ALIGNCLIP: SELF-GUIDED ALIGNMENT FOR REMOTE SENSING OPEN-VOCABULARY SEMANTIC SEGMENTATION

Anonymous authorsPaper under double-blind review

000

001

002

004 005 006

007

008 009 010

011 012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032

034

035

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Open-Vocabulary Semantic Segmentation (OVSS) for remote sensing imagery plays a crucial role in applications such as land cover mapping and environmental monitoring. Recently, Contrastive Language-Image Pre-training (CLIP) has advanced the training-free paradigm of OVSS while also inspiring its exploration in the remote sensing domain. However, directly applying CLIP to remote sensing leads to cross-modal mismatches. Prevalent methods focus on exploring attention mechanism of CLIP visual encoder or introducing vision foundation models to obtain more discriminative feature, but they often overlook the alignment between patches and textual representations. To address this issue, we propose a training-free framework named AlignCLIP. We find that, objects of the same category tend to exhibit a more compact distribution in remote sensing, this enables a single visual feature to effectively represent all objects within the category. Based on this observation, we design the Self-Guided Alignment (SGA) module, which leverages the most reliable image-specific visual prototypes to refine the text embeddings. To mitigate interference among irrelevant features, we further introduce the Cluster-Constrained Enhancement (CCE) module, which clusters semantically similar patch features, suppresses inter-cluster correlations, and updates the logits map via a constraint propagation mechanism. Experiments on eight remote sensing benchmarks demonstrate that AlignCLIP consistently outperforms state-of-the-art training-free OVSS methods, achieving an average gain of +2.2 mIoU and offering a robust adaptive solution for open-vocabulary semantic segmentation in remote sensing. All code will be released.

1 Introduction

Open-vocabulary semantic segmentation (OVSS) in remote sensing imagery serves as a fundamental task in land cover mapping and environmental monitoring. Using arbitrary textual descriptions, it enables pixel-level classification of remote sensing images. The remarkable success of Contrastive Language–Image Pre-training (CLIP) (Radford et al., 2021b) in zero-shot recognition has inspired the development of OVSS. Most prior studies have focused on fine-tuning CLIP (Liang et al., 2023; Peng et al., 2024; Zhang et al., 2025), but their progress is limited by the demand for large annotated datasets. Moreover, remote sensing imagery often contains categories beyond the training set due to seasonal changes, land use evolution, and geographic diversity, making these approaches difficult to generalize. Recently, several works (Wang et al., 2023a; Yang et al., 2024; Zhou et al., 2022; Lan et al., 2024c) have begun to explore *training-free* paradigms in natural image domain, which achieve OVSS by extracting image patches and textual representations and directly performing cross-modal matching. This paradigm has further inspired its exploration in the remote sensing domain.

Prevalent *training-free* approaches in natural image domain primarily focus on the image modality, and they explore the attention mechanism of the CLIP visual encoder or integrate advanced vision foundation models (VFMs) to obtain more discriminative features (Lan et al., 2024b; Shao et al., 2024; Kim et al., 2025b; Barsellotti et al., 2024). However, these methods largely overlook the alignment between image patches and textual representations. Although they generate masks that align well with object boundaries, the inherent semantic gap between visual and textual modalities

055

060

061

063

064

066 067 068

069

071

072

073

078

079

080

081

083

084

085

087

880

089

090

091

092

093

094

095

096

098

100

101

102

103 104 105

107

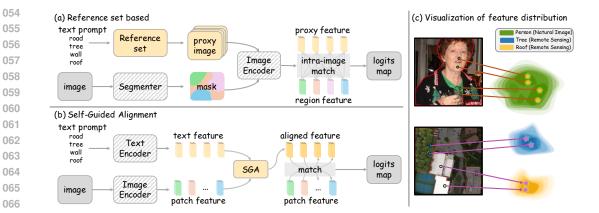


Figure 1: (a) Reference-set paradigm primarily focuses on constructing an accurate image-text matching set and performing matching based on proxy images. (b) Our SGA module refines the textual representation by selecting the most reliable visual prototypes from patch features and encourages mismatched patches to align with their corresponding textual representations. (c) We present feature visualization for both natural and remote sensing images and observe that, compared with natural images, objects of the same category in remote sensing imagery exhibit a more compact distribution.

is prone to causing mismatches between image patches and textual representations, i.e., cross-modal mismatches. To address this issue, another line of research leverages external image-text reference sets (Zermatten et al., 2025), specifically, these approaches transform text-image matching into intra-image matching by retrieving proxy images associated with category texts, as shown in Fig. 1(a), thereby mitigating cross-modal discrepancies. However, they heavily depend on the construction of cumbersome reference sets and exhibit limited generalization to unseen scenarios.

In this work, our experiments reveal that in natural images, objects belonging to the same category (e.g., humans) often exhibit diverse local feature distributions, and it is nearly impossible to identify a single representative feature that can capture all instances within the category. In contrast, remote sensing imagery demonstrates a markedly different characteristic, objects of the same category (e.g., trees, roofs) typically share similar textures and shapes, resulting in highly compact intra-class feature distribution, as shown in Fig. 1(c). This property enables the selection of a representative prototype feature that effectively characterizes all objects within the category. Based on this observation, we propose a simple yet effective training-free framework, termed AlignCLIP to mitigate cross-modal mismatches in OVSS of remote sensing imagery. We designed two key modules: (a) Self-Guided Alignment (SGA), which leverages the most reliable image-specific visual prototypes of the target image to refine textual embeddings, thereby bringing mismatched patches closer to their correct textual semantics. (b) Cluster-Constrained Enhancement (CCE), which clusters semantically similar patches while suppressing inter-cluster correlations, and updating logits map through constrained propagation.

Notably, AlignCLIP operates in a fully training-free manner, thereby eliminating the need for laborintensive reference sets construction. By relying solely on information inherent to the target image, it further ensures strong generalization across diverse scenarios. Extensive evaluations on eight remote sensing benchmarks demonstrate that AlignCLIP consistently outperforms state-of-the-art training-free OVSS methods, highlighting its robustness and adaptability to novel scenarios and unseen categories.

The main contributions of our work are as follows:

• We analyze the limitations of existing reference sets-based methods, and observe that objects of the same category in remote sensing imagery exhibit concentrated feature distribution. Leveraging this characteristic, we mitigate cross-modal mismatches while obviating the need for cumbersome reference sets construction.

- We propose AlignCLIP, a fully *training-free* framework that alleviates cross-modal mismatches. The framework incorporates the *Self-Guided Alignment (SGA)* module, which refines text embeddings using reliable image-specific prototypes, and the *Cluster-Constrained Enhancement (CCE)* module, which clusters image patches and suppresses the correlations between different clusters.
- Extensive experiments on eight remote sensing benchmarks demonstrate that AlignCLIP consistently outperforms state-of-the-art *training-free* OVSS methods, achieving both qualitative and quantitative improvements and exhibiting strong generalization to diverse scenarios and unseen categories.

2 Related work

Vision-Language Models. Vision-Language Models (VLMs) (Jia et al., 2021; Yuan et al., 2021) aim to align visual and textual representations within a shared semantic space, enabling zero-shot and open-vocabulary recognition. A landmark advancement in this field is CLIP (Radford et al., 2021b), a dual-encoder trained contrastively on image-text pairs with strong downstream generalization. However, CLIP is optimized for image-level classification, and its patch features are suboptimal for dense prediction (Cheng et al., 2022; Xu et al., 2022) due to limited spatial awareness and the absence of explicit spatial modeling. This issue is more prominent in the remote sensing domain, where high-resolution scenes exhibit fine spectral-textural details and large-scale layouts distinctly different from those of natural images (Cao et al., 2024; Zhang et al., 2025; Dutta et al., 2025). Although some works (e.g., RemoteCLIP (Liu et al., 2024), GeoRSCLIP (Zhang et al., 2024b)) have been adapted to remote sensing via prompt engineering or fine-tuning, such approaches typically require task-specific retraining or substantial labeled data, constraining their practicality for open-vocabulary semantic segmentation.

