# RADSeg: Unleashing Parameter and Compute Efficient Zero-Shot Open-Vocabulary Segmentation Using Agglomerative Models

Omar Alama[*1], Darshil Jariwala[*2], Avigyan Bhattacharya[*1],
Seungchan Kim[1], Wenshan Wang[1], Sebastian Scherer[1]

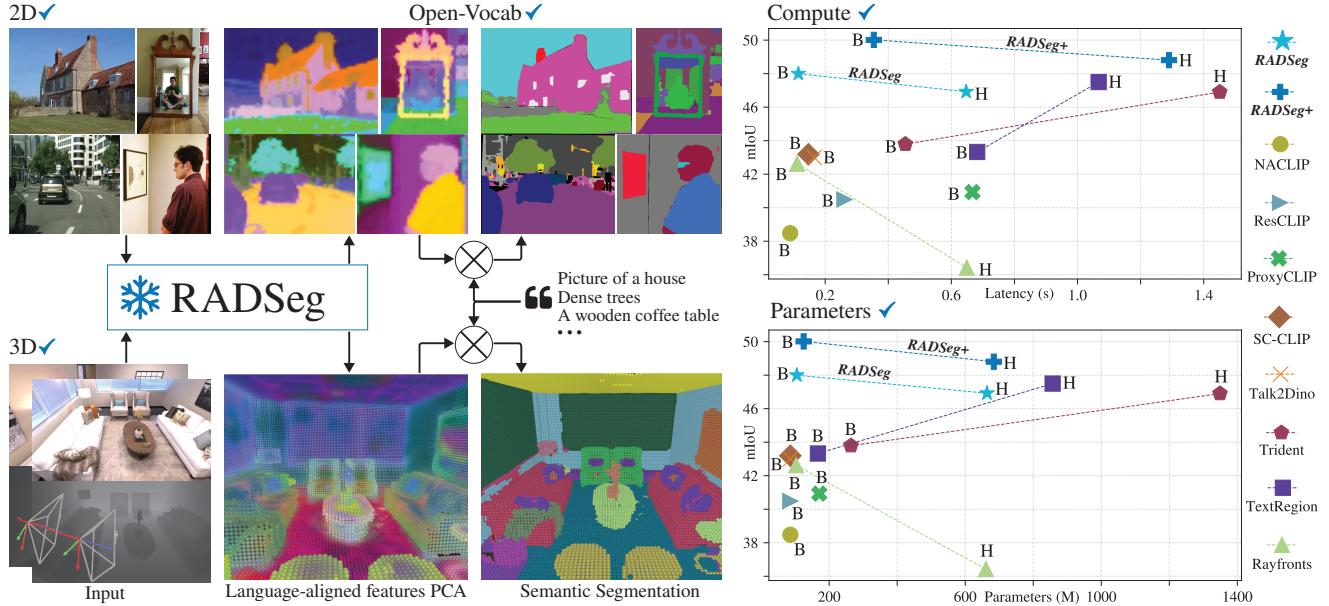[1]Carnegie Mellon University [2]IIIT Hyderabad

Figure 1. **RADSeg** is a dense, language-aligned feature encoder that enables low-parameter, low-latency open-vocabulary semantic segmentation in 2D and 3D. The efficiency plots report average latency, parameter count, and mIoU across five 2D datasets on a V100. By enhancing spatial locality of RADIO features, **RADSeg outperforms previous state-of-the-art methods in accuracy while remaining highly efficient in terms of parameter counts and inference speed.**

## Abstract

*Open-vocabulary semantic segmentation (OVSS) underpins many vision and robotics tasks that require generalizable semantic understanding. Existing approaches either rely on limited segmentation training data, which hinders generalization, or apply zero-shot heuristics to vision-language models (e.g CLIP), while the most competitive approaches combine multiple models to improve performance at the cost of high computational and memory demands. In this work, we leverage an overlooked agglomerative vision foundation model, RADIO, to improve zero-shot OVSS along three key axes simultaneously: mIoU, latency, and parameter efficiency. We present the first comprehensive study of RADIO for zero-shot OVSS and enhance its performance through self-correlating recursive attention, self-correlating global aggregation, and computationally effi-cient mask refinement. Our approach, **RADSeg**, achieves **6-30% mIoU improvement** in the base ViT class while being **3.95x faster** and using **2.5x fewer parameters**. Surprisingly, **RADSeg**-base (105M) outperforms previous combinations of huge vision models (850-1350M) in mIoU, achieving state-of-the-art accuracy with substantially lower computational and memory cost.*

## 1. Introduction

Semantic segmentation models often need to be deployed in open-world environments and support an open-set of tasks. This has led to growing interest in open-vocabulary semantic segmentation (OVSS), which enables segmentation based on arbitrary natural langauge descriptions, with applications ranging from medical analysis [24], robotic manipulation [30], to autonomous exploration and navigation [18, 46]. While vision language models (VLMs) such as CLIP [27] support open-vocabulary image classification,

---
[*]Equal contribution

their reliance on global image-text encoding fundamentally limits their capacity for fine-grained, pixel-level localization required for OVSS. Various approaches have since attempted to adapt CLIP to the OVSS task. Some adopt a training-based approach, either incorporating CLIP into a new architecture [21, 43, 49], or fine-tuning CLIP specifically for OVSS [29]. Both approaches, however, suffer from limited generalization, as existing semantic segmentation datasets [5, 7, 11, 25, 48] are orders of magnitude smaller than image-captioning datasets containing billions of image-text pairs [32].

Training-free, zero-shot approaches have emerged as a promising alternative, alleviating the limitations of training-based methods. Dense-inference strategies feed multiple image crops into the global CLIP encoder to produce dense, language-aligned feature maps [15, 41], but incur substantial computational overhead. Alternative methods attempt to restore the spatial alignment of the CLIP's patch-wise tokens via lightweight modifications but yield limited accuracy [13, 39, 44]. More recent approaches [20, 36, 42] improve segmentation by augmenting VLMs with combinations of other foundation models such as DINOv2 [26] and SAM [19], albeit at the cost of significantly increased model complexity, inference latency, and GPU memory overhead.

Given these limitations of existing OVSS methods, we investigate an underexplored agglomerative vision foundation model, RADIO [28]. By distilling knowledge from multiple foundation models into a single model, RADIO has shown strong performance across various vision benchmarks, in terms of generalization and efficiency. However, its potential for zero-shot OVSS remains largely unexplored. For instance, RayFronts [1] demonstrated RADIO's zero-shot capability for 3D semantic mapping, but evaluations of its broader OVSS abilities remain limited.

In this work, we provide the first comprehensive study of RADIO for zero-shot OVSS. Our analysis includes compute and parameter efficiency for competing baselines, revealing sharp trade-offs between mIoU, inference latency, and GPU memory, illustrated in Fig. 1. We also examine the ability of 2D OVSS methods to generate open-vocabulary semantic 3D maps, evaluating multi-view consistency.

Furthermore, we propose a novel and efficient pipeline, **RADSeg**, based on RADIOv3. **RADSeg** exploits an emergent property that enables alignment of RADIO's patch-wise features with language, enhances spatial locality by using Self-Correlating Recursive Attention (SCRA), and mitigates sliding-window artifacts through Self-Correlating Global Aggregation (SCGA). Finally, we leverage RADIOv3's efficient access to SAM-huge features to further refine the obtained masks using only 0.7% of SAM-huge parameters. **RADSeg** achieves **6-30% mIoU improvement** in the base ViT class over the next best baseline (Trident [36]) across datasets, while being **3.95x faster** and requiring **2.5x fewer parameters**. Remarkably, **RADSeg**-base **(105M)** surpasses the mIoU of previous compositions of huge vision models like Trident [36] **(1350M)** and TextRegion [42] **(850M)**, **achieving state-of-the-art mIoU with low computational and GPU memory requirements**. Our contributions can be summarized as:

- We conduct the first thorough empirical study of the foundational model RADIO for zero-shot open-vocabulary semantic segmentation (ZSOVSS), demonstrating its superiority to other backbones.
- We introduce a novel pipeline using RADIOv3 that advances the state-of-the-art in ZSOVSS while maintaining low parameter counts and computational efficiency.
- We perform comprehensive evaluations on five 2D and three 3D datasets across various resolutions and model sizes, showing that **RADSeg** consistently outperforms baselines at any resolution or model size budget.

## 2. Related work

### 2.1. Vision and Language Foundation Models

Vision-language foundation models (VLMs) trained on large amounts of image-text pairs have shown great generalization capabilties for various vision-language tasks [27, 38, 47]. VLMs align natural language and images into the same latent feature space, enabling zero-shot open-vocabulary classification by comparing text and image embeddings. However, such models are trained to align only a single CLS token that captures the global image context, hence termed global image-level VLMs. Consequently, the patch-wise/spatial tokens often exhibit poor spatial locality and language alignment. Recent research has sought to address these limitations to unlock open-vocabulary semantic segmentation (OVSS), which we discuss next.

### 2.2. Open-Vocabulary Semantic Segmentation

OVSS allows users to segment any image by simply providing natural language queries, removing the need for re-training when new categories are encountered in the wild. To achieve this, many existing methods fine-tune global image-level VLMs such as CLIP [27] on limited segmentation datasets at the expense of generalization ability [21, 29, 43, 49], i.e., training-based approaches.

In contrast, training-free approaches have attracted particular interest, as they generate dense, class-agnostic, language-aligned feature maps directly from images and generalize to arbitrary datasets in a zero-shot manner. Some of these methods adopt a 'dense-inference' strategy [15, 41], forwarding multiple crops or segments of an image through global VLMs and aggregating the results to obtain the map, yet incur substantial computational overhead. Others improve spatial locality and language alignment in global VLMs through lightweight attention modifications
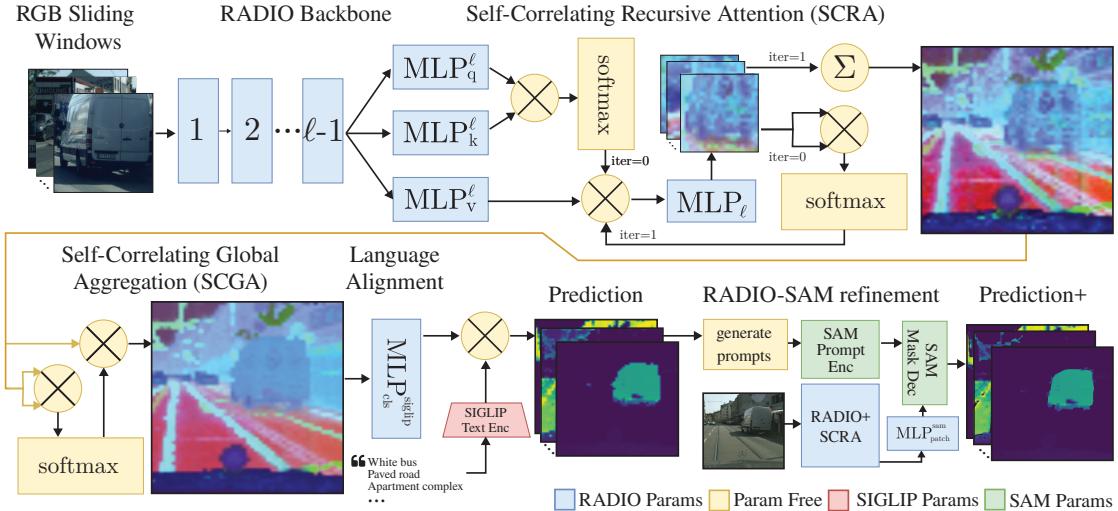
Figure 2. Overview of the **RADSeg** pipeline. RGB sliding windows are processed by the RADIO backbone. Self-Correlating Recursive Attention (SCRA) computes a similarity matrix from these outputs, which is recursively fed back into the last attention block of RADIO. Feature windows are aggregated into a feature map and refined through Self-Correlating Global Aggregation (SCGA) to reduce noise and windowing artifacts. Features are language-aligned with the SigLIP CLS token adaptor, and predictions are made by comparing them with text embeddings. Optionally, masks can be further refined using RADIO-SAM, requiring only **+20M** additional parameters.

