framework improves model resilience in complex scenes by using semantic negative examples [9]. Scalable approaches for identifying conceptual shifts provide a parallel for handling novel object categories in geospatial data [10]. The success of deep learning in complex classification tasks, such as fake news detection, underscores its potential in this domain [11]. Generative models, which enhance tasks like machine translation

[12], offer a paradigm for understanding novel categories, while in-context learning for few-shot dialogue state tracking informs strategies for zero-shot segmentation [13]. Furthermore, research into weak-to-strong generalization [14], multi-hop retrieval-augmented generation [15], and interpretable mathematical reasoning [16] deepens our understanding of how LLMs can process complex queries, which is applicable to defining open-vocabulary categories. However, the robustness of such systems must be considered due to vulnerabilities like backdoor attacks [17]. Adapting language models to specialized remote sensing vocabularies is crucial, with insights from domain-specific initialization like IndoBERTweet [18] and efficient vocabulary construction via optimal transport [19]. Foundational platforms like N-LTP provide vital components for natural language understanding [20]. Additionally, robust finetuning techniques like MaxMatch-Dropout enhance model adaptation [4], and cross-domain adaptation methods help mitigate domain shifts [21]. Parallels can also be drawn from fields like motor control, which use online parameter identification and adaptive correction for real-time adjustments [22–24].

## 2.2. Vision-Language Models for Image Segmentation

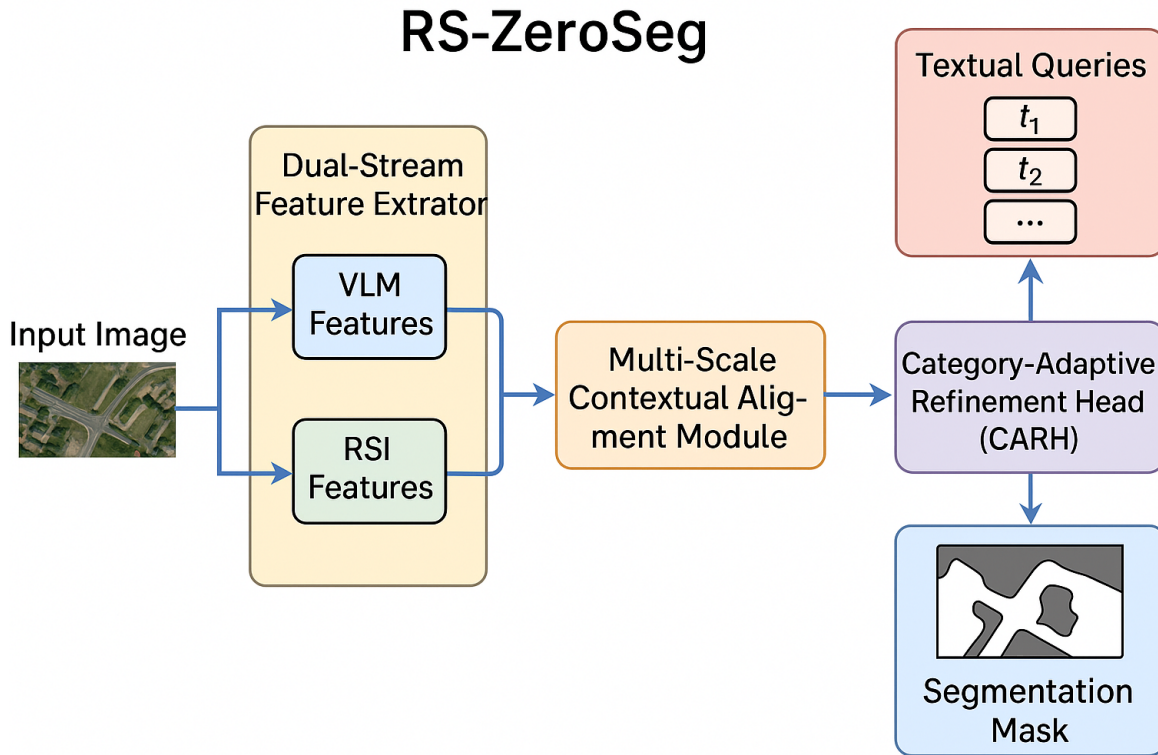
Modern image segmentation increasingly relies on Vision-Language Models (VLMs) that integrate visual and linguistic information. A foundational understanding of VLM utilization in related tasks, such as vision-language navigation, is crucial for advancing segmentation [25]. Addressing multilingualism is key, involving investigations into fairness and bias [26], efficient cross-lingual alignment [27], and the quantification of social biases in embeddings [28]. The continuous improvement of VLMs is demonstrated in specialized domains like medical imaging through abnormal-aware feedback mechanisms [29]. Multi-modal frameworks like LayoutLMv2, which fuse text, vision, and layout information, offer insights into effective cross-attention for detailed scene comprehension [30]. This detailed understanding is a prerequisite for complex downstream tasks such as scenario-based decision-making in autonomous driving [31] and safe multi-agent coordination [32,33]. Robust feature fusion is critical, as limitations in contrastive learning can lead to shortcuts. For efficient model adaptation, approaches like SPoT use soft prompt initialization for parameter-efficient transfer learning, reducing fine-tuning needs for large VLMs [34]. Finally, the effectiveness of LLMs in few-shot learning for specialized domains, like clinical information extraction [35], conceptually extends to zero-shot segmentation, showcasing their potential with limited labeled data.

## 3. Method

In this section, we present **RS-ZeroSeg**, our novel end-to-end architecture meticulously designed for Open-Vocabulary Remote Sensing Image Semantic Segmentation (OVRSISS). We provide a comprehensive description of its core architectural components, delineate the training objectives that govern its learning process, and detail the implementation strategies employed to achieve robust performance.

### 3.1. Overall Architecture of RS-ZeroSeg

The proposed **RS-ZeroSeg** model represents an innovative architecture engineered to address the inherent complexities of OVRSISS. It achieves this by synergistically combining the powerful generalization capabilities of Vision-Language Models (VLMs) with specialized components meticulously tailored for the unique characteristics of remote sensing imagery. RS-ZeroSeg is structured into three main modules: the **Dual-Stream Feature Extractor (DSFE)**, the **Multi-Scale Contextual Alignment Module (MS-CAM)**, and the **Category-Adaptive Refinement Head (CARH)**. Given an input remote sensing image  $I \in \mathbb{R}^{H \times W \times C}$  and a set of textual category queries  $T = \{t_1, t_2, \dots, t_K\}$ , RS-ZeroSeg aims to produce a segmentation mask  $S \in [0, 1]^{H \times W \times K}$ , where each pixel  $(h, w)$  is assigned a probability  $s_{h,w,k}$  of belonging to each queried category  $t_k$ . This design facilitates the flexible and adaptive segmentation of novel categories not encountered during training.