Vision Foundation Models. Vision Foundation Models (VFMs) (Caron et al., 2021; Oquab et al., 2023; Siméoni et al., 2025; Kirillov et al., 2023; Ravi et al., 2024) provide general visual representations across a wide range of tasks. One category of such models is DINO (Caron et al., 2021), which learns semantically rich and spatially coherent features via self-distillation. It can localize objects without explicit supervision, making it highly suitable for dense prediction tasks. Additionally, SAM (Kirillov et al., 2023) demonstrates strong image segmentation capabilities, supporting various segmentation prompts (*e.g.*, points, boxes, and masks) with excellent cross-domain generalization performance. In this work, we leverage the representations from VFMs to cluster semantically similar features and utilize inter-cluster correlations to update the logits map.

Training-free OVSS. *Training-free* open-vocabulary semantic segmentation (OVSS) labels pixels for arbitrary categories at inference by matching visual and textual embeddings from VLMs like CLIP via cross-modal similarities. prevalent works improve spatial awareness via attention modification (Yang et al., 2024; Wang et al., 2023a) or by integrating VFMs such as SAM (Zhang et al., 2024a; Lan et al., 2024c), but they overlook the correlations between patches and text representations. ReMe (Xuan et al., 2025) mitigate mismatches using curated reference sets, which are costly and difficult to generalize. In the remote sensing domain, SegEarth-OV (Li et al., 2025) represents the first *training-free* OVSS framework, which introduces an upsampling module to adapt CLIP, but still encounter cross-modal mismatches. Our work inherits the upsampling module of SegEarth-OV and designs two modules to alleviate the cross-modal mismatches, leveraging the characteristic of concentrated intra-class feature distribution in remote sensing imagery.

3 METHODOLOGY

In this section, we first introduce a preliminary of our framework in Sec. 3.1. Then, we introduce the *Self-Guided Alignment (SGA)* module in Sec. 3.2 and the *Cluster-Constrained Enhancement (CCE)* module in Sec. 3.3. Finally, we detail the integration with the upsampling module in Sec. 3.4. The overall framework is shown in Fig. 2.

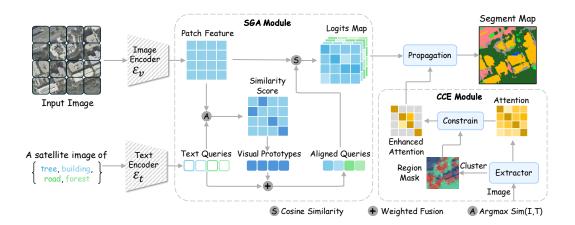


Figure 2: **The overall framework of AlignCLIP.** We propose a simple yet effective approach to alleviate cross-modal mismatches. We design two core modules (a) *Self-Guided Alignment (SGA)*, which refines textual embeddings using the most reliable image-specific visual prototypes. (b) *Cluster-Constrained Enhancement (CCE)*, which clusters semantically similar patches while suppressing inter-cluster correlations and updates the logits map through constrained propagation.

3.1 Preliminary

Given a remote sensing input image $I \in \mathbb{R}^{H \times W \times 3}$ and an open set of textual category names $\mathcal{T} = \{t_1, t_2, \dots, t_K\}$, where H, W denote the height and width of an image, K denotes the number of classes. The objective of open-vocabulary semantic segmentation (OVSS) is to assign each pixel in I to one of the categories in \mathcal{T} .

In the *training-free* setting, recent works adopt large-scale vision-language models (e.g., ViT-based CLIP) as the backbone for feature extraction and cross-modal matching. Specifically, the frozen CLIP image encoder \mathcal{E}_v divides the input image I into a grid of image patches and outputs a set of patch-level visual embeddings. For brevity, we omit the [CLS] token here:

$$\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N] \in \mathbb{R}^{N \times D}, \tag{1}$$

where D denotes the embedding dimension and $N=H_p\times W_p$ depends on the encoder's patch resolution. However, the CLIP model has limited capability in spatial awareness, previous studies have modified the attention score calculation in the last layer of self-attention in the CLIP visual encoder from *query-to-key* to *query-to-query* or *key-to-key*, which has significantly improved the performance of CLIP's dense prediction. Following the practice of prior works, we modified the calculation of the self-attention scores in the last layer of the visual encoder:

$$MSA(q, k, v) = \sum_{i \in \{q, k, v\}} softmax(\frac{i \cdot i^{T}}{\sqrt{d}}) \cdot v,$$
(2)

where q, k and v represent the *query*, key, and *value* matrices in self-attention, respectively, and d denotes the dimension of attention features. Meanwhile, to obtain more accurate text embeddings, we adopt a prompt template that is more suitable for remote sensing scenarios (e.g., "a satellite image of [CLS].") to incorporate contextual information, as opposed to the prompt template used for natural images. Subsequently, each text prompt is processed by the CLIP text encoder \mathcal{E}_t to obtain its corresponding textual embedding:

$$\mathbf{T} = [\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_K] \in \mathbb{R}^{K \times D}.$$
 (3)

Finally, we compute the logits map between each patch-level visual embedding $p_i \in \mathbf{P}$ and all textual embeddings \mathbf{T} using cosine similarity. The segmentation mask is obtained by applying the argmax operation to logits map:

$$S = sim(\mathbf{P}, \mathbf{T}), \quad S \in \mathbb{R}^{N \times K}.$$
 (4)

3.2 Self-Guided Alignment

In the process of cross-modal matching, the inherent gap between text and image leads to cross-modal mismatches. To address this, we design a *Self-Guided Alignment (SGA)* module, which exploits the intrinsic visual cues from the reliable image-specific visual prototypes to refine the textual embeddings. Formally, for each textual embedding $\mathbf{t}_k \in \mathbf{T}$, we compute its cosine similarity with all patch-level visual embeddings $\mathbf{p}_i \in \mathbf{P}$:

$$s_{i,k} = \frac{\mathbf{p}_i^{\mathsf{T}} \mathbf{t}_k}{\|\mathbf{p}_i\| \|\mathbf{t}_k\|}.$$
 (5)

We then select the most similar patch embedding \mathbf{p}_{i^*} for category k:

$$i^* = \arg\max_i \ s_{i,k}. \tag{6}$$

This selected patch embedding acts as an image-specific visual prototype extracted directly from the target image. We fuse the textual embedding \mathbf{t}_k with its corresponding visual prototype \mathbf{p}_{i^*} to obtain an aligned category embedding \mathbf{t}_k' :

$$\mathbf{t}_{k}' = (1 - \alpha) \cdot \mathbf{t}_{k} + \alpha \cdot \mathbf{p}_{i^{*}}, \tag{7}$$

where $\alpha \in [0,1]$ is a balancing hyperparameter controlling the contribution of textual and visual components.

Finally, the logits map for segmentation is computed by replacing the original textual embeddings with the aligned embeddings $\mathbf{T}' = \{\mathbf{t}'_1, \dots, \mathbf{t}'_K\}$:

$$S' = sim(\mathbf{P}, \mathbf{T}'). \tag{8}$$

In this way, those originally mismatched patches in remote sensing imagery will be closer to the aligned text embeddings due to the compact features of intra-class objects. Furthermore, since the prototypes are derived on-the-fly from the target image, the SGA module naturally adapts to new scenes without requiring any re-training or prebuilt reference sets.

3.3 Cluster-Constrained Enhancement

Although the SGA module mitigates cross-modal mismatches by refining text embeddings with image-specific visual prototypes, the image patches may still be disturbed by irrelevant patches. To address this issue, we introduce the *Cluster-Constrained Enhancement (CCE)* module, which aggregates semantically similar patches, suppresses interactions between irrelevant patches, and updates logits map through constrained propagation.

Specifically, we employ a VFM visual transformer (e.g., DINO, SAM) to extract a high discriminative visual feature map $\mathbf{F} \in \mathbb{R}^{H_p \times W_p \times D}$ from the target image. We reshape \mathbf{F} into N patch embeddings $\{\mathbf{f}_1, \ldots, \mathbf{f}_N\}$ and apply a clustering algorithm:

$$\{\mathcal{C}_1, \dots, \mathcal{C}_m, \dots \mathcal{C}_{K_c}\} = \text{Clustering}(\{\mathbf{f}_i\}_{i=1}^N),$$
 (9)

where C_m denotes the set of patch indices assigned to cluster m, and K_c is the total number of clusters—a hyperparameter controlling the granularity of segmentation refinement.