[13, 39], or residual connections [3, 44]. While these methods are computationally efficient, they yield lower segmentation accuracy. Another line of training-free methods augment VLMs with additional foundation models to produce dense language aligned features. ProxyClip [20] leverages DINO's [6] strong spatial locality to refine CLIP features, Trident [36] extends it by incorporating SAM [19] for refinement; and TextRegion [42] integrates SAM2 [31] with CLIP. While these methods leveraging multiple foundation models achieve impressive segmentation accuracy, they incur significant memory and computational costs.

Distinctively, our work incorporates RADIO [14, 28]–an agglomerative model that distills knowledge from SAM, CLIP, SigLIP [47], and DINOv2 [26]–within a novel pipeline. By leveraging this unified foundation model, our approach outperforms the methods that explicitly combine multiple foundation models in segmentation accuracy, while maintaining the computational and memory efficiency of lightweight attention modification methods.

### 2.3. 3D Open Vocabulary Semantic Segmentation

A growing body of work attempts to leverage 2D VLMs to achieve training-free/generalizable OVSS in 3D and create open-vocabulary 3D maps. These maps take various forms, ranging from semantic neural radiance fields (NeRFs) [17, 33], gaussian splats [16, 35], scene graphs [12, 40] to simple point clouds or voxels [1, 15]. The most common strategy for obtaining dense, language-aligned features is to detect bounding boxes with GroundingDINO [23], use the detections to prompt SAM [19] for segmenta-

tion, then cropping each segment and forwarding it through CLIP [27]. However, this pipeline is computationally and memory intensive. Recently, RayFronts [1] showed that using RADIOv2.5's SigLIP CLS-token adaptor to project RADIO's patch-wise tokens, combined with neighbor-aware attention [13], can yield dense, language-aligned feature maps–achieving impressive results on 3D OVSS datasets. Building on this insight, we conduct the first comprehensive evaluation of RADIOv2.5 and RADIOv3 on both 2D and 3D ZSOVSS, and introduce novel spatial locality improvements to enable fast, memory efficient open-vocabulary 3D mapping with state-of-the-art segmentation performance.

## 3. Method

This section overviews the **RADSeg** pipeline (Fig. 2), covering preliminaries on RADIO (Sec. 3.1), then **RADSeg**'s main components–dense language alignment (Sec. 3.2), SCRA (Sec. 3.3), SCGA (Sec. 3.4)–and the optional mask refinement process in **RADSeg**+ (Sec. 3.5).

### 3.1. RADIO Preliminary

RADIO [14, 28] is an agglomerative vision foundation model that unifies knowledge distilled from CLIP [27], SigLIP[47], SAM [19], and DINOv2 [26] into a single model. RADIO employs a vision transformer (ViT) [10]: an image $I \in \mathbb{R}^{H \times W \times 3}$ is divided into $N_{\text{patch}} = \frac{H \times W}{P_h \times P_w}$ non-overlapping patches $P \in \mathbb{R}^{N_{\text{patch}} \times P_h \times P_w \times 3}$, which are linearly projected into patch-wise tokens $T_{\text{patch}} \in \mathbb{R}^{N_{\text{patch}} \times D}$, where $D$ is the token embedding dimension. RADIO also

prepends four separate CLS tokens, one for each teacher model (CLIP, SigLIP, SAM, DINOv2), yielding $T_{\text{cls}} \in \mathbb{R}^{4 \times D}$. We denote the combined tokens as $T \in \mathbb{R}^{N \times D}$, where $N = N_{\text{patch}} + 4 + R$ ($R$ register tokens $T_{\text{reg}}$ are additionally introduced to mitigate high-norm non-local output tokens [9]). Patch tokens $T_{\text{patch}}$ model spatially covariant/local features and CLS tokens $T_{\text{cls}}$ model global features summarizing the image. Unlike other agglomerative models that only align patch-wise output tokens to teachers (e.g. Theia [34]), RADIO aligns both patch-wise and CLS output tokens to their teacher counterparts through lightweight two-layer MLPs ($\text{MLP}_{\text{patch}}$ and $\text{MLP}_{\text{cls}}$), yielding adapted tokens $\tilde{T}_{\text{patch}}$ and $\tilde{T}_{\text{cls}}$. Formally,

$$\tilde{T}_{\text{patch}}^{(m)} = \text{MLP}_{\text{patch}}^{(m)}(T_{l,\text{patch}}), \tilde{T}_{\text{cls}}^{(m)} = \text{MLP}_{\text{cls}}^{(m)}(T_{l,\text{cls}}^{(m)}) \quad (1)$$

where $m \in \{\text{clip}, \text{siglip}, \text{sam}, \text{dinov2}\}$[1] and the subscript $l$ denote *last* block for output tokens. Finally, RADIO demonstrates strong spatial locality in its patch-wise features, a property important for dense prediction tasks such as semantic segmentation. Since VLMs like CLIP and SigLIP only align the CLS tokens to language, RADIO learns a global image-language alignment.

## 3.2. Getting Dense Language Alignment

RADIO has only been trained to align its SigLIP CLS output token $T_{l,\text{cls}}^{\text{siglip}}$ to the CLS token of the SigLIP teacher:

$$\tilde{T}_{\text{cls}}^{(\text{siglip})} = \text{MLP}_{\text{cls}}^{(\text{siglip})}(T_{l,\text{cls}}^{(\text{siglip})}), \min_{\theta} \mathcal{L}(\tilde{T}_{\text{cls}}^{(\text{siglip})}, \hat{T}_{l,\text{cls}}^{(\text{siglip})}) \quad (2)$$

where $\tilde{T}$ are adapted tokens, $\hat{T}$ are original teacher tokens, and $\mathcal{L}(x, y)$ is a distance loss function that gets minimized. This training only aligns the *global summary* of an image with language, enabling **zero-shot open-vocabulary classification**. Similarly, RADIO aligns its patch output tokens $T_{l,\text{patch}}$ to SigLIP's output patch tokens $\hat{T}_{l,\text{patch}}^{\text{siglip}}$ through an MLP. However, SigLIP patch tokens exhibit poor spatial-language alignment as shown in Tab. 3. Consequently, RADIO trained a linear probe on small segmentation datasets, limiting generalization to a closed vocabulary.

In contrast, we build on an observed emergent property of RADIO that enables dense, language-aligned feature extraction [1]. By simply applying the SigLIP **CLS** token adaptor on RADIO's **patch** tokens,

$$\text{MLP}_{\text{cls}}^{(\text{siglip})}(T_{l,\text{patch}}),$$

we obtain dense, language-aligned features suitable for OVSS. We empirically show that this emergent property holds for RADIOv2.5 and the newer RADIOv3, across
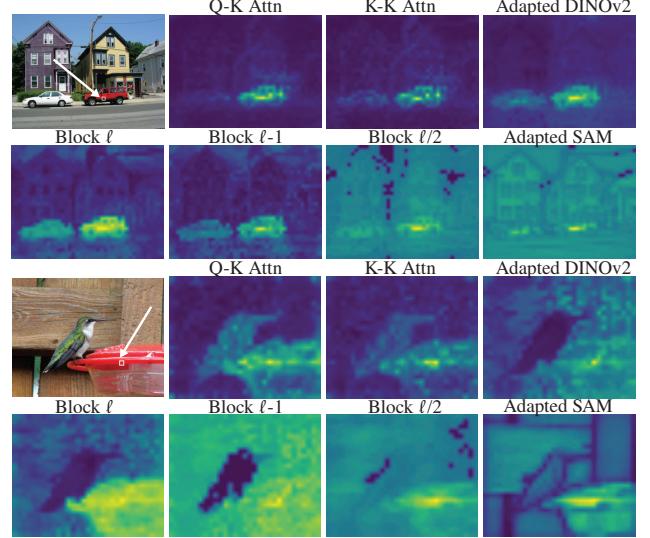


Figure 3. Qualitative comparison of last block attention and patch-wise similarity at different parts of the RADIO framework. The output of the RADIO backbone (Block $l$) can consistently attend to semantically similar patches, motivating our SCRA approach.

model sizes and for both SigLIP and CLIP adaptation. However, directly using the resulting feature maps yields suboptimal segmentation performance, motivating further alignment refinements in the following sections.

## 3.3. Self-Correlating Recursive Attention

A paradigm that has proven successful in previous training-free OVSS approaches [3, 13, 20, 36, 39, 44] is to exploit the last self-attention block, denoted as $\text{SA}_l$:

$$\text{SA}_l(T) = \text{Attn}(\text{MLP}_Q^l(T), \text{MLP}_K^l(T), \text{MLP}_V^l(T)) \quad (3)$$

, where $\text{Attn}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right)V$. The last self-attention block in encoders like CLIP and SigLIP compromises spatial locality by focusing on CLS token. This can be mitigated by encouraging attention to semantically similar patches. Several methods have proposed to recover this spatial alignment either through light-weight attention modifications yielding limited accuracy, or by incorporating additional backbones to provide better attention, achieving increased accuracy at the cost of memory and latency.

In contrast to these extremes, our goal is to achieve high accuracy *without* additional overheads. To this end, we examine what can be leveraged directly within RADIO. As shown in Fig. 3, we compare various attention and patch-similarity maps–computed without external backbones, including standard query-key and key-key attention, similarity maps from RADIO's intermediate blocks, and similarity maps derived from adapted DINOv2 and SAM patch tokens $\tilde{T}_{\text{patch}}^{(\text{dinov2})}, \tilde{T}_{\text{patch}}^{(\text{sam})}$. We find that RADIO's last block (Block

---

[1]A notable exception is that CLIP patch-wise and CLS outputs have different embedding dimensions as the CLS output token goes through another projection to be aligned to language.

$l$) offers the most stable spatial locality, providing a strong foundation for efficient zero-shot segmentation without external augmentations.