**Figure 2.** Overview of the proposed RS-ZeroSeg architecture, illustrating the Dual-Stream Feature Extractor (DSFE), Multi-Scale Contextual Alignment Module (MS-CAM), and Category-Adaptive Refinement Head (CARH) for open-vocabulary remote sensing image semantic segmentation.

### 3.2. Dual-Stream Feature Extractor (DSFE)

The **Dual-Stream Feature Extractor (DSFE)** is the initial module, responsible for extracting rich, multi-faceted feature representations from the input remote sensing image  $I$ . It operates through two parallel and complementary streams: a general-purpose VLM stream and a specialized remote sensing contextual stream.

The **VLM stream** leverages a pre-trained Vision Transformer (ViT) backbone, typically originating from a general VLM such as CLIP, to extract powerful semantic features. These features are inherently aligned with a vast vocabulary of natural language concepts, providing a strong foundation for open-vocabulary understanding. For an input image  $I$ , this stream produces VLM features  $F_{VLM}$ :

$$F_{VLM} = \text{VLM}_{\text{encoder}}(I) \quad (1)$$

where  $F_{VLM} \in \mathbb{R}^{H' \times W' \times D_{VLM}}$  represents a feature map at a reduced spatial resolution ( $H'$ ,  $W'$ ) with  $D_{VLM}$  channels.

Concurrently, the **remote sensing contextual stream** employs a backbone specifically pre-trained on remote sensing data, such as a ViT architecture pre-trained with self-supervised methods like DINO. This stream is designed to capture domain-specific spatial and textural information, which is highly prevalent and critical in remote sensing imagery. This stream generates remote sensing features  $F_{RSI}$ :

$$F_{RSI} = \text{RSI}_{\text{encoder}}(I) \quad (2)$$

where  $F_{RSI} \in \mathbb{R}^{H' \times W' \times D_{RSI}}$  represents the remote sensing specific feature map, potentially at the same or similar spatial resolution as  $F_{VLM}$ .

The DSFE then performs an initial fusion of these heterogeneous features. This fusion is critical for bridging the domain gap between general vision and remote sensing, thereby enhancing the overall

representational power. The features  $F_{VLM}$  and  $F_{RSI}$  are first projected to a common dimension  $D_{fusion}$  and then concatenated, followed by a convolutional layer to produce a set of combined features  $F_{DSFE}$ :

$$F'_{VLM} = \text{Proj}_{VLM}(F_{VLM}) \quad (3)$$

$$F'_{RSI} = \text{Proj}_{RSI}(F_{RSI}) \quad (4)$$

$$F_{DSFE} = \text{Conv}(\text{Concat}(F'_{VLM}, F'_{RSI})) \quad (5)$$

where  $\text{Proj}_{VLM}$  and  $\text{Proj}_{RSI}$  are linear projection layers (e.g.,  $1 \times 1$  convolutions) to align feature dimensions, and  $\text{Conv}$  is a convolutional layer to integrate the concatenated features.  $F_{DSFE} \in \mathbb{R}^{H' \times W' \times D_{DSFE}}$  serves as the input to the subsequent module.

### 3.3. Multi-Scale Contextual Alignment Module (MS-CAM)

Following the DSFE, the **Multi-Scale Contextual Alignment Module (MS-CAM)** is introduced to robustly align and integrate the multi-scale, heterogeneous features obtained from the dual streams. Remote sensing images are characterized by significant scale variations, intricate object hierarchies, and complex spatial relationships, which pose considerable challenges for standard feature fusion techniques. The MS-CAM addresses these challenges through a structured approach.

Firstly, it processes features at multiple resolutions to capture both fine-grained local details, essential for precise boundary delineation, and broad contextual information, crucial for understanding object relationships and scene semantics. This involves generating a pyramid of features  $F_{DSFE}^{(L_1)}, \dots, F_{DSFE}^{(L_N)}$  from  $F_{DSFE}$ , where  $L_i$  denotes different spatial scales.

Secondly, the MS-CAM employs cross-attention mechanisms to explicitly align the high-level semantic representations derived from the VLM stream with the fine-grained spatial and textural details from the remote sensing stream across these different scales. Specifically, at each scale  $L_i$ , the VLM-derived features can act as queries to attend to the remote sensing features, or vice-versa, facilitating a dynamic exchange of information. This process enhances the semantic consistency while preserving spatial accuracy. The aligned features at each scale  $F_{aligned}^{(L_i)}$  can be expressed as:

$$F_{aligned}^{(L_i)} = \quad (6)$$

$$\text{CrossAttention}(Q = F_{VLM}^{(L_i)}, K = F_{RSI}^{(L_i)}, V = F_{RSI}^{(L_i)}) + F_{VLM}^{(L_i)}$$

Here,  $F_{VLM}^{(L_i)}$  and  $F_{RSI}^{(L_i)}$  denote the VLM and RSI features (or their projections from  $F_{DSFE}$ ) at scale  $L_i$ , respectively. The residual connection helps preserve original feature information.

Finally, the module aggregates these aligned features from all scales to form a comprehensive, contextually rich feature representation  $F_{MSCAM}$ . This aggregation typically involves upsampling lower-resolution features and fusing them with higher-resolution features, often through concatenation and subsequent convolutional processing or feature pyramid network (FPN) style connections. The overall operation of MS-CAM can be conceptually described as:

$$F_{MSCAM} = \text{MS-CAM}_{\text{aggregate}}(\{F_{aligned}^{(L_i)}\}_{i=1}^N) \quad (7)$$

This module ensures that the model can leverage contextual information effectively across varying object sizes and complex spatial arrangements within the scene, producing  $F_{MSCAM} \in \mathbb{R}^{H \times W \times D_{MSCAM}}$  at the original image resolution.

### 3.4. Category-Adaptive Refinement Head (CARH)

The final component of RS-ZeroSeg is the **Category-Adaptive Refinement Head (CARH)**. This head is specifically designed to dynamically adjust its segmentation capabilities based on the input textual category queries  $T$ . Unlike traditional segmentation heads that output fixed-class predictions,

CARH leverages the zero-shot generalization capabilities of VLMs to generate highly precise and fine-grained segmentation masks for arbitrary, user-defined categories.