In addition to visual features, we also extract the self-attention matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ from the last layer of the visual transformer, which encodes pairwise affinities between patches. Subsequently, we constrain it using the clustering results. Specifically, a binary cluster mask matrix $\mathbf{M} \in \{0,1\}^{N \times N}$ is constructed as:

$$\mathbf{M}_{ij} = \begin{cases} 1, & \text{if } \mathcal{G}(i) = \mathcal{G}(j), \\ 0, & \text{otherwise,} \end{cases}$$
 (10)

where $\mathcal{G}(i)$ denotes the cluster assignment of patch i. The affinity matrix is then refined as:

$$\tilde{\mathbf{A}} = \mathbf{A} \odot \mathbf{M},\tag{11}$$

where \odot denoting element-wise multiplication. In this way, affinities are preserved only within the same cluster, while inter-cluster correlations are suppressed, logits map are propagated under the cluster-constrained affinities as follows:

$$\hat{\mathcal{S}}_i = \frac{1}{|\mathcal{C}_{\mathcal{G}(i)}|} \sum_{j \in \mathcal{C}_{\mathcal{G}(i)}} \tilde{\mathbf{A}}_{ij} \cdot \mathcal{S}'_j. \tag{12}$$

3.4 INTEGRATION WITH UPSAMPLING MODULE

To recover the fine-grained details critical for accurate segmentation in high-resolution remote sensing images, we inherit the upsample module from SegEarth-OV. Specifically, the visual feature map \mathbf{P} is first reshaped into a 2D feature representation $\mathbf{P} \in \mathbb{R}^{H_p \times W_p \times D}$, which is subsequently upsampled to the original image resolution. The upsampled features are then computed with the text embeddings \mathbf{T} via cosine similarity, yielding an upsampled logits map:

$$S_{up} = sim(featup(\mathbf{P}), \mathbf{T}). \tag{13}$$

And then, we interpolate the \hat{S} to match the spatial size of S_{up} , the two logits maps are then fused via a weighted combination:

$$S_{final} = \beta \cdot Interpolate(\hat{S}) + (1 - \beta) \cdot S_{up}, \tag{14}$$

where $\beta \in [0,1]$ is a fusion weight controlling the balance between CCE-refined logits map and upsampled logits map. Interpolate is a bilinear interpolation algorithm.

Finally, we apply an \mathbf{argmax} operation over \mathcal{S}_{final} to produce the final segmentation mask:

$$pred = \arg\max_{k} \mathcal{S}_{final}.$$
 (15)

4 EXPERIMENTS

4.1 SETTINGS

Datasets and Evaluation Metric. We conducted comprehensive experiments on eight widely used remote sensing semantic segmentation datasets. Among these, OpenEarthMap (Wang et al., 2023b), LoveDA (Wang et al., 2021), iSAID (Waqas Zamir et al., 2019), Potsdam (Gerke, 2014) and Vaihingen (Rottensteiner et al., 2014) are primarily composed of satellite images, while UAVid (Yang et al., 2020), UDD5 (Chen et al., 2018) and VDD (Pan et al., 2021) mainly consist of UAV images. These datasets collectively cover diverse spatial resolutions, imaging conditions, and scene types, thereby providing a comprehensive evaluation of model robustness. Each dataset contains multiple foreground categories along with a background category. Please refer to appendix A.1 for detailed dataset information. Following common practice in semantic segmentation, we report the mean Intersection over Union (mIoU) as the primary evaluation metric, which provides a balanced measure of classification accuracy across categories.

Baselines. We compared our AlignCLIP with a wide range of state-of-the-art *training-free* OVSS methods, including CLIP (Radford et al., 2021b), MaskCLIP (Zhou et al., 2022), SCLIP (Wang et al., 2023a), GEM (Bousselham et al., 2024), ClearCLIP (Lan et al., 2024a), NACLIP (Hajimiri et al., 2024), ResCLIP (Yang et al., 2024), ProxyCLIP (Lan et al., 2024c), CASS (Kim et al., 2025a), SC-CLIP (Bai et al., 2025), Trident (Shi et al., 2024) and CorrCLIP (Zhang et al., 2024a). These baselines represent different design paradigms such as attention modification and proxy-based adaptation. Furthermore, we evaluated SegEarth-OV, a method specifically tailored for remote sensing that employs a trained upsampling module to recover the lost detailed information in feature maps. It should be noted that the performance of reference-set-based methods (*e.g.*, ReME) largely depends on the scale and quality of the constructed reference set, making fair comparisons challenging. Therefore, we do not report evaluations of these methods in the experimental section.

Implementation Details. We provide two model variants of AlignCLIP, *i.e.*, AlignCLIP-D (integration with DINO) and AlignCLIP-S (integration with SAM). All experiments employ OpenCLIP (Radford et al., 2021a) to extract both image and text features. Unless otherwise specified, all models adopt ViT-B/16 as the default backbone. For the text encoder, we adopted a remotesensing-oriented prompt template, with the prompt list provided in appendix A.2. For the image encoder, we followed the settings of SegEarth-OV, input images were resized such that the long side was 448, and inference was conducted using a sliding window of size 224×224 with a stride of 112. For the clustering algorithm, we simply used the K-Means algorithm (Ikotun et al., 2023) with the number of clusters $K_c = 3$ as default. For the specific balance ratios α and fusion weights β of each dataset, please refer to appendix A.3. To isolate the effectiveness of our method, all post-processing techniques (*e.g.*, PAMR (Araslanov & Roth, 2020), denseCRF (Krähenbühl & Koltun,

Table 1: Quantitative comparison results on eight remote sensing datasets. **Bold** fonts indicate the optimal results, and <u>underlined</u> fonts indicate the suboptimal results. Avg. represents the average mIoU across the eight datasets.

Methods	OpenEarthMap	LoveDA	iSAID	Potsdam	Vaihingen	UAVid	UDD5	VDD	Avg.
CLIP _[ICML'21]	12.0	12.4	7.5	14.5	10.3	10.9	9.5	14.2	11.4
MaskCLIP[ECCV'22]	25.1	27.8	14.5	31.7	24.7	28.6	32.4	32.9	27.2
SCLIP _[ECCV'24]	29.3	30.4	16.1	36.6	28.4	31.4	38.7	37.9	31.1
GEM _[CVPR'24]	33.9	31.6	17.7	36.5	24.7	33.4	41.2	39.5	32.3
ClearCLIP[ECCV'24]	31.0	32.4	18.2	40.9	27.3	36.2	41.8	39.3	33.4
NACLIP _[WACV'25]	35.7	31.5	19.5	40.2	28.8	37.5	42.1	40.9	34.5
ResCLIP[CVPR'25]	34.2	31.2	20.0	42.6	28.2	37.6	42.3	40.3	34.6
ProxyCLIP[ECCV'24]	35.0	33.5	20.7	44.1	27.8	42.1	46.5	44.3	36.8
CASS _[CVPR'25]	34.6	34.0	20.6	42.9	31.5	38.6	39.0	40.9	35.3
SC-CLIP[ArXiv'24]	35.9	31.7	18.4	43.4	29.6	38.3	42.0	41.0	35.0
Trident[ICCV'25]	35.1	31.5	20.0	44.4	27.7	41.8	44.1	45.7	36.3
CorrCLIP _[ICCV'25]	35.4	32.7	16.9	42.6	24.7	38.1	40.1	37.7	33.5
SegEarth-OV _[CVPR'25]	39.8	36.9	21.7	<u>47.1</u>	29.1	42.5	50.6	45.3	39.1
AlignCLIP-D (Ours)	40.1	39.5	23.6	47.9	34.5	44.4	51.8	48.4	41.3
Alignetir-D (Ours)	(+0.3)	(+2.6)	(+1.9)	(+0.8)	(+3.0)	(+1.9)	(+1.2)	(+2.8)	(+2.2)
AlignCLIP-S (Ours)	40.1	39.5	23.4	47.8	34.6	44.4	51.8	48.1	41.2
Alignetir-5 (Ours)	(+0.3)	(+2.6)	(+1.7)	(+0.7)	(+3.1)	(+1.9)	(+1.2)	(+2.8)	(+2.1)

2011)) were disabled. Experiments were conducted on 8 RTX 3090 GPUs, and all the code of our implementation is based on **mmsegmentation** repository¹.