Building on these findings, we propose a novel attention enhancement method, called **self-correlating recursive attention (SCRA)**. As shown in Fig. 2, SCRA takes the output tokens $T_l \in \mathbb{R}^{N \times D}$ of the model, and computes a normalized correlation (cosine similarity) matrix $C = \check{T}_l \check{T}_l^\top$, where $\check{T}_l = T_l / \|T_l\|_2$. Cells with negative correlation between patches are set to $-\infty$, forcing patches to attend only to similar ones [20]. More formally,

$$T_l = \mathrm{MLP}_l\left(\mathrm{SA}\left(T_{l-1}\right)\right), \quad \check{T}_l = T_l / \|T_l\|_2 \quad (4)$$

$$C_{ij} = \begin{cases} \check{T}_l^{(i)} \cdot \check{T}_l^{(j)}, & \text{if } \check{T}_l^{(i)} \cdot \check{T}_l^{(j)} \geq 0 \\ -\infty, & \text{otherwise} \end{cases} \quad (5)$$

A softmax is then applied across each column of a scaled $C$ to ensure that each patch's attention sums to one, keeping the resulting features within the convex hull of the originals. Finally, the last block's activations are recomputed using these self-correlation weights. SCRA computes its output $T_{\mathrm{scra}}$ using RADIO's last block output $T_l$ as follows[1]:

$$T_{\mathrm{scra}} = \mathrm{MLP}_l\left(\mathrm{softmax}(\tau_{\mathrm{scra}} C) \cdot \mathrm{MLP}_V^l(T_l)\right) \quad (6)$$

where $\tau_{\mathrm{scra}} > 1$ is a temperature scaling factor that sharpens attention to only very similar patches. SCRA is able to improve the spatial-locality of RADIO features without introducing additional parameters and with minimal computational overhead of computing Eq. (5) and Eq. (6).

### 3.4. Self-Correlating Global Aggregation

When performing sliding window inference, we aggregate windows in feature space, which introduces noisy artifacts [36] demonstrated in Fig. 2. We propose a simple and effective **self-correlating global aggregation (SCGA)**. SCGA computes a self-correlation matrix from the aggregated feature map $G = \check{T}_{\mathrm{agg}} \check{T}_{\mathrm{agg}}^\top$, where $\check{T}_{\mathrm{agg}} = T_{\mathrm{agg}} / \|T_{\mathrm{agg}}\|_2$. Similar to SCRA, cells with negative correlation between patches are set to $-\infty$, a softmax is applied to each column of a scaled $G$. The correlation matrix $G$ is used as attention to aggregate the tokens:

$$G_{ij} = \begin{cases} \check{T}_{\mathrm{agg}}^{(i)} \cdot \check{T}_{\mathrm{agg}}^{(j)}, & \text{if } \check{T}_{\mathrm{agg}}^{(i)} \cdot \check{T}_{\mathrm{agg}}^{(j)} \geq 0 \\ -\infty, & \text{otherwise} \end{cases} \quad (7)$$

$$T_{\mathrm{scga}} = \mathrm{softmax}(\tau_{\mathrm{scga}} G) \cdot T_{\mathrm{agg}} \quad (8)$$

Qualitatively, Fig. 2 shows how SCGA suppresses the noise and windowing artifacts, while maintaining the sharpness of the feature map, all while not requiring additional backbones and maintaining efficiency.

---

[1]Ignoring standard ViT layer norms, multiple attention heads, and residual connections for simplicity.

Together, dense language alignment, SCRA, and SCGA constitute the core components of **RADSeg**. In the next subsection, we introduce an optional mask-refinement stage that extends the framework to **RADSeg+**.

### 3.5. RADIO-SAM Refinement

Most OVSS methods rely on ViT backbones, producing coarse feature maps with 14-16x downsampling. Even with overlapping windows, this resolution gap remains. Prior work closes this gap either with RGB-based refiners such as PAMR [2], which are computationally heavy, or with SAM-based refinement [22, 36], which improves quality but requires an additional backbone.

Within the RADIO framework, however, we can **leverage SAM-huge at only a +20M increase in parameters**. RADIO readily provides an adaptor that maps its output patch tokens $T_{l,\mathrm{patch}}$ into the SAM-huge token space $\tilde{T}_{\mathrm{patch}}^{(\mathrm{sam})} = \mathrm{MLP}_{\mathrm{patch}}^{(\mathrm{sam})}(T_{l,\mathrm{patch}})$. As a result, we only need the decoder and prompt encoders of SAM-huge, which constitute **0.7% of the total SAM-huge parameters**.

Following SAM-based refinement strategies [36, 42], we generate point, box, and mask prompts from the coarse predictions of $T_{\mathrm{scga}}$. Rather than simply reusing $T_{\mathrm{scga}}$ features, which we find degrades accuracy, we re-encode the full image with RADIO+SCRA (no sliding windows) to feed into SAM. As shown in Fig. 2, RADIO-SAM refinement sharpens boundaries and yields better final masks.

## 4. Experiments

### 4.1. 2D Zero-Shot Open Vocabulary Segmentation

**Benchmark Settings.** Consistent with previous methods [4, 13, 36, 42], we evaluate our method across five widely used 2D semantic segmentation benchmarks: PASCAL VOC [11], PASCAL Context [25], COCO-Stuff [5], Cityscapes [7], and ADE20K [48]. In prior work, different datasets and baselines were often evaluated at different resolutions, making it diffcult to disentangle performance gains from higher resolutions versus improvements due to the method itself. Thus, for a comprehensive evaluation, we evaluate three representative resolution standards set by previous methods: low (336-560) [13, 39, 44], mid (336-688) [36], and high (672-1344) [42], with details per dataset in the supplementary. We do not allow any method, to access a higher input resolution than its setting, yet, it can operate on lower resolutions. We do not employ heavy post-processing (e.g PAMR [2]) for any of the methods.

We follow standard OVSS protocols, each class is embedded in prompt templates, encoded with the text encoder, and used to obtain segmentation masks through cosine similarity with the feature maps. We ignore the ambiguous 'background' label. We report mean Intersection over Union (mIoU) on the validation split of each dataset.
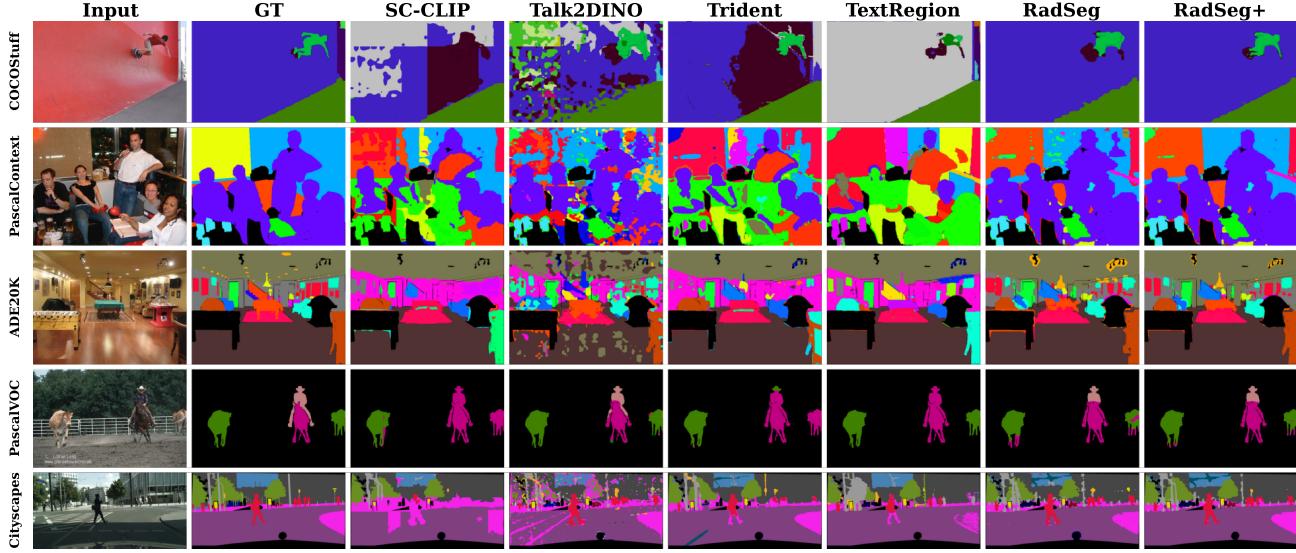
Figure 4. **Qualitative 2D Open-Vocabulary Semantic Segmentation Results.** For each of the five benchmark datasets, we show a representative example and compare **RADSeg** and **RADSeg+** with competitive baselines (SC-CLIP, Talk2DINO, Trident, and TextRegion). Both **RADSeg** and **RADSeg+** produce noticeably clearer and more accurate segmentation maps across all cases.

Table 1. 2D zero-shot open-vocabulary semantic segmentation mIoU results across resolution limits and model sizes. Ranking shown as first, second, and third. One ranking for base and huge is reported to highlight that **RADSeg**-base is superior to huge baselines. Underlined cells show results obtained at lower resolution. Numbers at specific resolutions are included in the supplementary.