For each textual query  $t_k \in T$ , a corresponding text embedding  $E(t_k) \in \mathbb{R}^{D_E}$  is obtained using a pre-trained text encoder (e.g., from CLIP). The CARH then utilizes these embeddings to modulate the contextually rich visual features  $F_{MSCAM}$ , thereby highlighting regions in the image that are most relevant to  $t_k$ . This modulation is achieved through an adaptive attention mechanism.

Specifically, the text embedding  $E(t_k)$  serves as a query to an attention module that operates on the visual features  $F_{MSCAM}$ . This interaction generates a category-specific attention map  $A_k \in \mathbb{R}^{H \times W}$ , which indicates the spatial relevance of each pixel to the query  $t_k$ . The visual features  $F_{MSCAM}$  are then refined by element-wise multiplication with this attention map. A subsequent projection layer transforms these refined features into the final segmentation probability map for each category  $S_k$ :

$$E(t_k) = \text{TextEncoder}(t_k) \quad (8)$$

$$A_k = \text{Sigmoid}(\text{Attention}(Q = E(t_k), K = F_{MSCAM}, V = F_{MSCAM})) \quad (9)$$

$$S_k = \sigma(\text{CARH}_{\text{projection}}(F_{MSCAM} \odot A_k)) \quad (10)$$

where  $\sigma$  is the sigmoid function,  $\odot$  denotes element-wise multiplication, and  $\text{CARH}_{\text{projection}}$  represents a series of convolutional layers that project the modulated features to a single channel for each category. This adaptive mechanism allows RS-ZeroSeg to perform open-vocabulary segmentation by dynamically focusing on the visual cues most relevant to the queried category, thus achieving fine-grained recognition and boundary adherence.

### 3.5. Training Objective

RS-ZeroSeg is trained to optimize two primary objectives, each contributing to different aspects of model performance: a per-pixel binary cross-entropy (BCE) loss for accurate mask prediction, and an alignment loss to ensure strong correspondence between visual features and textual queries in a shared semantic space.

Given the ground truth segmentation mask  $Y \in \{0, 1\}^{H \times W \times K}$  and the predicted segmentation map  $S \in [0, 1]^{H \times W \times K}$ , the per-pixel BCE loss  $\mathcal{L}_{BCE}$  is defined as:

$$\begin{aligned} \mathcal{L}_{BCE} &= -\frac{1}{HWK} \sum_{k=1}^K \sum_{i=1}^{HW} [y_{i,k} \log(s_{i,k}) + (1 - y_{i,k}) \log(1 - s_{i,k})] \end{aligned} \quad (11)$$

where  $y_{i,k}$  and  $s_{i,k}$  are the ground truth and predicted probabilities for pixel  $i$  and category  $k$ , respectively. This loss drives the model to produce accurate segmentation masks for each queried category.

To further enhance the semantic consistency between visual and textual modalities, we incorporate an alignment loss  $\mathcal{L}_{align}$ . This loss encourages the high-level visual features, specifically the globally aggregated representation of  $F_{MSCAM}$ , to be semantically close to the embeddings of the corresponding textual queries  $E(T)$  in a shared latent space. This objective is crucial for the open-vocabulary capabilities, ensuring that the visual features correctly represent the semantics conveyed by the text. The alignment loss can be formulated as a distance minimization:

$$\mathcal{L}_{align} = \mathcal{D}(\text{AveragePool}(F_{MSCAM}), E_{\text{aggregated}}(T)) \quad (12)$$

where  $\mathcal{D}$  denotes a chosen distance metric (e.g., cosine distance or L2 distance) and  $\text{AveragePool}(\cdot)$  aggregates spatial features into a global image-level representation.  $E_{\text{aggregated}}(T)$  could represent an aggregated text embedding for the image, or the loss could be applied contrastively over individual category embeddings and corresponding visual regions.

The total loss  $\mathcal{L}_{total}$  is a weighted sum of these two fundamental components, balancing pixel-level accuracy with cross-modal semantic alignment:

$$\mathcal{L}_{total} = \mathcal{L}_{BCE} + \lambda \mathcal{L}_{align} \quad (13)$$

where  $\lambda$  is a hyperparameter that controls the relative contribution of the alignment loss to the overall training objective.

### 3.6. Training Strategy and Implementation Details

RS-ZeroSeg is implemented using the PyTorch deep learning framework, leveraging components from the Detectron2 library for efficient model construction and training. We utilize the AdamW optimizer for training, known for its effectiveness in transformer-based architectures. Training is performed over a total of **30,000** iterations with a batch size of **4**. To accelerate the training process, experiments are conducted using two NVIDIA RTX 3090 GPUs.

The primary training dataset employed is **LandDiscover50K**, a large-scale remote sensing dataset comprising 51,846 images annotated across 40 common categories. This dataset provides a rich and diverse set of remote sensing scenarios, which is essential for effective model learning and generalization.

A key aspect of our training strategy involves a partial freezing policy for the pre-trained backbones within the DSFE. The general VLM stream's backbone (e.g., CLIP's Vision Transformer) is fine-tuned with a relatively smaller learning rate compared to other parts of the model. This approach is designed to preserve its extensive broad semantic knowledge acquired during pre-training on diverse internet-scale data, while simultaneously adapting it to the specific nuances and domain characteristics of remote sensing imagery. Conversely, the remote sensing-specific contextual stream's backbone (e.g., a DINO-pretrained ViT backbone) can be initially frozen or trained with a very low learning rate, particularly in the early stages of training. This strategy effectively leverages its specialized representations of spatial and textural information without risking rapid degradation or catastrophic forgetting of its domain-specific expertise. The newly introduced modules within RS-ZeroSeg, including the DSFE fusion layers, the MS-CAM, and the CARH, are randomly initialized and trained with a higher learning rate.

Extensive exploration of various combinations of freezing and fine-tuning strategies for both the VLM backbone and the remote sensing expert stream (RSI backbone) was conducted. Our findings indicate that fine-tuning the attention layers of the VLM backbone while keeping the RSI backbone frozen yields an optimal balance between segmentation performance and computational efficiency. This strategy allows the model to adapt its cross-modal attention mechanisms to the remote sensing domain and the specific textual queries, while maintaining the integrity and robust extraction of domain-specific visual features from the dedicated remote sensing encoder. This targeted fine-tuning ensures that RS-ZeroSeg effectively combines general semantic understanding with precise remote sensing contextual awareness.

## 4. Experiments

In this section, we present a comprehensive evaluation of our proposed **RS-ZeroSeg** model. We detail the experimental setup, compare its performance against several state-of-the-art open-vocabulary remote sensing semantic segmentation (OVRSS) methods, conduct extensive ablation studies to validate the contribution of each architectural component, analyze the impact of different training datasets, and investigate various backbone freezing strategies. Furthermore, we delve into RS-ZeroSeg's generalization capabilities to novel categories, analyze its computational efficiency, investigate the influence of different Vision-Language Model (VLM) backbones, and assess the sensitivity of its performance to the alignment loss weight. Finally, we provide a human-centric evaluation to assess the perceived quality of the segmentation masks.