4.2 RESULTS

 Quantitative Evaluation. As shown in Table 1, AlignCLIP achieves the overall best performance across all eight remote sensing benchmarks, achieving a highest average mIoU of 41.3%, which outperforms all compared *training-free* OVSS methods. The improvements are particularly remarkable on datasets such as LoveDA (+2.6%), Vaihingen (+3.1%), and VDD (+2.8%). On the remaining datasets, including OpenEarthMap, iSAID, Potsdam, UAVid, and UDD5, our method also achieves steady gains over existing approaches. Moreover, compared to SegEarth-OV, AlignCLIP still achieves a substantial improvement (+2.2% on average). Interestingly, the two model variants based on DINO and SAM yield comparable results, indicating that our approach performs robustly across different VFM architectures. The above experimental results demonstrate that our method achieves consistent improvements across different scenarios and model architectures.

Qualitative Evaluation. As illustrated in Fig. 3, we present the qualitative visualization results of AlignCLIP-D and other representative methods. The results demonstrate that compared with CASS, Trident, CorrCLIP, and SegEarth-OV, our AlignCLIP generates more accurate and spatially coherent segmentation results across various datasets. Existing methods often suffer from category confusion (*e.g.*, walls vs. roofs in VDD). While SegEarth-OV improves boundary smoothness, it still exhibits matching errors in fine-grained structures. In contrast, AlignCLIP effectively mitigates cross-modal mismatches, producing clearer regions and sharper object boundaries. For more qualitative comparisons, refer to Appendix A.7.

4.3 ABLATION STUDIES

In this section, we conduct a comprehensive ablation study to evaluate the effectiveness of the proposed components in AlignCLIP and to examine the impact of key hyperparameters. For clarity, the hyperparameter sensitivity analysis reports results on four representative datasets—OpenEarthMap (OEM), Potsdam (PD), UAVid (UAV) and VDD, while the complete experimental results are provided in Appendix A.4. Unless otherwise specified, we use the AlignCLIP-D as the default model for our primary analysis.

¹https://github.com/open-mmlab/mmsegmentation

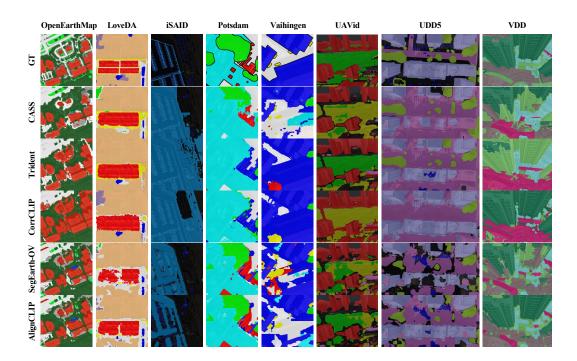


Figure 3: Qualitative comparison of different training-free OVSS methods on eight remote sensing datasets.

Component ablation analysis. We first investigate the effectiveness of our proposed SGA and CCE modules in AlignCLIP through component-wise ablation, as reported in Table 2. The 1^{st} row reports the performance of the baseline method SegEarth-OV, the 2^{nd} and 3^{rd} rows report the results of introducing the SGA and CCE modules respectively, while the 4^{th} row reports the performance of the complete model. We can find that: i) incorporating the SGA module alone yields a 1.0% improvement, demonstrating that alleviating cross-modal mismatches can substantially enhance segmentation performance. ii) applying the CCE module alone provides only a modest 0.3% gain, we attribute this to the fact that the CCE module operates on the aligned logits map produced by the SGA module, and propagating optimization on a poorly aligned logits map offers limited benefit.

Table 2: Ablation analysis of different components.

SGA	CCE	OpenEarthMap	LoveDA	iSAID	Potsdam	Vaihingen	UAVid	UDD5	VDD	Avg.
=	=	39.8	36.9	21.7	47.1	29.1	42.5	50.6	45.3	39.1
\checkmark		39.8	37.8	22.7	47.2	32.3	43.3	50.6	46.8	40.1 ↑1.0
	\checkmark	39.8	34.9	21.1	47.3	31.6	43.2	51.2	45.8	39.4 ↑0.3
\checkmark	\checkmark	40.1	39.5	23.6	47.9	34.5	44.4	51.8	48.4	41.3 ↑2.2

Effect of the balance ratios α . We further study the effect of the balance ratios α , which control the relative contributions of visual and textual features, a larger α assigns greater weight to the visual features. As shown in Table 3a, we observe that increasing α does not lead to a monotonic performance gain, instead, the performance generally rises initially and then declines (e.g., when $\alpha=0.3$ on the PD dataset). We attribute this phenomenon to the fact that, as the contribution of patch features increases, the text features can better align with the image features. However, beyond a certain threshold—determined by the feature distribution of the dataset, the image features begin to compromise the general representational capacity of the text features.

Effect of the fusion weights β . We further investigate the effect of the fusion weights β , which control the balance between our logits map and the upsampling logits map, a larger β indicates a smaller contribution of the upsampled logits map. As shown in Table 3b, smaller values (e.g.,

 $\beta=0.1$) achieve the best performance on the OEM and PD datasets, while larger values (e.g., $\beta=0.3$ and $\beta=0.4$) yield better results on the UAV and VDD datasets. We attribute this result to the differences in dataset scales, the OEM and PD datasets consist of large-scale satellite images, where detail information is more likely to be lost during feature extraction, thus requiring more contributions from the upsampled logits map to compensate. In contrast, UAV and VDD datasets contain small-scale UAV aerial images, where detail information is relatively preserved, making the contribution of the upsampled logits map relatively limited.

Effect of the cluster number K_c . We further analyze the effect of varying number of the clusters K_c used in the CCE module (i.e., $K_c = 3, 6, 9, 12, 15$), with the results summarized in Table 3c. The results indicate that $K_c = 3$ achieves the best performance across all datasets. Increasing the number of clusters leads to a slight performance decrease, with a drop of no more than 0.5%, suggesting that our method is not sensitive to the choice of K_c . We attribute this to the fact that $K_c = 3$ is sufficient for the model to distinguish irrelevant features, and further increasing the number of clusters does not yield significant performance gains.

Table 3: Sensitivity analysis of various hyperparameters across different datasets.

	(a) Balance ratios α					(b) Fusion weights β						(c) Cluster numbers K_c					
α	OEM	PD	UAV	VDD	β	OEM	PD	UAV	VDD		K_c	OEM	PD	UAV	VDD		
0.1	40.1	47.5	44.0	47.7	0.1	40.1	47.9	43.4	46.7		3	40.1	47.9	44.4	48.4		
0.2	40.1	47.6	44.4	48.4	0.2	40.0	47.8	44.0	47.9		6	40.1	47.6	44.2	48.2		
0.3	39.8	47.9	44.3	48.3	0.3	39.6	47.5	44.4	48.4		9	40.0	47.6	44.2	48.1		
0.4	39.3	47.8	43.9	47.9	0.4	39.1	46.9	44.4	48.4		12	40.0	47.6	44.1	48.0		
0.5	38.8	47.7	43.3	47.2	0.5	38.6	46.0	44.3	48.2		15	40.0	47.5	44.0	47.9		

Effect of different top-n visual prototypes. We further investigate the effect of varying the number of visual prototypes on text embeddings. Specifically, we compute the cosine distance between each patch feature and the text feature, select the n closest patch features (i.e., n=1,2,3,4,5,6), and average them before fusing with the text features. As illustrated in Fig. 4, our experiments reveal that increasing the number of patch features leads to a consistent performance decline across all eight datasets. This indicates that excessive prototypes not only fail to improve patch—text alignment but also introduce mismatched features, thereby increasing noise. Therefore, we select only the most similar feature to pursue the best performance.

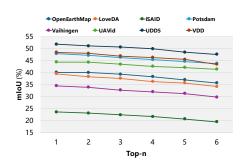


Figure 4: The effect of different numbers of top-n visual prototypes.