| Methods | ≤ Low Resolution | | | | | | ≤ Mid Resolution | | | | | | ≤ High Resolution | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | CTX | VOC | Stuff | ADE | City | Avg | CTX | VOC | Stuff | ADE | City | Avg | CTX | VOC | Stuff | ADE | City | Avg |
| **Base Models** | | | | | | | | | | | | | | | | | | |
| NACLIP [13] | 35.17 | 79.71 | 23.30 | 17.42 | 35.49 | 38.22 | 35.17 | 80.29 | 23.64 | 17.78 | 35.98 | 38.48 | 35.17 | 80.29 | 23.64 | 17.80 | 36.74 | 38.48 |
| ResCLIP [44] | 36.80 | 82.32 | 24.70 | 18.03 | 35.85 | 39.54 | 36.80 | 85.89 | 25.07 | 18.55 | 36.19 | 40.49 | 36.80 | 85.89 | 25.07 | 18.57 | 36.88 | 40.49 |
| RayFronts [1] | 36.81 | 72.59 | 25.34 | 23.38 | 36.29 | 38.88 | 38.07 | 80.12 | 27.39 | 24.63 | 40.47 | 42.14 | 38.07 | 80.12 | 27.39 | 24.63 | 40.47 | 42.14 |
| ProxyCLIP [20] | 38.74 | 78.18 | 26.11 | 19.71 | 39.69 | 40.49 | 38.74 | 80.31 | 26.41 | 19.71 | 40.39 | 40.93 | 38.74 | 80.31 | 26.41 | 19.71 | 40.39 | 40.93 |
| SC-CLIP [3] | 40.12 | 84.29 | 26.62 | 20.06 | 41.02 | 42.42 | 40.12 | 87.67 | 27.25 | 20.68 | 41.02 | 43.19 | 40.12 | 87.67 | 27.25 | 20.68 | 41.24 | 43.19 |
| Talk2Dino [4] | 39.43 | 85.68 | 27.43 | 20.42 | 38.05 | 42.20 | 40.31 | 85.68 | 27.89 | 21.67 | 39.47 | 43.00 | 40.31 | 85.68 | 27.89 | 21.70 | 41.15 | 43.00 |
| Trident [36] | 41.62 | 83.74 | 28.22 | 21.17 | 41.86 | 43.32 | 42.16 | 84.50 | 28.24 | 21.98 | 42.08 | 43.79 | 42.16 | 84.50 | 28.24 | 21.98 | 42.19 | 43.79 |
| TextRegion [42] | 41.53 | 83.83 | 27.39 | 21.21 | 40.49 | 43.17 | 42.29 | 84.21 | 27.39 | 21.74 | 41.26 | 43.33 | 42.29 | 84.21 | 28.62 | 22.55 | 42.15 | 43.88 |
| **RADSeg** | 44.49 | 87.24 | 29.79 | 27.16 | 42.04 | 46.14 | 45.64 | 89.28 | 30.76 | 28.96 | 45.35 | 48.00 | 45.64 | 89.28 | 30.76 | 28.96 | 48.79 | 48.00 |
| **RADSeg+** | 48.23 | 88.69 | 31.66 | 29.44 | 45.59 | 48.72 | 48.59 | 90.35 | 32.45 | 30.86 | 47.82 | 50.01 | 48.59 | 90.35 | 32.45 | 30.86 | 50.72 | 50.01 |
| **Huge Models** | | | | | | | | | | | | | | | | | | |
| RayFronts [1] | 31.67 | 72.59 | 23.31 | 20.46 | 28.98 | 35.40 | 32.29 | 73.36 | 23.44 | 21.19 | 31.84 | 36.42 | 32.29 | 73.36 | 23.44 | 21.19 | 34.27 | 36.42 |
| ProxyCLIP [20] | 39.16 | 78.02 | 26.19 | 23.90 | 43.64 | 42.18 | 39.16 | 83.03 | 27.76 | 24.05 | 43.92 | 43.21 | 39.16 | 83.03 | 27.76 | 24.05 | 43.92 | 43.21 |
| Trident [36] | 43.16 | 87.97 | 28.55 | 25.64 | 46.87 | 46.44 | 44.32 | 88.67 | 28.55 | 26.70 | 46.87 | 46.90 | 44.32 | 88.67 | 28.55 | 27.02 | 47.34 | 46.90 |
| TextRegion [42] | 44.04 | 89.53 | 30.19 | 24.53 | 47.35 | 47.13 | 44.76 | 89.53 | 30.19 | 26.42 | 47.35 | 47.50 | 46.13 | 89.53 | 31.22 | 27.30 | 47.35 | 48.18 |
| **RADSeg** | 42.27 | 89.58 | 28.30 | 25.96 | 38.85 | 44.99 | 44.80 | 89.74 | 28.93 | 28.21 | 42.93 | 46.92 | 45.01 | 89.74 | 29.46 | 28.21 | 47.75 | 47.58 |
| **RADSeg+** | 45.44 | 90.04 | 30.12 | 27.75 | 41.66 | 47.00 | 47.74 | 90.14 | 30.44 | 29.92 | 45.78 | 48.80 | 47.92 | 90.14 | 30.84 | 29.92 | 50.01 | 49.57 |

We report our two methods, **RADSeg** and **RADSeg+**, without and with RADIO-SAM refinement. We evaluate both the base and huge variants of RADIOv3, using SigLIP2 as the feature adaptor with $\tau_{\mathrm{scra}} = 10, \tau_{\mathrm{scga}} = 10$.

**Baselines**. We compare our approach against a diverse set of SOTA training-free OVSS methods. This includes CLIP-based adaptations such as NACLIP [13], ResCLIP [44], and SC-CLIP [3], which refine attention blocks within the CLIP architecture; hybrid multi-model fusion methods such as ProxyCLIP [20], Talk2DINO [4], Trident [36], and TextRegion [42], which leverage multiple vision foundation models including CLIP, DINO, and SAM; and the only RADIO-based method Rayfronts [1].

**Quantitative and Qualitative Results.** Tab. 1 shows how well methods perform within a resolution limit across datasets. **RADSeg** establishes a clear performance margin over prior zero-shot OVSS methods. In the *base* setting, **RADSeg** achieves the highest average mIoU at all resolu-
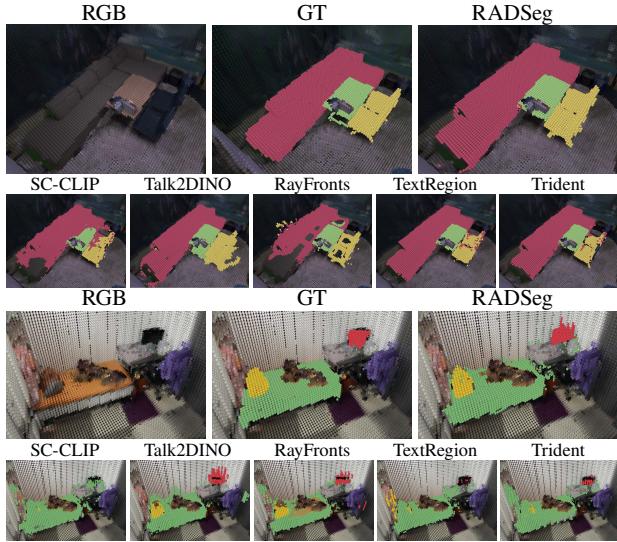
6

Figure 5. **Qualitative 3D Open-Vocabulary Semantic Segmentation Results.** We show two scenes: one from Replica ("chair", "table", "couch" classes), and one from ScanNet++ ("bed", "pillow", "monitor" classes). Segmented voxels are overlaid on the RGB for visualization. Across all 3D baselines, **RADSeg** provides more accurate segmentations with far fewer outlier voxels.

tion budgets, outperforming the closest baselines by **+3–5% mIoU**. **RADSeg+** strengthens this lead further, improving average mIoU by up to **+7% mIoU** over the next best baseline. Remarkably, our base model already surpasses several huge-model baselines, underscoring the effectiveness of RADIO and our pipeline.

We show qualitative comparisons of **RADSeg** and **RADSeg+** against baselines, using base models, in Fig. 4. **RADSeg** demonstrates superior semantic consistency and more precise segmentation boundaries, especially evident in challenging categories like "person" and in complex scenes in Pascal Context and Cityscapes.

### 4.2. 3D Zero-Shot Open Vocabulary Segmentation

While 2D segmentation evaluates per-frame semantic understanding, it fails to capture multi-view semantic consistency, essential for embodied perception. Thus, we extend our evaluation to 3D, examining how different encoders generalize within a 3D reconstruction pipeline. We focus on evaluating encoders rather than 3D scene representations. Thus, we adopt a simple setup across methods, in which 2D outputs are unprojected onto point clouds and aggregated within voxels, averaging coincident point features.

**Benchmark Settings.** We test two 2D-to-3D lifting strategies: projecting the features, or projecting the segmentation probabilities. Probability projection prevents feature collisions and supports post-segmentation refinement, yet forces closed-set outputs. Feature projections retain open-vocabulary flexibility at the cost of possible aliasing.

We evaluate on Replica [37] and ScanNet [8] following prior work [1], and additionally on ScanNet++ [45]. For a fair comparison, we extend the top-performing 2D baselines Trident, SC-CLIP, TextRegion, Talk2DINO, and RayFronts, to 3D using the same projection and fusion procedure. Since 3D dense feature maps require much more GPU memory, we only evaluate base models for all baselines and our method.

**Quantitative and Qualitative Results.** Across all benchmarks, **RADSeg** establishes a clear performance margin over the baselines. Under feature-space aggregation, our approach achieves the highest mIoU on all the datasets, improving over the strongest baselines by **+3–5% mIoU** on average. When evaluated in probability space, which allows for the evaluation of **RADSeg+**, it again ranks first or second in every setting: without refinement, it secures the **second-highest** metrics across datasets, while with refinement it achieves SOTA results in terms of mIoU/f-mIoU with **33.05%/59.21%** on Replica, **38.58%/48.49%** on ScanNet, and **33.27%/48.35%** on ScanNet++. Notably, while RayFronts struggles in 2D, its use of RADIO allows it to be competitive in 3D, showing the importance of a complementary 3D evaluation. Overall, **RADSeg** consistently outperforms CLIP-based, DINO-based, and multi-model baselines showing its strong multi-view consistency.

Fig. 5 presents qualitative comparisons between **RADSeg** and all the baselines in our evaluation. Across a diverse set of queries, **RADSeg** consistently yields cleaner and more semantically coherent predictions, producing significantly fewer visual outliers than competing methods.

### 4.3. Ablation Studies

**Exploring RADIO's language alignment.** Tab. 3 summarizes the performance of different RADIO language adapters. Following RADIO's training regimen of projecting RADIO patch features to SigLIP/CLIP patch features shows no language alignment with a dimensionality mismatch in case of CLIP, while using CLS adaptors on patch features yields high OVSS mIoU (34-44%) across datasets, highlighting RADIO's emergent dense language alignment. RADIOv3 with SigLIP2 CLS shows the highest alignment which we adopt in our method. Notably, RADIO's language alignment excels with base model sizes, making it all the more suitable for efficient zero-shot OVSS.

**Effectiveness of RADSeg components**. Tab. 4 summarizes 2D OVSS performance across datasets for evaluating our attention and refinement components. Notably, SCRA and SCGA provide consistent improvements with an average +2%, and +1.4% respectively, at negligible computational overhead. RADIO-SAM refinement improves by another +2% yet slows down the pipeline by 3×. Regardless, **RADSeg+** (which includes RADIO-SAM refinement) remains significantly faster than competing methods.