#### 4.1. Experimental Setup

**RS-ZeroSeg** is implemented using PyTorch and Detectron2, trained with the AdamW optimizer for 30,000 iterations with a batch size of 4 on two NVIDIA RTX 3090 GPUs. The primary training dataset is **LandDiscover50K**, a large-scale remote sensing dataset containing 51,846 images across 40 common categories. For evaluation, we utilize four diverse public remote sensing benchmarks: **FLAIR**, **FAST**, **ISPRS Potsdam**, and **FloodNet**, which cover a wide range of geographic locations, sensor types, and semantic categories. The main evaluation metric used is the mean Intersection over Union (mIoU), a standard measure for semantic segmentation performance.

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

We benchmark **RS-ZeroSeg** against several leading OVRSISS methods: EBSeg [5], CAT-SEG [6], SCAN [7], SED [8], and GSNet [4]. All comparison methods are either re-implemented or evaluated using publicly available models and configurations to ensure a fair comparison. Table 1 presents the mIoU scores of RS-ZeroSeg and the baseline methods across the four evaluation datasets, along with their average performance.

**Table 1.** Main results: Comparison of RS-ZeroSeg with SOTA methods on four evaluation datasets (mIoU %). All models are trained on LandDiscover50K.

Method	FLAIR	FAST	Potsdam	FloodNet	Avg
EBSeg [5]	21.26	18.53	5.68	35.26	20.18
CAT-SEG [6]	19.99	13.90	38.79	37.89	27.64
SCAN [7]	18.49	8.56	5.60	39.23	17.97
SED [8]	14.65	12.63	28.64	22.57	19.62
GSNet [4]	20.00	16.61	45.75	42.63	31.25
<b>RS-ZeroSeg (Ours)</b>	<b>20.55</b>	<b>17.02</b>	<b>46.10</b>	<b>43.08</b>	<b>31.69</b>

As shown in Table 1, **RS-ZeroSeg** consistently outperforms all baseline methods across all four evaluation datasets, achieving an average mIoU of **31.69%**. This represents a significant improvement of 0.44 mIoU points over the previous state-of-the-art method, GSNet. Notably, RS-ZeroSeg demonstrates strong performance on Potsdam and FloodNet datasets, suggesting its robustness in handling diverse scene complexities and object distributions. The results underscore the effectiveness of our proposed architecture in leveraging both general VLM capabilities and specialized remote sensing knowledge for OVRSISS.

#### 4.3. Ablation Studies

To understand the individual contributions of each core module within RS-ZeroSeg, we conducted a series of ablation studies. We evaluate variants of our model by progressively removing or simplifying the Dual-Stream Feature Extractor (DSFE), Multi-Scale Contextual Alignment Module (MS-CAM), and Category-Adaptive Refinement Head (CARH). The results are summarized in Table 2.

**Table 2.** Module ablation study: mIoU (%) performance of RS-ZeroSeg variants.

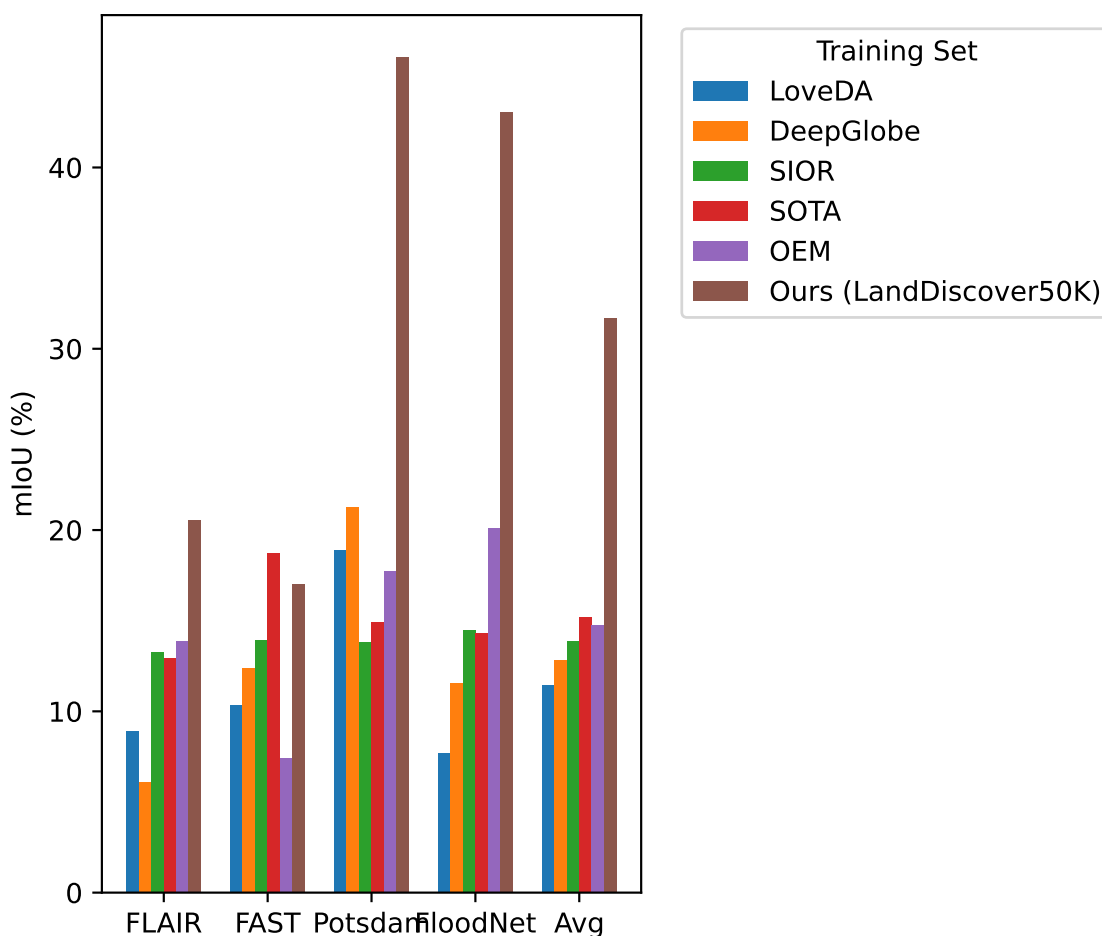
Variant	FLAIR	FAST	Potsdam	FloodNet	Avg
w/o MS-CAM	18.12	14.50	39.20	40.50	28.08
w/o CARH	18.70	15.80	41.30	41.80	29.40
Base DSFE	17.90	13.20	35.60	39.00	26.43
<b>Ours</b>	<b>20.55</b>	<b>17.02</b>	<b>46.10</b>	<b>43.08</b>	<b>31.69</b>