5 Conclusion

In this work, we presented AlignCLIP, a novel *training-free* framework for open-vocabulary semantic segmentation in the remote sensing domain. We find that features of intra-class objects in remote sensing tend to be compact. Based on this observation, we design two modules to alleviate cross-modal mismatches between image patches and textual representations. Specifically, the *Self-Guided Alignment (SGA)* module leverages the most reliable image-specific visual prototypes to refine textual embeddings, and the *Cluster-Constrained Enhancement (CCE)* clusters semantically similar patches while suppressing inter-cluster correlations, and updating logits map through constrained propagation. Extensive experiments across eight remote sensing benchmarks demonstrated that AlignCLIP consistently outperforms state-of-the-art approaches, We hope this work can inspire future related research and bring new possibilities to *training-free* open-vocabulary semantic segmentation in the remote sensing domain.

REFERENCES

- Nikita Araslanov and Stefan Roth. Single-stage semantic segmentation from image labels. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 4253–4262, June 2020.
- Sule Bai, Yong Liu, Yifei Han, Haoji Zhang, and Yansong Tang. Self-calibrated clip for training-free open-vocabulary segmentation, 2025. URL https://arxiv.org/abs/2411.15869.
- Luca Barsellotti, Roberto Amoroso, Lorenzo Baraldi, and Rita Cucchiara. Fossil: Free openvocabulary semantic segmentation through synthetic references retrieval. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 1464–1473, 2024.
- Walid Bousselham, Felix Petersen, Vittorio Ferrari, and Hilde Kuehne. Grounding everything: Emerging localization properties in vision-language transformers. In 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 3828–3837, 2024. doi: 10.1109/CVPR52733.2024.00367.
- Qinglong Cao, Yuntian Chen, Chao Ma, and Xiaokang Yang. Open-vocabulary remote sensing image semantic segmentation. arXiv preprint arXiv:2409.07683, 2024.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In <u>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</u>, pp. 9650–9660, October 2021.
- Yu Chen, Yao Wang, Peng Lu, Yisong Chen, and Guoping Wang. Large-scale structure from motion with semantic constraints of aerial images. In <u>Chinese Conference on Pattern Recognition and Computer Vision (PRCV)</u>, pp. 347–359. Springer, 2018.
- Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In <u>Proceedings of the IEEE/CVF</u> conference on computer vision and pattern recognition, pp. 1290–1299, 2022.
- Saikat Dutta, Akhil Vasim, Siddhant Gole, Hamid Rezatofighi, and Biplab Banerjee. Aeroseg: Harnessing sam for open-vocabulary segmentation in remote sensing images. In <u>Proceedings of the Computer Vision and Pattern Recognition Conference</u>, pp. 2254–2264, 2025.
- Markus Gerke. The isprs 2d semantic labeling contest potsdam dataset. In <u>ISPRS Annals of</u> the Photogrammetry, Remote Sensing and Spatial Information Sciences, 2014. URL https://www.isprs.org/education/benchmarks/UrbanSemLab/default.aspx.
- Sina Hajimiri, Ismail Ben Ayed, and Jose Dolz. Pay attention to your neighbours: Training-free open-vocabulary semantic segmentation. arXiv preprint, 2024. URL https://arxiv.org/abs/2404.08181.
- Abiodun M. Ikotun, Absalom E. Ezugwu, Laith Abualigah, Belal Abuhaija, and Jia Heming. Kmeans clustering algorithms: A comprehensive review, variants analysis, and advances in the era of big data. Information Sciences, 622:178–210, 2023. ISSN 0020-0255. doi: https://doi.org/10.1016/j.ins.2022.11.139. URL https://www.sciencedirect.com/science/article/pii/S0020025522014633.
- 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, pp. 4904–4916. PMLR, 2021.
- Chanyoung Kim, Dayun Ju, Woojung Han, Ming-Hsuan Yang, and Seong Jae Hwang. Distilling spectral graph for object-context aware pen-vocabulary semantic segmentation. In <u>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</u>, 2025a.
- Chanyoung Kim, Dayun Ju, Woojung Han, Ming-Hsuan Yang, and Seong Jae Hwang. Distilling spectral graph for object-context aware open-vocabulary semantic segmentation. In <u>Proceedings</u> of the Computer Vision and Pattern Recognition Conference, pp. 15033–15042, 2025b.

- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. arXiv:2304.02643, 2023.
 - Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In Proceedings of the 25th International Conference on Neural Information Processing Systems, NIPS'11, pp. 109–117, Red Hook, NY, USA, 2011. Curran Associates Inc. ISBN 9781618395993.
 - Mengcheng Lan, Chaofeng Chen, Yiping Ke, Xinjiang Wang, Litong Feng, and Wayne Zhang. Clearclip: Decomposing clip representations for dense vision-language inference. In <u>European</u> Conference on Computer Vision, pp. 143–160. Springer, 2024a.
 - Mengcheng Lan, Chaofeng Chen, Yiping Ke, Xinjiang Wang, Litong Feng, and Wayne Zhang. Proxyclip: Proxy attention improves clip for open-vocabulary segmentation. In <u>European Conference</u> on Computer Vision, pp. 70–88. Springer, 2024b.
 - Mengcheng Lan, Chaofeng Chen, Yiping Ke, Xinjiang Wang, Litong Feng, and Wayne Zhang. Proxyclip: Proxy attention improves clip for open-vocabulary segmentation. In <u>European Conference</u> on Computer Vision, pp. 70–88. Springer, 2024c.
 - 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. In Proceedings of the Computer Vision and Pattern Recognition Conference, pp. 10545–10556, 2025.
 - 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, pp. 7061–7070, 2023.
 - Yuting Lin, Kumiko Suzuki, and Shinichiro Sogo. Practical techniques for vision-language segmentation model in remote sensing. The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences, 48:203–210, 2024.
 - Fan Liu, Delong Chen, Zhangqingyun Guan, Xiaocong Zhou, Jiale Zhu, Qiaolin Ye, Liyong Fu, and Jun Zhou. Remoteclip: A vision language foundation model for remote sensing. <u>IEEE Transactions on Geoscience and Remote Sensing</u>, 62:1–16, 2024. doi: 10.1109/TGRS.2024. 3390838. URL https://doi.org/10.1109/TGRS.2024.3390838.
 - Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2023.
 - X. Pan, Y. Li, J. Chen, and Z. Wang. Vdd: A new benchmark dataset for semantic segmentation of uav imagery. Remote Sensing, 13(7):1302, 2021. doi: 10.3390/rs13071302.
 - Zelin Peng, Zhengqin Xu, Zhilin Zeng, Changsong Wen, Yu Huang, Menglin Yang, Feilong Tang, and Wei Shen. Understanding fine-tuning clip for open-vocabulary semantic segmentation in hyperbolic space-supplementary material.
- Zelin Peng, Zhengqin Xu, Zhilin Zeng, Yu Huang, Yaoming Wang, and Wei Shen. Parameter-efficient fine-tuning in hyperspherical space for open-vocabulary semantic segmentation. In Proceedings of the Computer Vision and Pattern Recognition Conference, pp. 15009–15020, 2025.
 - Alec Radford, Jong Wook Kim, Chris Hallacy, A. Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In ICML, 2021a.

- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In Marina Meila and Tong Zhang (eds.), Proceedings of the 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pp. 8748–8763. PMLR, 18–24 Jul 2021b. URL https://proceedings.mlr.press/v139/radford21a.html.
- Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos. arXiv:2408.00714, 2024. URL https://arxiv.org/abs/2408.00714.
- Franz Rottensteiner, Gunho Sohn, Jaewan Jung, Markus Gerke, Caroline Baillard, Silvia Benitez, and Uwe Breitkopf. The isprs semantic labeling benchmark (vaihingen dataset). In ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences, 2014. URL https://www.isprs.org/education/benchmarks/UrbanSemLab/default.aspx.
- Tong Shao, Zhuotao Tian, Hang Zhao, and Jingyong Su. Explore the potential of clip for training-free open vocabulary semantic segmentation. In <u>European Conference on Computer Vision</u>, pp. 139–156. Springer, 2024.
- Yuheng Shi, Minjing Dong, and Chang Xu. Harnessing vision foundation models for high-performance, training-free open vocabulary segmentation. arXiv preprint arXiv:2411.09219, 2024.
- Oriane Siméoni, Huy V. Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, Francisco Massa, Daniel Haziza, Luca Wehrstedt, Jianyuan Wang, Timothée Darcet, Théo Moutakanni, Leonel Sentana, Claire Roberts, Andrea Vedaldi, Jamie Tolan, John Brandt, Camille Couprie, Julien Mairal, Hervé Jégou, Patrick Labatut, and Piotr Bojanowski. DINOv3, 2025. URL https://arxiv.org/abs/2508.10104.
- Feng Wang, Jieru Mei, and Alan Yuille. Sclip: Rethinking self-attention for dense vision-language inference. arXiv preprint arXiv:2312.01597, 2023a.
- Jiwei Wang, Shunping Zhang, Wei Wang, Kun Fu, Zhiyong Li, Zhenwei Shi, Wei Wei, and Wei Liu. Loveda: A remote sensing land-cover dataset for domain adaptive semantic segmentation. In Advances in Neural Information Processing Systems (NeurIPS), 2021. URL https://arxiv.org/abs/2110.08733.
- Xiaoyang Wang, Yue Wu, Htoo Htoo Aung, Xiaoxiang Liu, and Xiao Xiang Zhu. Openearthmap: A benchmark dataset for global high-resolution land cover mapping. In <u>Advances in Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track</u>, 2023b. URL https://arxiv.org/abs/2307.15062.
- Syed Waqas Zamir, Aditya Arora, Akshita Gupta, Salman Khan, Guolei Sun, Fahad Shahbaz Khan, Fan Zhu, Ling Shao, and Gui-Song Xia. isaid: A large-scale dataset for instance segmentation in aerial images. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops, 2019. URL https://arxiv.org/abs/1905.12886.
- Yixuan Wei, Han Hu, Zhenda Xie, Ze Liu, Zheng Zhang, Yue Cao, Jianmin Bao, Dong Chen, and Baining Guo. Improving clip fine-tuning performance. In <u>Proceedings of the IEEE/CVF</u> International Conference on Computer Vision, pp. 5439–5449, 2023.
- Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit: Semantic segmentation emerges from text supervision. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 18134–18144, 2022.
- Xiwei Xuan, Ziquan Deng, and Kwan-Liu Ma. Reme: A data-centric framework for training-free open-vocabulary segmentation, 2025. URL https://arxiv.org/abs/2506.21233.

- Yuhang Yang, Jinhong Deng, Wen Li, and Lixin Duan. Resclip: Residual attention for training-free dense vision-language inference. arXiv preprint arXiv:2411.15851, 2024.
- Zhenhua Yang, Liang Wang, Yue Zhang, Licheng Wang, and Zhaoxiang Zhang. Uavid: A semantic segmentation dataset for uav imagery. <u>ISPRS Journal of Photogrammetry and Remote Sensing</u>, 165:108–119, 2020. doi: 10.1016/j.isprsjprs.2020.02.017.
- Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for computer vision. arXiv preprint arXiv:2111.11432, 2021.
- Quan-Sheng Zeng, Yunheng Li, Daquan Zhou, Guanbin Li, Qibin Hou, and Ming-Ming Cheng. Maskclip++: A mask-based clip fine-tuning framework for open-vocabulary image segmentation. 2024.
- Valérie Zermatten, Javiera Castillo-Navarro, Diego Marcos, and Devis Tuia. Learning transferable land cover semantics for open vocabulary interactions with remote sensing images. <u>ISPRS Journal</u> of Photogrammetry and Remote Sensing, 220:621–636, 2025.
- Dengke Zhang, Fagui Liu, and Quan Tang. Corrclip: Reconstructing patch correlations in clip for open-vocabulary semantic segmentation. arXiv preprint arXiv:2411.10086, 2024a.
- Qiang Zhang, Decheng Wang, and Xiao Yu. Rlita: A region-level image–text alignment method for remote sensing foundation model. Remote Sensing, 17(10):1661, 2025.
- Zilun Zhang, Tiancheng Zhao, Yulong Guo, and Jianwei Yin. Rs5m and georsclip: A large scale vision-language dataset and a large vision-language model for remote sensing. <u>IEEE Transactions</u> on Geoscience and Remote Sensing, pp. 1–1, 2024b. doi: 10.1109/TGRS.2024.3449154.
- Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from clip. In <u>European</u> Conference on Computer Vision (ECCV), 2022.

A APPENDIX

A.1 DATASET DESCRIPTION

OpenEarthMap contains 5,000 aerial and satellite remote sensing images, including 8 foreground classes and 1 background class, with a spatial resolution of 0.25-0.5 meters, covering 97 regions in 44 countries / territories on six continents. we use the validation set for evaluation.

LoveDA includes remote sensing images with a spatial resolution of 0.3 meters covering multiple cities, totaling 5,987 images. These images are annotated with 6 foreground classes and one background class. we use the validation set for evaluation.

iSAID consists of images captured by the JL-1 satellite and GF-2 satellite. It includes 15 foreground classes and one background class. Following the data processing pipeline of MMSegmentation, we cropped the images into rectangles with a size of 896 and an overlapping area of 384. Finally, 33,978 images were generated for training and 11,644 for validation. In this study, the validation set was used for evaluation.

Potsdam comprises 38 image patches with a spatial resolution of 0.05 meters, with an average size of 6,000×6,000 pixels. It includes 5 foreground categories and 1 background category. Following the data processing pipeline of MMSegmentation, we used the validation set for evaluation.

Vaihingen comprises 33 image patches with a spatial resolution of 0.09 meters, with an average size of 2,494×2,064 pixels. It includes 5 foreground categories and 1 background category. Following the data processing pipeline of MMSegmentation, we used the validation set for evaluation.

UAVid is a 4K semantic segmentation video dataset for urban scenes, which contains a large number of street views and is annotated with 6 foreground classes and 1 background class. In this study, its test set was used for evaluation.

UDD5 consists of images collected by unmanned aerial vehicles (UAVs), including 4 foreground classes and 1 background class. In this study, we used its validation set for evaluation.

VDD is a collection of UAV images featuring diverse scenes, camera angles, and varying weather/lighting conditions. It provides high-resolution annotated images at the 400-pixel scale. With 6 foreground classes and 1 background class, its test set was used for evaluation in this study.

A.2 REMOTE SENSING PROMPT TEMPLATE

To obtain more effective text embeddings for remote sensing scenarios, we carefully designed 80 prompt templates tailored to remote sensing scenarios to replace the prompt templates oriented to natural images. As shown in Table 4, five representative examples are presented. For each category, these prompt templates are used to generate corresponding text features, which are then averaged to obtain a semantically rich category feature representation for semantic segmentation.

Table 4: Examples of remote sensing prompt templates for generating text descriptions.

Remote sensing prompt templates a low-quality aerial image of [class]. a cropped remote sensing image of [class]. a remote sensing interpretation map of [class]. a satellite image containing hardly recognizable [class]. a low-resolution remote sensing image of [class].

A.3 HYPERPARAMETER SETTING

We provide the detailed hyperparameter settings for each dataset corresponding to the two model variants, *i.e.*, AlignCLIP-D and AlignCLIP-S, as shown in Table 5. The balance ratio α denotes

the contribution of text features and visual features, while the fusion weight β is used to control the fusion weight between the CCE-refined logits and the upsampled logits.

Table 5: Hyperparameter settings of AlignCLIP across different datasets.

hyperparameters	OpenEarthMap	LoveDA	iSAID	Potsdam	Vaihingen	UAVid	UDD5	VDD				
Integration with DINO												
α	0.1	0.7	0.5	0.3	0.3	0.2	0.3	0.2				
β	0.1	0.2	0.2	0.1	0.5	0.3	0.2	0.3				
	Integration with SAM											
α	0.1	0.5	0.5	0.3	0.3	0.2	0.3	0.2				
β	0.1	0.2	0.2	0.2	0.5	0.4	0.2	0.4				

A.4 SENSITIVITY ANALYSIS DETAILS

In this section, we present the detailed sensitivity analysis of the hyperparameters involved in the two model variants (i.e., AlignCLIP-D and AlignCLIP-S) across eight datasets, as shown in Table 6-8. Specifically, the balance ratios α is used to control the contribution of text features and visual features in the SGA module (see Sec. 3.2), the fusion weights β is employed to control the fusion balance between the CCE-refined logits map and the upsampled logits map (see Sec. 3.4), and K_c represents the number of clusters in the clustering algorithm within the CCE module, which is used to control the granularity of segmentation refinement (see Sec. 3.3).