Table 2. 3D zero-shot open-vocabulary semantic segmentation mIoU results. Ranking shown as <mark>first</mark>, <mark>second</mark>, and <mark>third</mark>. **RADSeg** consistently shows superior performance across datasets emphasizing its multi-view consistency.

| | Feature Space 3D Aggregation | | | | | | Probability Space 3D Aggregation | | | | | |
| | Replica [37] | | ScanNet [8] | | ScanNet++ [45] | | Replica [37] | | ScanNet [8] | | ScanNet++ [45] | |
| Methods | mIoU | f-mIoU | mIoU | f-mIoU | mIoU | f-mIoU | mIoU | f-mIoU | mIoU | f-mIoU | mIoU | f-mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SC-CLIP [3] | 21.33 | 43.57 | 28.04 | 40.57 | 19.24 | 37.68 | 21.98 | 44.31 | 28.59 | 40.89 | 20.97 | 38.10 |
| Talk2Dino [4] | 23.12 | 35.07 | 28.58 | 32.88 | 24.88 | 30.96 | 23.94 | 36.86 | 29.13 | 34.47 | 25.43 | 32.37 |
| Trident [36] | 23.62 | 47.18 | 32.23 | 43.12 | 23.11 | 41.56 | 25.44 | 47.48 | 32.78 | 43.67 | 24.21 | 42.23 |
| Trident (w/refine) | – | – | – | – | – | – | 27.38 | 49.81 | 32.63 | 44.17 | 26.24 | 44.51 |
| TextRegion [42] | 27.15 | 50.50 | 34.42 | 44.47 | 24.08 | 45.86 | 28.35 | 50.90 | 35.83 | 46.09 | 24.59 | 45.91 |
| RayFronts [1] | 28.39 | 52.05 | 31.02 | 42.14 | 26.34 | 38.91 | 31.28 | 53.28 | 31.05 | 42.30 | 26.86 | 40.55 |
| **RADSeg** | 32.29 | 58.96 | 36.29 | 45.64 | 31.95 | 46.53 | 32.68 | 57.02 | 36.41 | 46.92 | 32.9 | 47.76 |
| **RADSeg+** | – | – | – | – | – | – | 33.05 | 59.21 | 38.58 | 48.49 | 33.27 | 48.35 |

Table 3. Ablation of different RADIO language adapters across model sizes. Mid resolution is used. Average mIoU reported across 2D datasets.

| Methods | Base | Large | Huge |
|---|---|---|---|
| Rv2.5-SigLIP$_{cls}$ | 39.96 | **38.55** | 36.97 |
| Rv2.5-SigLIP$_{patch}$ | 0.32 | 0.31 | 0.36 |
| Rv2.5-CLIP$_{cls}$ | 39.43 | 37.79 | 34.85 |
| Rv2.5-CLIP$_{patch}$ | NA | NA | NA |
| Rv3-SigLIP2$_{cls}$ | **44.62** | 0.86 | **41.24** |
| Rv3-SigLIP2$_{patch}$ | 0.73 | 0.76 | 0.25 |
| Rv3-CLIP$_{cls}$ | 42.32 | 0.75 | 36.79 |
| Rv3-CLIP$_{patch}$ | NA | NA | NA |

Table 4. Ablation study of components of **RADSeg** on mid resolution. Average mIoU and latency on a V100 across datasets.

| Methods | mIoU (%) | | | | | | (s) |
| | CTX | VOC | Stuff | ADE | City | Avg | |
|---|---|---|---|---|---|---|---|
| Base | 40.4 | 88.07 | 27.41 | 27.3 | 39.94 | 44.62 | 0.107 |
| + SCRA | 43.07 | 88.03 | 29.87 | 27.69 | 44.5 | 46.63 | 0.111 |
| + SCGA | 45.64 | 89.28 | 30.76 | 28.96 | 45.35 | 48.00 | 0.115 |
| + R-SAM | **48.48** | **90.14** | **32.48** | **30.83** | **48.10** | **50.01** | 0.354 |

Table 5. Ablation study showing RADIO's ability to empower existing approaches. Average mIoU across datasets, number of parameters for base models and latency on a V100. RADIO can unlock efficiency and mIoU gains for different baselines

| Methods | mIoU (%) | Params (M) | Latency (s) |
|---|---|---|---|
| NACLIP [13] | 38.48 | 86M | 0.089 |
| NARADIO | 44.54 | 105M | 0.107 |
| ProxyCLIP [20] | 40.87 | 171M | 0.667 |
| ProxyRADIO-D | 45.44 | 106M | 0.109 |
| ProxyRADIO-S | 46.05 | 105M | 0.109 |
| ResCLIP [44] | 40.33 | 86M | 0.261 |
| ResRADIO | 45.18 | 105M | 0.123 |
| **RADSeg** | **48.00** | 105M | 0.115 |

**Can RADIO improve baseline approaches?** To demonstrate RADIO's broader applicability in OVSS beyond **RADSeg**, we adapt some of the previous approaches

that can readily benefit from the backbone. Specifically, we use RADIOv3 with adapted SigLIP2 features and incorporate neighbor aware [13], residual [44], and proxy attention [20]. Tab. 5 highlights how replacing CLIP with RADIO can provide significant consistent improvements in mIoU, and efficiency, **showing RADIO's potential as a backbone for OVSS**.

### 4.4. Parameter and Compute Efficiency

To better contextualize the segmentation accuracy improvements, we perform a thorough analysis of each method's latency and parameter efficiency. For a holistic view that reflects all computations that enabled a method to obtain its mIoU results, we report average latency across all datasets at mid resolution. As illustrated in Fig. 1, **RADSeg**-base surpasses the strongest baselines, TextRegion Huge and Trident Huge, achieving 9.27× lower latency while requiring 8.57× fewer parameters. Furthermore, **RADSeg**+, which enhances **RADSeg** with SAM-based refinement, delivers more fine-grained segmentation maps while remaining 3× faster and using 7.26× fewer parameters.

### 5. Conclusion

We present **RADSeg**, a dense language-aligned encoder that leverages the RADIO agglomerative framework. **RADSeg** effectively employs self correlations for attention and window aggregation, and RADIO SAM adapted features for efficient mask refinement. We evaluate **RADSeg** on 8 datasets in 2D and 3D showing significant improvements in zero-shot open-vocabulary semantic segmentation (ZSOVSS) across mIoU, latency, and parameter efficiency. In addition, we provide the first empirical study on the use of RADIO for ZSOVSS highlighting its efficiency and strong emergent dense language alignment. We expect this work to motivate broader exploration of agglomerative models towards efficient, accurate, and generalizable open-vocabulary semantic segmentation. Our future work aims to incorporate instance differentiation, and multi-label segmentation.

# References

[1] Omar Alama, Avigyan Bhattacharya, Haoyang He, Seungchan Kim, Yuheng Qiu, Wenshan Wang, Cherie Ho, Nikhil Keetha, and Sebastian Scherer. Rayfronts: Open-set semantic ray frontiers for online scene understanding and exploration. *arXiv preprint arXiv:2504.06994*, 2025. 2, 3, 4, 6, 7, 8, 5

[2] Nikita Araslanov and Stefan Roth. Single-stage semantic segmentation from image labels. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4253–4262, 2020. 5

[3] Sule Bai, Yong Liu, Yifei Han, Haoji Zhang, and Yansong Tang. Self-calibrated clip for training-free open-vocabulary segmentation. *arXiv preprint arXiv:2411.15869*, 2024. 3, 4, 6, 8, 2, 5

[4] Luca Barsellotti, Lorenzo Bianchi, Nicola Messina, Fabio Carrara, Marcella Cornia, Lorenzo Baraldi, Fabrizio Falchi, and Rita Cucchiara. Talking to dino: Bridging self-supervised vision backbones with language for open-vocabulary segmentation. *arXiv preprint arXiv:2411.19331*, 2024. 5, 6, 8, 1, 2

[5] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Cocostuff: Thing and stuff classes in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1209–1218, 2018. 2, 5

[6] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 3

[7] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 2, 5

[8] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, pages 5828–5839, 2017. 7, 8

[9] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. *arXiv preprint arXiv:2309.16588*, 2023. 4

[10] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 3

[11] Mark Everingham, Luc Gool, Christopher K. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *Int. J. Comput. Vision*, 88(2): 303–338, 2010. 2, 5

[12] Qiao Gu, Ali Kuwajerwala, Sacha Morin, Krishna Murthy Jatavallabhula, Bipasha Sen, Aditya Agarwal, Corban Rivera, William Paul, Kirsty Ellis, Rama Chellappa, et al. Conceptgraphs: Open-vocabulary 3d scene graphs for perception and planning. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5021–5028. IEEE, 2024. 3, 2

[13] Sina Hajimiri, Ismail Ben Ayed, and Jose Dolz. Pay attention to your neighbours: Training-free open-vocabulary semantic segmentation. In *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 5061–5071. IEEE, 2025. 2, 3, 4, 5, 6, 8, 1

[14] Greg Heinrich, Mike Ranzinger, Hongxu Yin, Yao Lu, Jan Kautz, Andrew Tao, Bryan Catanzaro, and Pavlo Molchanov. Radiov2.5: Improved baselines for agglomerative vision foundation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22487–22497, 2025. 3

[15] Krishna Murthy Jatavallabhula, Alihusein Kuwajerwala, Qiao Gu, Mohd Omama, Tao Chen, Shuang Li, Ganesh Iyer, Soroush Saryazdi, Nikhil Keetha, Ayush Tewari, Joshua B. Tenenbaum, Celso Miguel de Melo, Madhava Krishna, Liam Paull, Florian Shkurti, and Antonio Torralba. Conceptfusion: Open-set multimodal 3d mapping. *Robotics: Science and Systems (RSS)*, 2023. 2, 3

[16] Kim Jun-Seong, GeonU Kim, Kim Yu-Ji, Yu-Chiang Frank Wang, Jaesung Choe, and Tae-Hyun Oh. Dr. splat: Directly referring 3d gaussian splatting via direct language embedding registration. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 14137–14146, 2025. 3

[17] Chung Min Kim, Mingxuan Wu, Justin Kerr, Ken Goldberg, Matthew Tancik, and Angjoo Kanazawa. Garfield: Group anything with radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21530–21539, 2024. 3

[18] Seungchan Kim, Omar Alama, Dmytro Kurdydyk, John Keller, Nikhil Keetha, Wenshan Wang, Yonatan Bisk, and Sebastian Scherer. Raven: Resilient aerial navigation via open-set semantic memory and behavior adaptation. *arXiv preprint arXiv:2509.23563*, 2025. 1

[19] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 2, 3

[20] Mengcheng Lan, Chaofeng Chen, Yiping Ke, Xinjiang Wang, Litong Feng, and Wayne Zhang. Proxyclip: Proxy attention improves clip for open-vocabulary segmentation. In *European Conference on Computer Vision*, pages 70–88. Springer, 2024. 2, 3, 4, 5, 6, 8

[21] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. *arXiv preprint arXiv:2201.03546*, 2022. 2

[22] Yuqi Lin, Hengjia Li, Wenqi Shao, Zheng Yang, Jun Zhao, Xiaofei He, Ping Luo, and Kaipeng Zhang. Samrefiner: Taming segment anything model for universal mask refinement. *arXiv preprint arXiv:2502.06756*, 2025. 5

[23] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 3

9

[24] Andrew F Luo, Jacob Yeung, Rushikesh Zawar, Shaurya Dewan, Margaret M Henderson, Leila Wehbe, and Michael J Tarr. Brain mapping with dense features: Grounding cortical semantic selectivity in natural images with vision transformers. *arXiv preprint arXiv:2410.05266*, 2024. 1

[25] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014. 2, 5

[26] 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. 2, 3

[27] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021. 1, 2, 3

[28] Mike Ranzinger, Greg Heinrich, Jan Kautz, and Pavlo Molchanov. Am-radio: Agglomerative vision foundation model reduce all domains into one. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12490–12500, 2024. 2, 3

[29] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou, and Jiwen Lu. Denseclip: Language-guided dense prediction with context-aware prompting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18082–18091, 2022. 2

[30] Adam Rashid, Satvik Sharma, Chung Min Kim, Justin Kerr, Lawrence Yunliang Chen, Angjoo Kanazawa, and Ken Goldberg. Language embedded radiance fields for zero-shot task-oriented grasping. In *7th Annual Conference on Robot Learning*, 2023. 1