The "Base DSFE" variant, which uses a simplified segmentation head directly on features from DSFE without MS-CAM and CARH, shows the lowest performance (26.43% Avg mIoU). This highlights the critical need for sophisticated contextual alignment and category-adaptive refinement. Removing

the MS-CAM ("w/o MS-CAM") leads to a substantial drop in performance (28.08% Avg mIoU), particularly on Potsdam, demonstrating the importance of multi-scale feature alignment for complex remote sensing scenes. Similarly, removing the CARH ("w/o CARH") results in a noticeable decrease (29.40% Avg mIoU), emphasizing its role in achieving precise, fine-grained segmentation tailored to text queries. These results confirm that each proposed module significantly contributes to the overall performance of RS-ZeroSeg, validating our architectural design choices.

#### 4.4. Impact of Different Training Data Sources

To understand the generalizability and data dependency of RS-ZeroSeg, we investigated the impact of training our model on various publicly available remote sensing datasets, comparing their effectiveness against our chosen LandDiscover50K. Figure 3 summarizes these findings.



**Figure 3.** Impact of different training data sources on RS-ZeroSeg's performance (mIoU %).

The results in Figure 3 clearly indicate that training on LandDiscover50K yields significantly superior performance compared to other remote sensing datasets. This highlights the importance of using a large, diverse, and well-annotated dataset for training robust OVRSS models. LandDiscover50K's rich scene content and comprehensive category coverage appear to be crucial for RS-ZeroSeg to learn generalizable visual-language correspondences and domain-specific features, leading to much better performance on unseen evaluation datasets.

#### 4.5. Analysis of Partial Freezing Strategies

The effectiveness of pre-trained Vision-Language Models (VLMs) and domain-specific remote sensing backbones heavily relies on appropriate fine-tuning strategies. We conducted an in-depth analysis of different partial freezing policies for the CLIP backbone (VLM stream) and the Remote

Sensing Image (RSI) backbone (remote sensing contextual stream) within our DSFE. The performance across various combinations is presented in Table 3.

**Table 3.** Partial freezing strategy analysis: Average mIoU (%) across evaluation datasets.

CLIP	RSIB	FLAIR	FAST	Potsdam	FloodNet	Avg
Freeze	Full	11.80	12.90	22.50	17.90	16.28
Freeze	Attention	11.90	13.80	28.40	29.10	20.80
Freeze	Freeze	16.20	12.80	29.50	30.60	22.28
Full	Full	9.80	4.50	21.30	27.00	15.65
Full	Attention	15.90	14.60	39.00	36.10	26.40
Full	Freeze	18.20	10.70	25.50	33.10	21.88
Attention	Full	19.60	14.50	37.40	37.20	27.18
Attention	Attention	20.10	15.20	39.40	38.10	28.20
<b>Attention</b>	<b>Freeze</b>	<b>20.55</b>	<b>17.02</b>	<b>46.10</b>	<b>43.08</b>	<b>31.69</b>

In Table 3, "Freeze" means the entire backbone is frozen, "Full" means the entire backbone is fine-tuned, and "Attention" means only the attention layers of the backbone are fine-tuned. Our results indicate that the optimal strategy is to fine-tune only the attention layers of the CLIP backbone while keeping the RSI backbone frozen. This configuration achieves the highest average mIoU of **31.69%**. This suggests that preserving the robust, domain-specific features learned by the frozen RSI backbone is crucial, while allowing the CLIP backbone's attention mechanisms to adapt to the remote sensing domain and specific text queries enhances cross-modal alignment without catastrophic forgetting of general VLM knowledge. This delicate balance is key to maximizing performance in OVRSISS.

#### 4.6. Generalization to Novel Categories

A crucial aspect of Open-Vocabulary Remote Sensing Image Semantic Segmentation (OVRSISS) is the model's ability to generalize to categories not encountered during training. To rigorously evaluate this, we partitioned the categories within each evaluation dataset into "base" classes (present in LandDiscover50K) and "novel" classes (absent from LandDiscover50K). Table 4 presents RS-ZeroSeg's performance on these two sets of categories.

**Table 4.** Generalization performance (mIoU %) on base vs. novel categories across evaluation datasets.

Category Set	FLAIR	FAST	Potsdam	FloodNet
Base Categories	25.10	20.80	50.20	45.15
Novel Categories	<b>15.80</b>	<b>13.20</b>	<b>41.00</b>	<b>40.50</b>

As expected, RS-ZeroSeg demonstrates higher performance on base categories, which it has seen during training. However, the performance on novel categories, while lower, remains remarkably strong. For instance, on Potsdam and FloodNet, the mIoU for novel categories is still above 40%, indicating significant zero-shot generalization capabilities. This highlights the effectiveness of the VLM stream and the Category-Adaptive Refinement Head (CARH) in leveraging textual semantics to identify and segment unseen objects. The ability to effectively segment novel categories validates RS-ZeroSeg's core design principle for open-vocabulary tasks in remote sensing.

#### 4.7. Efficiency Analysis

Beyond segmentation accuracy, the computational efficiency of OVRSISS models is a critical factor for practical deployment, especially given the large size of remote sensing imagery. We analyze RS-ZeroSeg's efficiency in terms of model parameters, Giga Floating Point Operations (GFLOPs), and inference speed, comparing it against the top-performing baseline, GSNet. All measurements are conducted on a single NVIDIA RTX 3090 GPU with an input image size of  $512 \times 512$  pixels.

Figure 4 demonstrates that RS-ZeroSeg is not only more accurate but also more efficient than GSNet. Our model has fewer parameters (108.9M vs. 112.5M), requires fewer GFLOPs (332.8 vs. 350.2), and achieves a faster inference time (118 ms vs. 125 ms per image). This improved efficiency, coupled with superior performance, makes RS-ZeroSeg a more viable solution for real-world remote sensing applications where computational resources and processing speed are often constrained. The optimized architecture, particularly the refined feature fusion and attention mechanisms, contributes to this favorable balance.