Table 6: Sensitivity analysis of different balance ratios α across different datasets.

α	OpenEarthMap	LoveDA	iSAID	Potsdam	Vaihingen	UAVid	UDD5	VDD	Avg.			
	Integration with DINO											
0.1	40.1	36.2	21.7	47.5	32.7	44.0	51.5	47.7	40.2			
0.2	40.1	37.3	22.3	47.6	33.9	44.4	51.7	48.4	40.7			
0.3	39.8	38.2	22.9	47.9	34.5	44.3	51.8	48.3	40.9			
0.4	39.3	38.8	23.4	47.8	34.4	43.9	51.8	47.9	40.9			
0.5	38.8	39.2	23.6	47.7	34.0	43.3	51.7	47.2	40.7			
	'		Int	egration wit	h SAM							
0.1	40.1	36.1	21.6	47.5	32.7	44.0	51.4	47.4	40.1			
0.2	40.0	37.2	22.2	47.7	33.9	44.4	51.7	48.1	40.7			
0.3	39.9	38.1	22.7	47.8	34.6	44.4	51.8	48.1	40.9			
0.4	39.3	38.8	23.2	47.8	34.5	43.9	51.8	47.7	40.9			
0.5	38.7	39.2	23.4	47.8	34.1	43.4	51.7	47.0	40.7			

Table 7: Sensitivity analysis of different fusion weights β across different datasets.

β	OpenEarthMap	LoveDA	iSAID	Potsdam	Vaihingen	UAVid	UDD5	VDD	Avg.				
	Integration with DINO												
0.1	40.1	38.5	22.2	47.9	31.1	43.4	51.4	46.7	40.1				
0.2	40.0	39.5	23.6	47.8	32.5	44.0	51.8	47.9	40.9				
0.3	39.6	39.3	23.6	47.5	33.4	44.4	51.8	48.4	41.0				
0.4	39.1	38.4	22.2	46.9	34.1	44.4	51.4	48.4	40.6				
0.5	38.6	37.4	20.4	46.0	34.5	44.3	50.9	48.2	40.0				
			Int	egration wit	h SAM								
0.1	40.1	38.4	22.1	47.8	31.1	43.4	51.4	46.6	40.1				
0.2	40.0	39.5	23.4	48.0	32.5	44.0	51.8	47.7	40.9				
0.3	39.6	39.3	23.4	47.8	33.4	44.4	51.8	48.1	41.0				
0.4	39.1	38.5	22.1	47.2	34.1	44.5	51.4	48.2	40.6				
0.5	38.6	37.4	20.4	46.3	34.6	44.4	50.9	47.9	40.1				

Table 8: Sensitivity analysis of different cluster numbers K_c across different datasets.

K_c	OpenEarthMap	LoveDA	iSAID	Potsdam	Vaihingen	UAVid	UDD5	VDD	Avg.			
	Integration with DINO											
3	40.1	39.5	23.6	47.9	34.5	44.4	51.8	48.4	41.3			
6	40.1	39.3	23.5	47.6	34.4	44.2	51.8	48.2	41.1			
9	40.0	39.2	23.5	47.6	34.3	44.2	51.7	48.1	41.1			
12	40.0	39.1	23.4	47.6	34.1	44.1	51.6	48.0	41.0			
15	40.0	39.0	23.3	47.5	34.1	44.0	51.6	47.9	41.0			
			Inte	egration wit	h SAM							
3	40.1	39.5	23.4	47.8	34.6	44.4	51.8	48.1	41.2			
6	40.1	39.3	23.4	47.7	34.4	44.3	51.7	48.1	41.1			
9	40.0	39.2	23.3	47.6	34.3	44.2	51.6	48.0	41.0			
12	40.0	39.1	23.3	47.6	34.2	44.2	51.6	47.9	41.0			
15	40.0	39.0	23.3	47.6	34.1	44.1	51.5	47.8	40.9			

A.5 SEAMLESS INTEGRATION INTO OTHER METHODS

In this section, we further validate the generality of the proposed approach by integrating the SGA module into other representative frameworks and conducting comprehensive evaluations on eight remote sensing benchmark datasets. Specifically, we incorporate SGA as an independent, plug-and-play component into existing methods. Since SGA aligns only the text embeddings without altering the remaining architecture, it can be seamlessly integrated into a variety of CLIP-based frameworks. For the experimental setup, we select ProxyCLIP, SC-CLIP, and CorrCLIP as baseline methods, with the hyperparameter α uniformly set to 0.1.

As shown in Table 9. The experimental results demonstrate that, across different baseline models, incorporating our method enables the text embeddings to achieve more precise alignment with the patch features. This finding not only confirms that enhancing image—text alignment can significantly improve semantic segmentation performance, but also highlights the generality of the SGA module, which can be seamlessly integrated into other CLIP-based frameworks.

Table 9: The proposed SGA module is integrated as a plugin into other methods.

Methods	OpenEarthMap	LoveDA	iSAID	Potsdam	Vaihingen	UAVid	UDD5	VDD	Avg.
ProxyCLIP	35.0	33.5	20.7	44.1	27.8	42.1	46.5	44.3	36.8
+SGA	38.6	34.2	21.6	44.6	32.7	42.4	48.3	45.3	38.5 ↑1.7
SC-CLIP	35.9	31.7	18.4	43.4	29.6	38.3	42.0	41.0	35.0
+SGA	39.8	32.8	19.6	43.6	31.6	39.6	46.0	42.3	36.9 ↑1.9
CorrCLIP	35.4	32.7	16.9	42.6	24.7	38.1	40.1	37.7	33.5
+SGA	36.6	33.5	18.6	43.8	27.9	39.9	41.1	39.6	35.1 ↑1.6

A.6 PSEUDO CODE OF OUR ALIGNCLIP

To clearly present the implementation details of our method and ensure reproducibility, we provide pseudo code for the two core modules of AlignCLIP, *i.e.*, SGA and CCE, in Algorithm 1 and Algorithm 2, respectively. In addition, the full implementation of our method (based on PyTorch), is provided in the supplementary materials, and the complete code will be publicly released after curation.

886 887 888

913 914

915

916

917

Algorithm 1 Pseudo code for Self-Guided Alignment in a PyTorch-like style.

```
865
866 1
        def self_guided_alignment(image_features, query_features, visual_query_alpha):
867
    3
            Self-Guided Alignment (SGA) module.
    4
868
           Args:
869
    6
               image_features: [num_patches, feature_dim]
               query_features: [num_queries, feature_dim]
870
               visual_query_alpha: balance ratio (0~1)
871 9
872 10
           Returns:
            Aligned query features: [num_queries, feature_dim]
    11
873 12
874 13
           num_queries = len(query_features)
    14
875 15
           # Similarity between image patches and query features
876 16
           _, index = similarity.topk(1, dim=-1)
877 18
878 19
            # Gather top patch features and average
    20
           visual_query_features = torch.gather(
879 21
               image_features.unsqueeze(0).repeat(num_queries, 1, 1),
               dim=1,
880 22
    23
               index=index.unsqueeze(-1).repeat(1, 1, image_features.shape[-1])
881 24
           ).mean(dim=1)
882 25
   26
            # Fuse with visual features
883 27
           aligned_query_features = visual_query_alpha * visual_query_features + \
884 28
                               (1 - visual_query_lambda) * query_features
            return aligned_query_features / aligned_query_features.norm(dim=-1, keepdim=True)
885
```

Algorithm 2 Pseudo code for Cluster-Constrained Enhancement in a PyTorch-like style.

```
889
         def cluster_constrained_enhancement(vfm_features, logits_map, cluster_num):
890
891 3
             Cluster-Constrained Enhancement (CCE) Module.
892
             Args:
                 vfm_features: Feature map for clustering, shape [num_patches, feature_dim].
893 6
                 logits_map: original logits map, shape [num_patches, num_classes].
894
                 cluster num: Number of clusters to group patches.
895 9
896 11
             Returns:
                 refined logits map: [num patches, num classes].
897 12
    13
898
             # Cluster the features
    14
             _, cluster_ids = perform_clustering(vfm_features, n_clusters=cluster_num)
899 15
900 16
             # Calculate patch attention
901 18
            vfm_attn = vfm_features @ vfm_features.T
902 20
             % Calculate masked attn based on clustering results
903 21
            masked_attn = torch.zeros_like(vfm_attn)
904 \begin{array}{c} 22 \\ 23 \end{array}
             for cluster_id in np.unique(cluster_ids):
905 24
                 # Create mask for current clust
906 25 26
                 mask = (cluster_ids == cluster_id)
907 27
                 # Aggregate attention within cluster
908 28
                 masked_attn[mask] = vfm_attn[mask, :] * mask[None, :] # element-wise masking
909 30
             #Propagate attention to refine logits map
910 \frac{31}{32}
             refined_logits = propagate_aff(logits_map, aff=final_attn)
911 33
             return refined_logits
912
```

A.7 ADDITIONAL QUALITATIVE RESULTS

We provide additional visualization analysis results for eight datasets to further validate the effectiveness of our proposed method, as illustrated in Fig. 5-12.