[31] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 3

[32] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294, 2022. 2

[33] Nur Muhammad Mahi Shafiullah, Chris Paxton, Lerrel Pinto, Soumith Chintala, and Arthur Szlam. Clip-fields: Weakly supervised semantic fields for robotic memory. *arXiv preprint arXiv:2210.05663*, 2022. 3

[34] Jinghuan Shang, Karl Schmeckpeper, Brandon B May, Maria Vittoria Minniti, Tarik Kelestemur, David Watkins, and Laura Herlant. Theia: Distilling diverse vision foundation models for robot learning. *arXiv preprint arXiv:2407.20179*, 2024. 4

[35] Jin-Chuan Shi, Miao Wang, Hao-Bin Duan, and Shao-Hua Guan. Language embedded 3d gaussians for open-vocabulary scene understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5333–5343, 2024. 3

[36] 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. 2, 3, 4, 5, 6, 8, 1

[37] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. The replica dataset: A digital replica of indoor spaces. *arXiv:1906.05797*, 2019. 7, 8

[38] Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, et al. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. *arXiv preprint arXiv:2502.14786*, 2025. 2

[39] Feng Wang, Jieru Mei, and Alan Yuille. Sclip: Rethinking self-attention for dense vision-language inference. In *European Conference on Computer Vision*, pages 315–332. Springer, 2024. 2, 3, 4, 5, 1

[40] Abdelrhman Werby, Chenguang Huang, Martin Büchner, Abhinav Valada, and Wolfram Burgard. Hierarchical open-vocabulary 3d scene graphs for language-grounded robot navigation. In *First Workshop on Vision-Language Models for Navigation and Manipulation at ICRA 2024*, 2024. 3, 2

[41] Monika Wysoczańska, Michaël Ramamonjisoa, Tomasz Trzciński, and Oriane Siméoni. Clip-diy: Clip dense inference yields open-vocabulary semantic segmentation for-free. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1403–1413, 2024. 2

[42] Yao Xiao, Qiqian Fu, Heyi Tao, Yuqun Wu, Zhen Zhu, and Derek Hoiem. Textregion: Text-aligned region tokens from frozen image-text models. *arXiv preprint arXiv:2505.23769*, 2025. 2, 3, 5, 6, 8, 1

[43] Mengde Xu, Zheng Zhang, Fangyun Wei, Han Hu, and Xiang Bai. Side adapter network for open-vocabulary semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2945–2954, 2023. 2

[44] Yuhang Yang, Jinhong Deng, Wen Li, and Lixin Duan. Resclip: Residual attention for training-free dense vision-language inference. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 29968–29978, 2025. 2, 3, 4, 5, 6, 8, 1

[45] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. Scannet++: A high-fidelity dataset of 3d indoor scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12–22, 2023. 7, 8

[46] Naoki Yokoyama, Sehoon Ha, Dhruv Batra, Jiuguang Wang, and Bernadette Bucher. Vlfm: Vision-language frontier maps for zero-shot semantic navigation. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 42–48. IEEE, 2024. 1

[47] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023. 2, 3

[48] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017. 2, 5

[49] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from clip. In *European conference on computer vision*, pages 696–712. Springer, 2022. 2

# RADSeg: Unleashing Parameter and Compute Efficient Zero-Shot Open-Vocabulary Segmentation Using Agglomerative Models

## Supplementary Material

## Limitations

While `RADSeg`-base delivers strong mIoU gains with a lightweight vision encoder, it relies on a huge text encoder, unlike CLIP-based baselines in the same model size class. However, text features are computed once per query, so their cost is amortized across many images. Furthermore, Although we demonstrate `RADSeg`'s strong performance on eight 2D and 3D OVSS datasets, current OVSS benchmarks remain limited. They primarily test whether each pixel is more similar to query A, B, **or** C. A more challenging setting would prevent access to all class labels and evaluate queries individually (is this pixel similar to query A). Moreover, support for multi-label segmentation is not explored. Developing and evaluating on such challenging settings is left for future work.

## Acknowledgments

## A1. Contribution Statement

**Omar Alama** Led and shaped the research, conceived the initial ideas for the encoder, wrote the majority of the manuscript, designed key figures and table formats. Additionally, developed the 3D evaluation pipeline and provided coding support and debugging throughout the project.

**Darshil Jariwala** Developed the encoder code and refined the initial ideas, conducted all evaluations and ablations for 2D OVSS – including porting and optimizing 2D baselines, generated qualitative 2D results, and helped with paper writing.

**Avigyan Bhattacharya** Conducted all evaluations for 3D OVSS, including porting 3D baselines, generated 3D qualitative figures, adapted ResCLIP to RADIO, wrote the 2D and 3D experimental sections. Contributed to paper writing and refinement.

**Seungchan Kim** Incorporated and processed the ScanNet++ dataset for 3D OVSS evaluation, and played a key role in refining the manuscript, improving clarity, structure, and presentation.

**Wenshan Wang** Provided valuable feedback on research design, manuscript writing, and method presentation.

**Sebastian Scherer** Provided valuable feedback on research direction, manuscript clarity, and figure presentation.

## A2. Additional 2D Evaluation Details

Table A.1. Adopted resolution and sliding window settings. Cell format: Shorter-side-Crop-Stride. TextRegion's stride always equals the crop size as per their method.

| Dataset | Low [13, 39, 44] | Mid [36] | High [42] |
|---------|------------------|----------|-----------|
| VOC | 336-224-112 | 336-336-112 | 672-336-336 |
| Stuff | 336-224-112 | 448-336-224 | 896-336-336 |
| CTX/ADE | 336-224-112 | 576-336-224 | 672-336-336 |
| City | 560-224-112 | 688-336-224 | 1344-336-336 |

**Emphasis on resolution standardization**. By examining the settings used for evaluation in previous works, we observe varying resolution standards, not only across methods and datasets, but even within different sub-modules of a method itself. These inconsistencies can conflate performance gains from increased resolutions with method improvements. For example, while NACLIP [13] and ResCLIP [44] choose the low resolution setting in their evaluations, Tab. 1 shows that they benefit from increased resolutions at the mid setting, making it unfair to compare their low-resolution performance with mid or high resolution performances of other approaches, which is done in many previous works [36, 42]. Furthermore, previous works have shown that increasing resolution can hurt performance [4, 36]. Thus, for a fair standard comparison, we rerun all baselines with 3 different representative resolution and sliding window standards (low, mid, high) shown in Tab. A.1. Note that evaluation resolution, at which we compute metrics, always remains unchanged. To avoid penalizing approaches that perform better when using lower resolutions, **we report the maximum performance at a resolution limit** in Tab. 1. Tab. 1 answers what the best performance is for a method given a **resolution limit**, while its source data, shown in Tab. A.2, answers what the performance is at a **particular resolution**. While Tab. 1 is more relevant for method comparisons, we report Tab. A.2 for completeness.

Table A.2. 2D zero-shot open-vocabulary semantic segmentation mIoU results across resolutions and model sizes. Ranking shown as first , second , and third . One ranking for base and huge is reported to highlight that `RADSeg`-base is superior to huge baselines.

| Methods | = Low Resolution | | | | | | = Mid Resolution | | | | | | = High Resolution | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | CTX | VOC | Stuff | ADE | City | Avg | CTX | VOC | Stuff | ADE | City | Avg | CTX | VOC | Stuff | ADE | City | Avg |
| **Base Models** | | | | | | | | | | | | | | | | | | |
| NACLIP [13] | 35.17 | 79.71 | 23.3 | 17.42 | 35.49 | 38.22 | 34.7 | 80.29 | 23.64 | 17.78 | 35.98 | 38.48 | 33.86 | 72.57 | 20.04 | 17.8 | 36.74 | 36.20 |
| ResCLIP [44] | 36.8 | 82.32 | 24.7 | 18.03 | 35.85 | 39.54 | 36.75 | 85.89 | 25.07 | 18.55 | 36.19 | 40.49 | 35.98 | 76.6 | 22.01 | 18.57 | 36.88 | 38.01 |
| RayFronts [1] | 36.81 | 72.59 | 25.34 | 23.38 | 36.29 | 38.88 | 38.07 | 80.12 | 27.39 | 24.63 | 40.47 | 42.14 | 36.99 | 67.39 | 23.92 | 24.18 | 39.13 | 38.32 |
| ProxyCLIP [20] | 38.74 | 78.18 | 26.11 | 19.71 | 39.69 | 40.49 | 37.88 | 80.31 | 26.41 | 19.67 | 40.39 | 40.93 | 34.85 | 70.44 | 21.63 | 19.18 | 38.84 | 36.99 |
| SC-CLIP [3] | 40.12 | 84.29 | 26.62 | 20.06 | 41.02 | 42.42 | 39.87 | 87.67 | 27.25 | 20.68 | 40.49 | 43.19 | 38.68 | 77.67 | 22.89 | 20.34 | 41.24 | 40.16 |
| Talk2Dino [4] | 39.43 | 85.68 | 27.43 | 20.42 | 38.05 | 42.20 | 40.31 | 85.66 | 27.89 | 21.67 | 39.47 | 43.00 | 40.23 | 84.15 | 27.06 | 21.70 | 41.15 | 42.86 |
| Trident [36] | 41.62 | 83.74 | 28.22 | 21.17 | 41.86 | 43.32 | 42.16 | 84.5 | 28.24 | 21.98 | 42.08 | 43.79 | 40.70 | 80.75 | 25.45 | 20.81 | 42.19 | 41.98 |
| TextRegion [42] | 41.53 | 83.83 | 27.39 | 21.21 | 40.49 | 43.17 | 42.29 | 84.21 | 27.13 | 21.74 | 41.26 | 43.33 | 42.88 | 83.19 | 28.62 | 22.55 | 42.15 | 43.88 |
| **RADSeg** | 44.49 | 87.24 | 29.79 | 27.16 | 42.04 | 46.14 | 45.64 | 89.28 | 30.76 | 28.96 | 45.35 | 48.00 | 45.14 | 84.07 | 29.05 | 28.93 | 48.79 | 47.20 |
| **RADSeg+** | 48.23 | 88.69 | 31.66 | 29.44 | 45.59 | 48.72 | 48.59 | 90.35 | 32.45 | 30.86 | 47.82 | 50.01 | 48.01 | 86.17 | 30.40 | 30.71 | 50.72 | 49.20 |
| **Huge Models** | | | | | | | | | | | | | | | | | | |
| RayFronts [1] | 31.67 | 72.59 | 23.31 | 20.46 | 28.98 | 35.40 | 32.29 | 73.36 | 23.44 | 21.19 | 31.84 | 36.42 | 32.19 | 67.64 | 22.13 | 21.08 | 34.27 | 35.46 |
| ProxyCLIP [20] | 39.16 | 78.02 | 26.19 | 23.90 | 43.64 | 42.18 | 38.31 | 83.03 | 27.76 | 24.05 | 43.92 | 43.21 | 37.03 | 71.66 | 21.81 | 23.65 | 43.09 | 39.45 |
| Trident [36] | 43.16 | 87.97 | 28.55 | 25.64 | 46.87 | 46.44 | 44.32 | 88.67 | 28.52 | 26.70 | 46.30 | 46.90 | 43.77 | 87.22 | 25.72 | 27.02 | 47.34 | 46.21 |
| TextRegion [42] | 44.04 | 89.53 | 30.19 | 24.53 | 47.35 | 47.13 | 44.76 | 89.42 | 29.85 | 26.42 | 47.04 | 47.50 | 46.13 | 89.36 | 31.22 | 27.30 | 46.88 | 48.18 |
| **RADSeg** | 42.27 | 89.58 | 28.3 | 25.96 | 38.85 | 44.99 | 44.8 | 89.74 | 28.93 | 28.21 | 42.93 | 46.92 | 45.01 | 88.52 | 29.46 | 27.15 | 47.75 | 47.58 |
| **RADSeg+** | 45.44 | 90.04 | 30.12 | 27.75 | 41.66 | 47.00 | 47.74 | 90.14 | 30.44 | 29.92 | 45.78 | 48.80 | 47.92 | 89.24 | 30.84 | 29.82 | 50.01 | 49.57 |