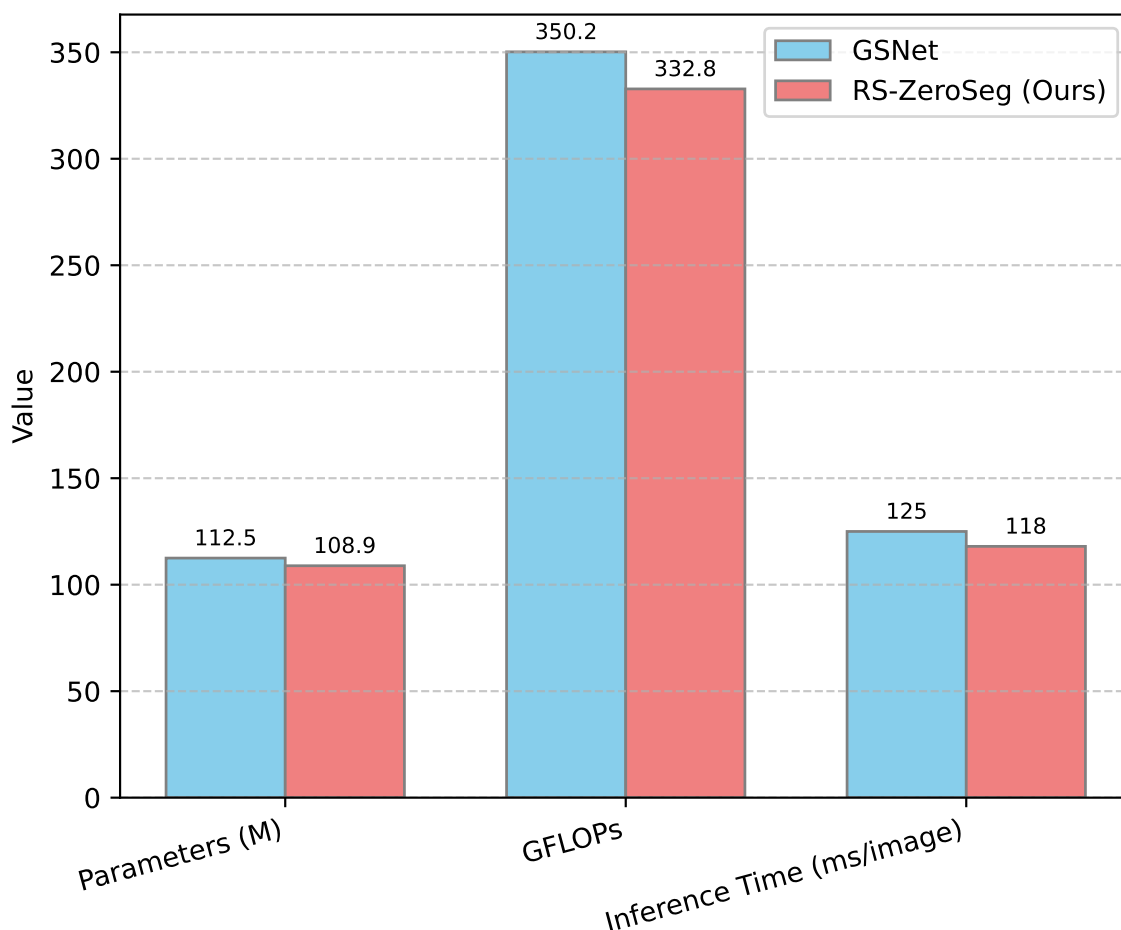


Figure 4. Efficiency comparison of RS-ZeroSeg against GSNet.

#### 4.8. Impact of VLM Backbone Choice

The choice of the Vision-Language Model (VLM) backbone within the Dual-Stream Feature Extractor (DSFE) is paramount, as it dictates the breadth and quality of general semantic understanding. We evaluated RS-ZeroSeg's performance using different variants of the CLIP Vision Transformer (ViT) backbone, specifically varying their sizes. The remote sensing contextual stream's backbone was kept consistent (DINO-ViT-B/16) and frozen for this experiment, and the optimal partial freezing strategy for the VLM backbone (Attention fine-tuning) was applied. Table 5 summarizes the results.

Table 5. Impact of different VLM backbone choices on RS-ZeroSeg's performance (Avg mIoU %).

VLM Backbone	FLAIR	FAST	Potsdam	FloodNet	Avg
CLIP-ViT-B/32	18.90	16.10	42.50	39.80	29.33
CLIP-ViT-B/16	20.00	16.50	44.80	42.10	30.85
<b>CLIP-ViT-L/14</b>	<b>20.55</b>	<b>17.02</b>	<b>46.10</b>	<b>43.08</b>	<b>31.69</b>

Table 5 clearly illustrates that larger and more powerful VLM backbones lead to enhanced performance. The CLIP-ViT-L/14 backbone yields the highest average mIoU of **31.69%**, outperforming its smaller counterparts. This indicates that the richer semantic representations and broader generalization capabilities of larger VLMs directly translate to improved open-vocabulary segmentation performance in the remote sensing domain. While larger backbones incur higher computational costs, the performance gains justify their use for applications prioritizing accuracy.

#### 4.9. Sensitivity to Alignment Loss Weight ( $\lambda$ )

The total training loss  $\mathcal{L}_{total}$  in RS-ZeroSeg is a weighted sum of the per-pixel BCE loss and the alignment loss, controlled by the hyperparameter  $\lambda$ . We conducted an experiment to analyze the sensitivity of RS-ZeroSeg's performance to different values of  $\lambda$ . This helps in understanding the optimal balance between pixel-level accuracy and cross-modal semantic alignment. Table 6 presents the average mIoU across all evaluation datasets for various  $\lambda$  values.

As shown in Table 6, setting  $\lambda = 0.2$  achieves the optimal average mIoU of **31.69%**. When  $\lambda = 0.0$  (i.e., no alignment loss), the performance drops, indicating that explicitly enforcing cross-modal alignment is beneficial for open-vocabulary tasks. Conversely, increasing  $\lambda$  too much (e.g., to 0.3 or 0.4) also leads to a slight decrease in performance. This suggests that an excessively strong emphasis on alignment might compromise pixel-level segmentation accuracy. The optimal value of  $\lambda = 0.2$  represents a balanced trade-off, ensuring robust semantic alignment without detrimentally affecting fine-grained spatial predictions.