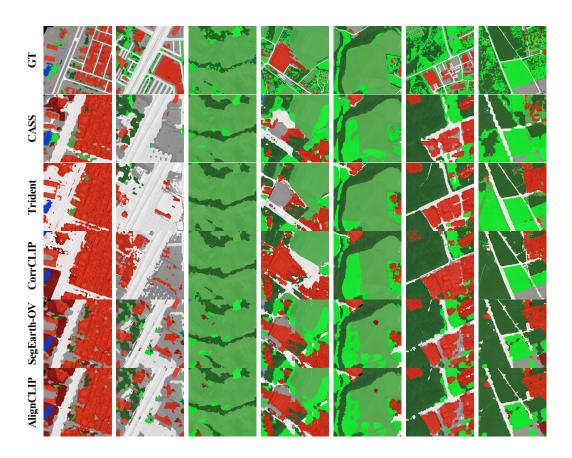


Figure 5: Qualitative comparison of different training-free OVSS methods on OpenEarthMap.

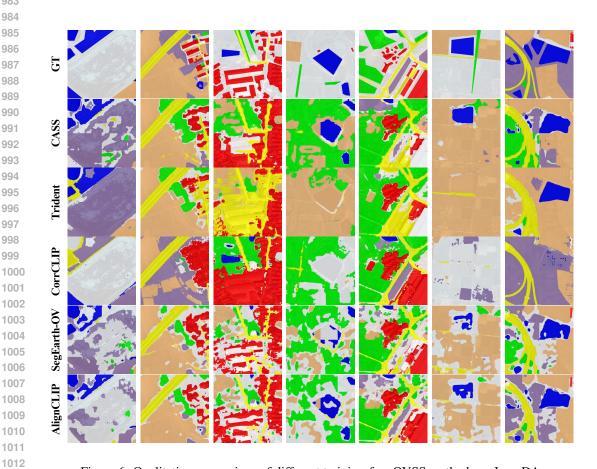


Figure 6: Qualitative comparison of different training-free OVSS methods on LoveDA.

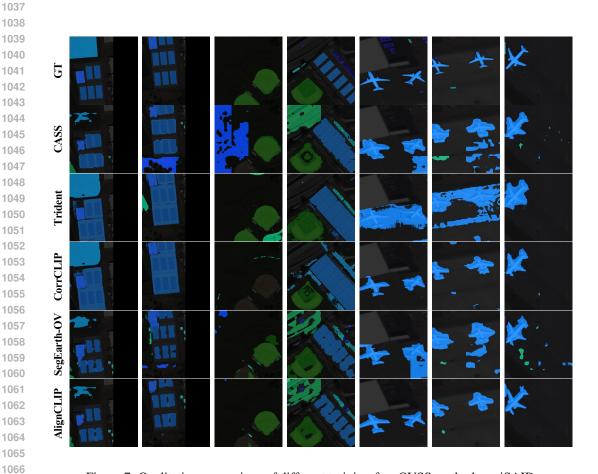


Figure 7: Qualitative comparison of different training-free OVSS methods on iSAID.

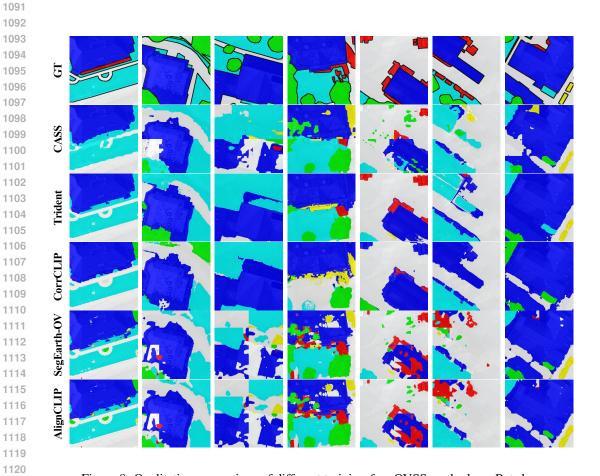


Figure 8: Qualitative comparison of different training-free OVSS methods on Potsdam.

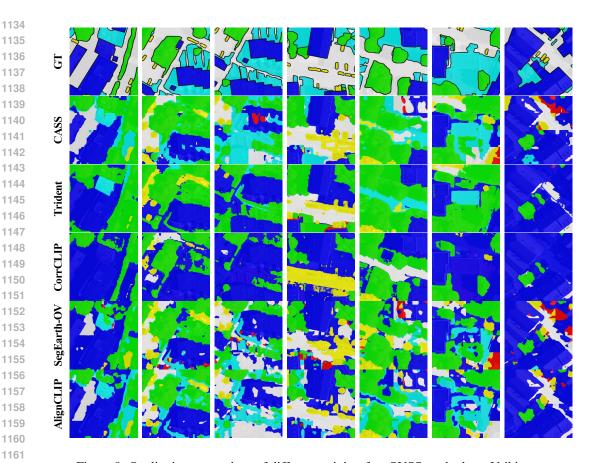


Figure 9: Qualitative comparison of different training-free OVSS methods on Vaihingen.

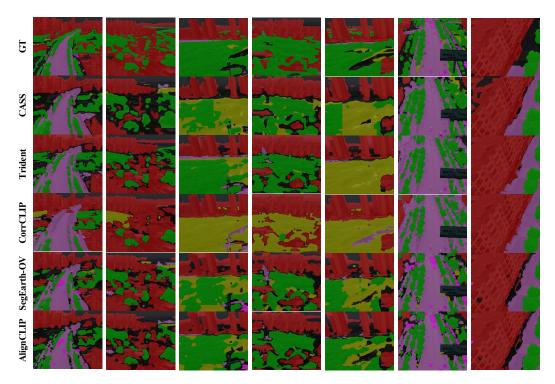


Figure 10: Qualitative comparison of different training-free OVSS methods on UAVid.

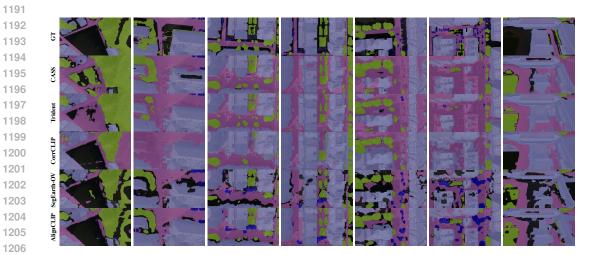


Figure 11: Qualitative comparison of different training-free OVSS methods on UDD5.

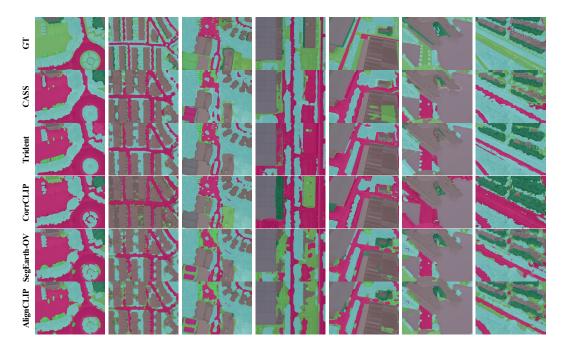


Figure 12: Qualitative comparison of different training-free OVSS methods on VDD.