## A3. Additional 3D Evaluation Details

**3D mapping for evaluation.** We follow RayFronts' [1] approach to construct the 3D map and perform the evaluation. Given a pose $P_t \in SE(3)$, a corresponding depth map $D_t \in \mathbb{R}^{H \times W}$, and a feature map $F_t \in \mathbb{R}^{H \times W \times D}$, feature pixels are back projected to a 3D point cloud $P_t^{\text{local}} = \{(p_i, f_i)\}_{i=1}^{M}$, where $p_i \in \mathbb{R}^3$ denotes 3D position, and $f_i \in \mathbb{R}^{D+1}$ represents the concatenated feature vector, and hit count (initialized to 1 per point). Local updates are accumulated over frames and voxelized at a resolution of $\alpha$ to form the global semantic voxel map $P_t^{\text{global}} = \{(p_i, f_i)\}_{i=1}^{N}$. During voxelization, points that lie in the same voxel get their features averaged, and their hit counts summed to use as weights for subsequent averaging. In case of probability space 3D aggregation, the embedding dimension D in the feature map and semantic voxel map is equal to the number of classes as segmentation probabilities are projected instead of embeddings.

**Datasets.** Following prior works [1, 12, 15, 40] for our 3D OVSS evaluation, we choose the scenes, `office[0-4]`, `room[0-2]` from Replica and `scene[0011,0050,0231,0378,0518]` from ScanNet. Additionally, to assess performance on a cleaner real-world dataset, we evaluate on Scan-Net++ and select nine diverse scenes. The selected scenes are `scene[00777c41d4, bcd2436daf, a5114ca13d, 2b1dc6d6a5, d551dac194, f9f95681fd, ea42cd27e6, 20ff72df6e, bf6e439e38]`. We use all 101 classes from Replica, the standard ScanNet-to-NYU40 label mapping provided with the dataset itself for ScanNet, and the top 100 classes from ScanNet++ as defined in its official semantic segmentation
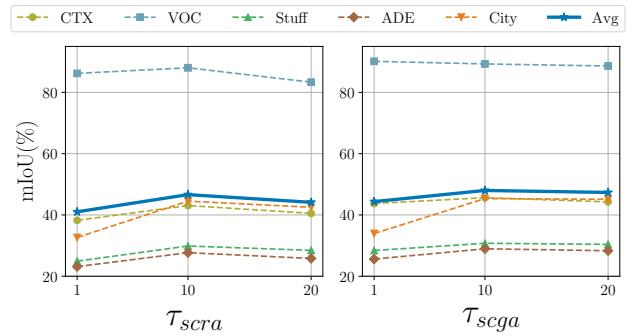


Figure A.1. Ablation study on different temperature factors $\tau_{scra}$ and $\tau_{scga}$. Left plot shows performance as we change $\tau_{scra}$ without SCGA. Right plot uses $\tau_{scra} = 10$ and varies $\tau_{scga}$. Overall $\tau_{scra} = \tau_{scga} = 10$ yield the best results on average.

benchmark. In ScanNet, we assign three of the forty classes ('otherprop,' 'otherstructure,' and 'otherfurniture') as ignore classes due to their ambiguity.

**Implementation details.** For all the datasets, we resize each image by setting the shorter side to 640 and keeping the original aspect ratio. A frame skip of 10 and 5cm voxels are used for constructing the 3D map. We use external ground-truth voxels for Replica (following [1], we use those provided by HOV-SG) and for ScanNet++. For Scan-Net, however, we derive the ground truth by lifting the 2D semantic segmentation annotations into 3D.

## A4. Additional Ablations and Detailed Tables

We ablate temperature parameters $\tau_{scra}$ and $\tau_{scga}$ for SCRA and SCGA at mid resolution across 2D datasets.

Table A.3. Expanded version of Tab. 3. Ablation of different RADIO language adapters across model sizes. Mid resolution is used. **RADIOv3-base with the SigLIP2 CLS adaptor has the best performance across datasets and model sizes**.

| Methods | Base | | | | | | Large | | | | | | Huge | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | CTX | VOC | Stuff | ADE | City | Avg | CTX | VOC | Stuff | ADE | City | Avg | CTX | VOC | Stuff | ADE | City | Avg |
| Rv2.5-SigLIP$_{cls}$ | 34.24 | 84.99 | 24.66 | 23.11 | 32.78 | 39.96 | 31.70 | 83.88 | 23.68 | 21.00 | 32.49 | 38.55 | 31.12 | 79.04 | 23.00 | 21.92 | 29.78 | 36.97 |
| Rv2.5-CLIP$_{cls}$ | 33.42 | 86.55 | 23.56 | 22.37 | 31.24 | 39.43 | 31.09 | 84.57 | 22.68 | 20.10 | 30.53 | 37.79 | 27.81 | 79.70 | 21.08 | 20.03 | 25.65 | 34.85 |
| Rv2.5-SigLIP$_{patch}$ | 0.31 | 1.00 | 0.06 | 0.14 | 0.09 | 0.32 | 0.27 | 0.97 | 0.06 | 0.12 | 0.14 | 0.31 | 0.37 | 1.06 | 0.06 | 0.17 | 0.13 | 0.36 |
| **Rv3-SigLIP$_{cls}$** | **40.40** | **88.07** | **27.41** | **27.30** | **39.94** | **44.62** | **1.00** | **1.46** | **0.52** | **0.15** | **1.18** | **0.86** | **35.23** | **88.35** | **24.55** | **24.34** | **33.74** | **41.24** |
| Rv3-CLIP$_{cls}$ | 37.90 | 86.89 | 26.36 | 24.48 | 35.99 | 42.32 | 0.67 | 1.71 | 0.07 | 0.12 | 1.17 | 0.75 | 29.48 | 84.07 | 21.85 | 21.39 | 27.16 | 36.79 |
| Rv3-SigLIP$_{patch}$ | 0.11 | 2.47 | 0.05 | 0.14 | 0.89 | 0.73 | 0.12 | 2.50 | 0.04 | 0.14 | 0.98 | 0.76 | 0.12 | 0.84 | 0.05 | 0.12 | 0.14 | 0.25 |

Table A.4. Expanded version of Tab. 5. Ablation study showing RADIOv3's ability to empower existing approaches. Mid resolution is used. ProxyRADIO-D/S refer to using DINOv2 and SAM adapted feature maps respectively for the proxy attention. **RADIO can improve mIoU for different CLIP-based baselines**.

| Methods | CTX | VOC | Stuff | ADE | City | Avg |
|---|---|---|---|---|---|---|
| NACLIP [13] | 34.70 | 80.28 | 23.65 | 17.78 | 35.98 | 38.48 |
| NARADIO | 40.85 | 86.25 | 27.93 | 27.19 | 40.47 | 44.54 |
| ProxyCLIP [20] | 37.76 | 80.24 | 26.4 | 19.65 | 40.28 | 40.87 |
| ProxyRADIO-D | 41.62 | 86.52 | 29.28 | 26.90 | 42.87 | 45.44 |
| ProxyRADIO-S | 41.99 | 87.90 | 29.33 | 27.00 | 44.04 | 46.05 |
| ResCLIP [44] | 36.52 | 85.79 | 24.91 | 18.43 | 36.02 | 40.33 |
| ResRADIO | 41.79 | 87.96 | 28.74 | 26.63 | 40.77 | 45.18 |
| **RADSeg** | **45.64** | **89.28** | **30.76** | **28.96** | **45.35** | **48.00** |

Fig. A.1 highlights the importance of scaling the attention correlation matrix to sharpen the attention to semantically similar patches and similarly for global aggregation. $\tau_{scra} = \tau_{scga} = 10$ yield the best results and is what we use for all experiments.

To demonstrate the effectiveness of the components of **RADSeg**, we augment Tab. 4 with a qualitative visualization. Fig. A.2 qualitatively illustrates the effectiveness of SCRA and SCGA in suppressing noise in the feature map and segmentation as well as reducing windowing artifacts. The effect is particularly pronounced in higher resolution datasets like Cityscapes. And while RADIO-SAM refinement is a post-segmentation process that cannot refine features, it is able to provide higher fidelity masks in many cases. Overall, the proposed **RADSeg** and **RADSeg**+ components demonstrate quantitative and qualitative improvements.

Finally, for reference, we provide per-dataset details for Tab. 3 in Tab. A.3 and per-dataset details for Tab. 5 in Tab. A.4.

## A5. Additional Efficiency Evaluation Details

**Improving baselines efficiency.** To have a robust comparison, we modify the inference code of NACLIP [13], ResCLIP [44], SC-CLIP [3], and ProxyCLIP [20] to use batched sliding window inference as opposed to iterating over each window. This significantly improves the baselines latency and gives us stronger comparison points. In our experiments, all baselines use batched sliding windows with a batch size equal to the number of windows.

**What latency to report ?** Input resolutions vary across datasets and even across samples since only the shorter image side is fixed. Some methods adjust hyperparameters per dataset, affecting total computation. Furthermore, the latency of methods (including **RADSeg**+) can vary with the content of an image as the number of masks varies. To provide a holistic measure that reflects all computations contributing to a method's final segmentation accuracy, we report the average latency across all validation samples in each dataset, as shown in Tab. A.5. The table reveals substantial latency differences between datasets, underscoring that measuring latency on a single image or at a fixed resolution does not capture overall performance.

Latencies are measured using mid resolution, FP32 precision, on a Tesla V100-32GB GPU. In addition, the number of parameters of the vision encoders of each method is reported. Fig. 1 visualizes the mIoU vs number of parameters and mIoU vs latency tradeoffs. Notably, at mid resolution, **RADSeg-base surpasses huge Trident and TextRegion baselines while being 9.3x-12.6x faster and having 8.1x-12.8x fewer parameters**.