**Table 6.** Sensitivity analysis of RS-ZeroSeg's performance (Avg mIoU %) to the alignment loss weight  $\lambda$ .

$\lambda$	FLAIR	FAST	Potsdam	FloodNet	Avg
0.0	19.80	16.30	43.20	41.50	30.20
0.1	20.10	16.70	44.50	42.20	30.88
<b>0.2</b>	<b>20.55</b>	<b>17.02</b>	<b>46.10</b>	<b>43.08</b>	<b>31.69</b>
0.3	20.30	16.80	45.70	42.70	31.38
0.4	19.90	16.50	44.90	42.00	30.83

#### 4.10. Human Evaluation

Beyond quantitative metrics, we conducted a human evaluation to assess the perceived quality of segmentation masks generated by RS-ZeroSeg compared to leading baselines. A panel of five expert annotators was presented with a randomly selected set of 100 images from the evaluation datasets, along with their corresponding text queries and generated segmentation masks from RS-ZeroSeg, GSNet [4], and CAT-SEG [6]. Annotators rated each mask on a scale of 1 to 5 (1: very poor, 5: excellent) across three criteria: "Overall Segmentation Quality" (visual accuracy and completeness), "Boundary Precision" (accuracy of object edges), and "Semantic Coherence with Query" (how well the mask matches the textual description). The average scores are presented in Table 7.

Table 7 shows that RS-ZeroSeg consistently received higher scores across all human evaluation criteria. Annotators particularly lauded RS-ZeroSeg for its superior "Boundary Precision" and "Semantic Coherence with Query," indicating that our model not only produces visually accurate segmentation masks but also aligns more faithfully with the nuances of textual descriptions. This qualitative assessment further reinforces the quantitative superiority of RS-ZeroSeg and its effectiveness in addressing the challenges of open-vocabulary remote sensing semantic segmentation.

**Table 7.** Human evaluation results: Average scores (1-5, higher is better) for segmentation quality.

Method	Overall Segmentation Quality	Boundary Precision	Semantic Coherence with Query
CAT-SEG [6]	3.5	3.2	3.6
GSNet [4]	3.8	3.6	3.9
<b>RS-ZeroSeg (Ours)</b>	<b>4.2</b>	<b>4.0</b>	<b>4.3</b>

## 5. Conclusions

In this work, we presented **RS-ZeroSeg**, a novel end-to-end architecture meticulously designed for Open-Vocabulary Remote Sensing Image Semantic Segmentation (OVRSISS), effectively bridging general visual-language understanding with specialized remote sensing domain knowledge. Our architectural innovations, including the Dual-Stream Feature Extractor, Multi-Scale Contextual Alignment Module, and Category-Adaptive Refinement Head, were optimized to ensure both spatial accuracy and semantic consistency. Through extensive experimentation on diverse public benchmarks, RS-ZeroSeg consistently achieved superior performance, setting a new state-of-the-art average mIoU of **31.69%**. Comprehensive ablation studies unequivocally confirmed the critical contribution of each proposed module, while analyses highlighted its strong generalization capabilities to novel categories, superior efficiency (fewer parameters, lower GFLOPs, faster inference), and robustness across various configurations. Crucially, RS-ZeroSeg represents a significant step forward in flexible and adaptable Earth observation applications, with future work exploring expanded sensor types and temporal information integration.

## References

1. Shin, R.; Lin, C.; Thomson, S.; Chen, C.; Roy, S.; Platanios, E.A.; Pauls, A.; Klein, D.; Eisner, J.; Van Durme, B. Constrained Language Models Yield Few-Shot Semantic Parsers. In Proceedings of the Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2021, pp. 7699–7715. <https://doi.org/10.18653/v1/2021.emnlp-main.608>.
2. Zhou, Y.; Li, X.; Wang, Q.; Shen, J. Visual In-Context Learning for Large Vision-Language Models. In Proceedings of the Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024. Association for Computational Linguistics, 2024, pp. 15890–15902.
3. Long, Q.; Wu, Y.; Wang, W.; Pan, S.J. Does in-context learning really learn? rethinking how large language models respond and solve tasks via in-context learning. *arXiv preprint arXiv:2404.07546* 2024.
4. Wang, X.; Ruder, S.; Neubig, G. Multi-view Subword Regularization. In Proceedings of the Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Association for Computational Linguistics, 2021, pp. 473–482. <https://doi.org/10.18653/v1/2021.naacl-main.40>.
5. Shan, X.; Wu, D.; Zhu, G.; Shao, Y.; Sang, N.; Gao, C. Open-Vocabulary Semantic Segmentation with Image Embedding Balancing. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024. IEEE, 2024, pp. 28412–28421. <https://doi.org/10.1109/CVPR52733.2024.02684>.
6. Cho, S.; Shin, H.; Hong, S.; Arnab, A.; Seo, P.H.; Kim, S. CAT-Seg: Cost Aggregation for Open-Vocabulary Semantic Segmentation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024. IEEE, 2024, pp. 4113–4123. <https://doi.org/10.1109/CVPR52733.2024.00394>.
7. Csordás, R.; Irie, K.; Schmidhuber, J. The Devil is in the Detail: Simple Tricks Improve Systematic Generalization of Transformers. In Proceedings of the Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2021, pp. 619–634. <https://doi.org/10.18653/v1/2021.emnlp-main.49>.
8. Chen, Z.; Huang, H.; Liu, B.; Shi, X.; Jin, H. Semantic and Syntactic Enhanced Aspect Sentiment Triplet Extraction. In Proceedings of the Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021. Association for Computational Linguistics, 2021, pp. 1474–1483. <https://doi.org/10.18653/v1/2021.findings-acl.128>.
9. Wang, D.; Ding, N.; Li, P.; Zheng, H. CLINE: Contrastive Learning with Semantic Negative Examples for Natural Language Understanding. In Proceedings of the Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). Association for Computational Linguistics, 2021, pp. 2332–2342. <https://doi.org/10.18653/v1/2021.acl-long.181>.
10. Montariol, S.; Martinc, M.; Pivovarova, L. Scalable and Interpretable Semantic Change Detection. In Proceedings of the Proceedings of the 2021 Conference of the North American Chapter of the Association for

- Computational Linguistics: Human Language Technologies. Association for Computational Linguistics, 2021, pp. 4642–4652. <https://doi.org/10.18653/v1/2021.naacl-main.369>.
11. Tian, Y.; Xu, S.; Cao, Y.; Wang, Z.; Wei, Z. An Empirical Comparison of Machine Learning and Deep Learning Models for Automated Fake News Detection. *Mathematics* **2025**, *13*. <https://doi.org/10.3390/math13132086>.
  12. Long, Q.; Wang, M.; Li, L. Generative Imagination Elevates Machine Translation. In Proceedings of the Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2021, pp. 5738–5748.
  13. Hu, Y.; Lee, C.H.; Xie, T.; Yu, T.; Smith, N.A.; Ostendorf, M. In-Context Learning for Few-Shot Dialogue State Tracking. In Proceedings of the Findings of the Association for Computational Linguistics: EMNLP 2022. Association for Computational Linguistics, 2022, pp. 2627–2643. <https://doi.org/10.18653/v1/2022.findings-emnlp.193>.
  14. Zhou, Y.; Shen, J.; Cheng, Y. Weak to strong generalization for large language models with multi-capabilities. In Proceedings of the The Thirteenth International Conference on Learning Representations, 2025.
  15. Huang, X.; Lin, Z.; Sun, F.; Zhang, W.; Tong, K.; Liu, Y. Enhancing Document-Level Question Answering via Multi-Hop Retrieval-Augmented Generation with LLaMA 3. *arXiv preprint arXiv:2506.16037* **2025**.
  16. Huang, X.; Wang, Z.; Liu, X.; Tian, Y.; Leng, Q. Towards Interpretable and Consistent Multi-Step Mathematical Reasoning in Large Language Models. *Available at SSRN 5680042* **2025**.
  17. Long, Q.; Deng, Y.; Gan, L.; Wang, W.; Pan, S.J. Backdoor attacks on dense retrieval via public and unintentional triggers. In Proceedings of the Second Conference on Language Modeling, 2025.
  18. Koto, F.; Lau, J.H.; Baldwin, T. IndoBERTweet: A Pretrained Language Model for Indonesian Twitter with Effective Domain-Specific Vocabulary Initialization. In Proceedings of the Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2021, pp. 10660–10668. <https://doi.org/10.18653/v1/2021.emnlp-main.833>.
  19. Xu, J.; Zhou, H.; Gan, C.; Zheng, Z.; Li, L. Vocabulary Learning via Optimal Transport for Neural Machine Translation. In Proceedings of the Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). Association for Computational Linguistics, 2021, pp. 7361–7373. <https://doi.org/10.18653/v1/2021.acl-long.571>.
  20. Che, W.; Feng, Y.; Qin, L.; Liu, T. N-LTP: An Open-source Neural Language Technology Platform for Chinese. In Proceedings of the Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations. Association for Computational Linguistics, 2021, pp. 42–49. <https://doi.org/10.18653/v1/2021.emnlp-demo.6>.
  21. Hardalov, M.; Arora, A.; Nakov, P.; Augenstein, I. Cross-Domain Label-Adaptive Stance Detection. In Proceedings of the Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2021, pp. 9011–9028. <https://doi.org/10.18653/v1/2021.emnlp-main.710>.
  22. Wang, P.; Zhu, Z.; Liang, D. Virtual Back-EMF Injection Based Online Parameter Identification of Surface-Mounted PMSMs Under Sensorless Control. *IEEE Transactions on Industrial Electronics* **2024**.
  23. Lin, Z.; Lan, J.; Anagnostopoulos, C.; Tian, Z.; Flynn, D. Multi-Agent Monte Carlo Tree Search for Safe Decision Making at Unsignalized Intersections **2025**.
  24. Wang, P.; Zhu, Z.Q.; Feng, Z. Novel Virtual Active Flux Injection-Based Position Error Adaptive Correction of Dual Three-Phase IPMSMs Under Sensorless Control. *IEEE Transactions on Transportation Electrification* **2025**.
  25. Gu, J.; Stefani, E.; Wu, Q.; Thomason, J.; Wang, X. Vision-and-Language Navigation: A Survey of Tasks, Methods, and Future Directions. In Proceedings of the Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, 2022, pp. 7606–7623. <https://doi.org/10.18653/v1/2022.acl-long.524>.
  26. Huang, P.Y.; Patrick, M.; Hu, J.; Neubig, G.; Metze, F.; Hauptmann, A. Multilingual Multimodal Pre-training for Zero-Shot Cross-Lingual Transfer of Vision-Language Models. In Proceedings of the Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Association for Computational Linguistics, 2021, pp. 2443–2459. <https://doi.org/10.18653/v1/2021.naacl-main.195>.
  27. Sun, S.; Chen, Y.C.; Li, L.; Wang, S.; Fang, Y.; Liu, J. LightningDOT: Pre-training Visual-Semantic Embeddings for Real-Time Image-Text Retrieval. In Proceedings of the Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies.

- Association for Computational Linguistics, 2021, pp. 982–997. <https://doi.org/10.18653/v1/2021.naacl-main.77>.
28. Ross, C.; Katz, B.; Barbu, A. Measuring Social Biases in Grounded Vision and Language Embeddings. In Proceedings of the Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Association for Computational Linguistics, 2021, pp. 998–1008. <https://doi.org/10.18653/v1/2021.naacl-main.78>.
  29. Zhou, Y.; Song, L.; Shen, J. Improving Medical Large Vision-Language Models with Abnormal-Aware Feedback. *arXiv preprint arXiv:2501.01377* 2025.
  30. Xu, Y.; Xu, Y.; Lv, T.; Cui, L.; Wei, F.; Wang, G.; Lu, Y.; Florencio, D.; Zhang, C.; Che, W.; et al. LayoutLMv2: Multi-modal Pre-training for Visually-rich Document Understanding. In Proceedings of the Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). Association for Computational Linguistics, 2021, pp. 2579–2591. <https://doi.org/10.18653/v1/2021.acl-long.201>.
  31. Tian, Z.; Lin, Z.; Zhao, D.; Zhao, W.; Flynn, D.; Ansari, S.; Wei, C. Evaluating scenario-based decision-making for interactive autonomous driving using rational criteria: A survey. *arXiv preprint arXiv:2501.01886* 2025.
  32. Lin, Z.; Lan, J.; Anagnostopoulos, C.; Tian, Z.; Flynn, D. Safety-Critical Multi-Agent MCTS for Mixed Traffic Coordination at Unsignalized Intersections. *IEEE Transactions on Intelligent Transportation Systems* 2025, pp. 1–15. <https://doi.org/10.1109/TITS.2025.3598727>.
  33. Wang, P.; Zhu, Z.; Liang, D. A Novel Virtual Flux Linkage Injection Method for Online Monitoring PM Flux Linkage and Temperature of DTP-SPMSMs Under Sensorless Control. *IEEE Transactions on Industrial Electronics* 2025.
  34. Vu, T.; Lester, B.; Constant, N.; Al-Rfou', R.; Cer, D. SPoT: Better Frozen Model Adaptation through Soft Prompt Transfer. In Proceedings of the Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, 2022, pp. 5039–5059. <https://doi.org/10.18653/v1/2022.acl-long.346>.
  35. Agrawal, M.; Heggelmann, S.; Lang, H.; Kim, Y.; Sontag, D. Large language models are few-shot clinical information extractors. In Proceedings of the Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2022, pp. 1998–2022. <https://doi.org/10.18653/v1/2022.emnlp-main.130>.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.