## A6. Additional Qualitative Results

We provide additional 2D and 3D qualitative comparisons in Fig. A.3 and Fig. A.4, further illustrating the strengths of **RADSeg** over existing open-vocabulary segmentation methods. In 2D, **RADSeg** yields cleaner boundaries and more accurate segmentations in both cluttered indoor scenes and complex urban layouts. In 3D, it shows strong multi-view consistency, producing coherent semantic voxels with far fewer outliers and mislabeled regions than the baselines. These visualizations reinforce the ability of **RADSeg**, with its proposed components, to suppress noise and enhance object localizations across various scenarios.
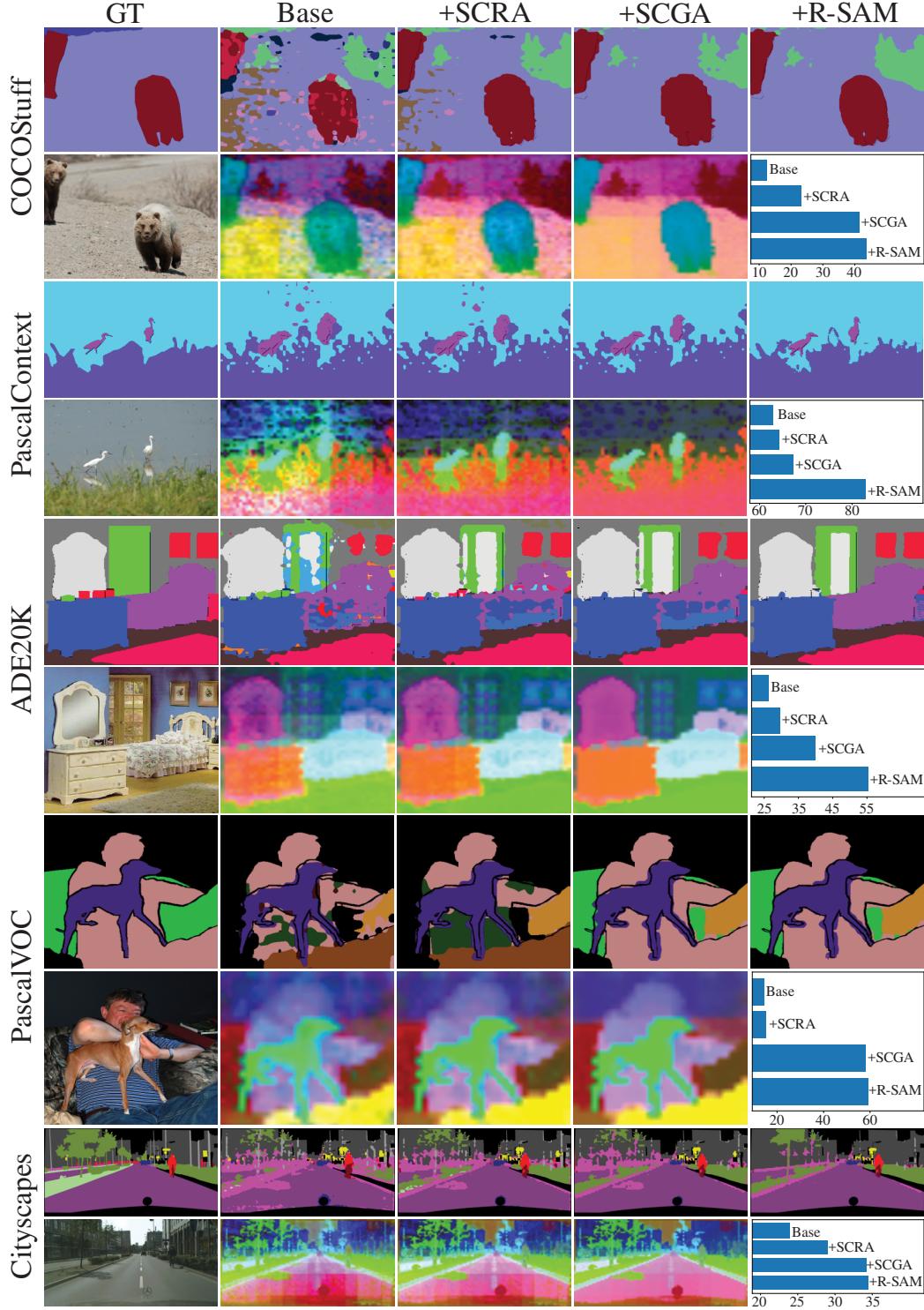
Figure A.2. Qualitative comparison of the contribution of each **RADSeg** component to feature and prediction quality. For each 2D dataset we show two rows: the first demonstrates how adding **RADSeg** components progressively improves segmentation output, while the second (Showing first 3 PCA components) illustrates how SCRA and SCGA enhance feature map quality and mitigate windowing artifacts. The accompanying bar plot links these visual trends to per-sample mIoU scores. Overall, the proposed **RADSeg** components yield clear improvements in both segmentation accuracy and feature-map fidelity.

Table A.5. Compute and parameter efficiency across datasets on a V100 GPU at mid resolution and FP32 precision. The table reports vision encoder parameters, average mIoU across datasets, and per-dataset latency to capture differences arising from varying input resolutions—both across and within datasets. Illustrated visually in Fig. 1. At mid resolution, **RADSeg-base surpasses huge Trident and TextRegion baselines while being 9.3x-12.6x faster and having 8.1x-12.8x fewer parameters**.

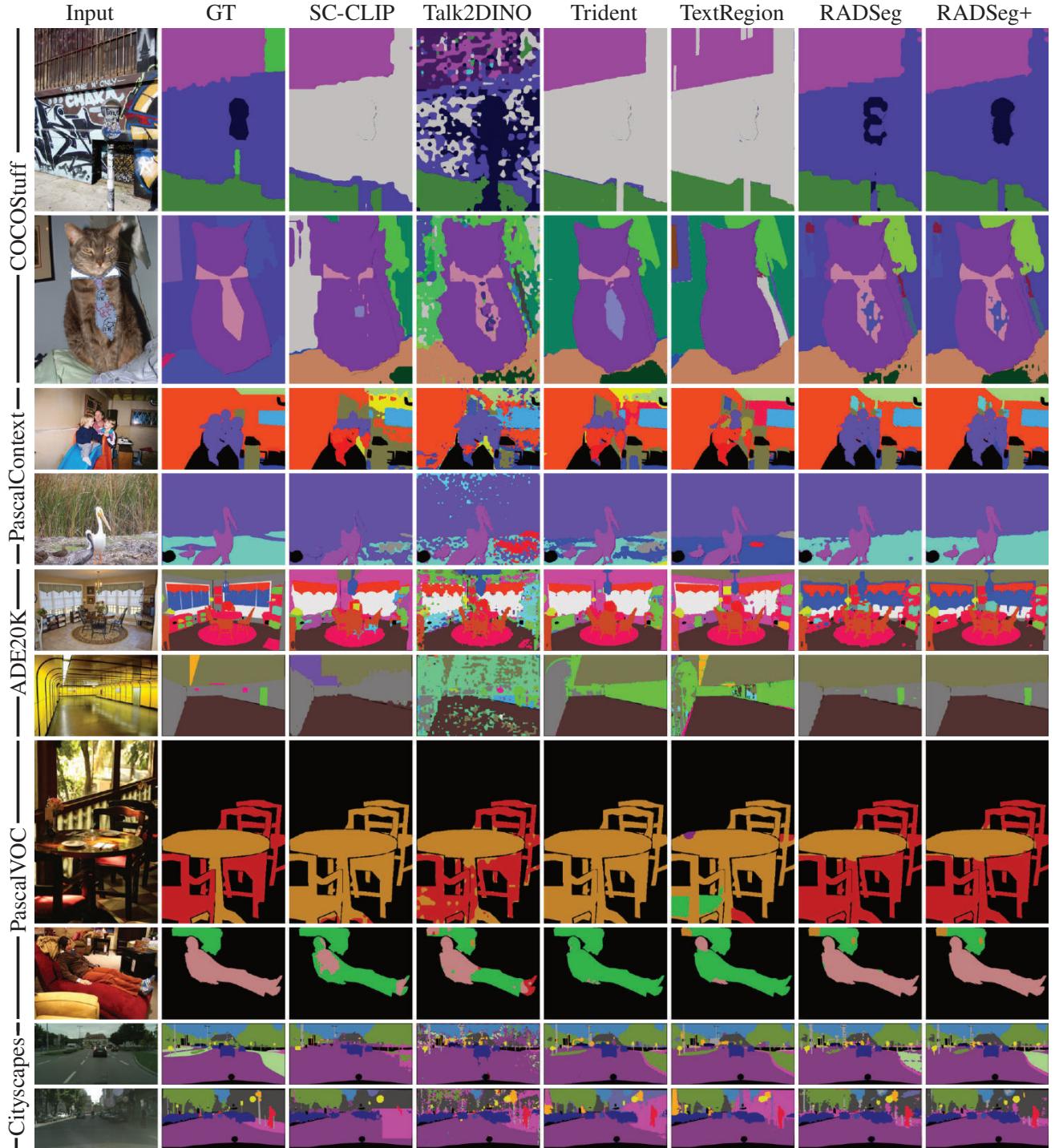| Methods | Lat. (s) | | | | | Lat. (s) | Params (M) | mIoU (%) |
|---|---|---|---|---|---|---|---|---|
| | CTX | VOC | Stuff | ADE | City | | Avg | |
| **Base Models** | | | | | | | | |
| NACLIP [13] | 0.091 | 0.025 | 0.074 | 0.111 | 0.144 | 0.089 | 86 | 38.48 |
| ResCLIP [44] | 0.252 | 0.152 | 0.266 | 0.325 | 0.311 | 0.261 | 86 | 40.33 |
| RayFronts [1] | 0.106 | 0.035 | 0.087 | 0.109 | 0.218 | 0.111 | 103 | 42.14 |
| ProxyCLIP [20] | 0.689 | 0.262 | 0.247 | 0.693 | 1.446 | 0.667 | 171 | 40.87 |
| SC-CLIP [3] | 0.148 | 0.059 | 0.111 | 0.174 | 0.243 | 0.147 | 86 | 41.86 |
| Trident [36] | 0.353 | 0.229 | 0.348 | 0.414 | 0.923 | 0.454 | 264 | 43.79 |
| TextRegion [42] | 0.410 | 0.403 | 0.438 | 0.428 | 1.729 | 0.682 | 167 | 43.33 |
| Talk2Dino [4] | 0.150 | 0.043 | 0.155 | 0.188 | 0.289 | 0.165 | 86 | 43.00 |
| **RADSeg** | **0.109** | **0.036** | **0.088** | **0.112** | **0.228** | **0.115** | **105** | **48.00** |
| **RADSeg+** | **0.261** | **0.156** | **0.286** | **0.360** | **0.709** | **0.354** | **125** | **50.01** |
| **Huge Models** | | | | | | | | |
| RayFronts [1] | 0.619 | 0.204 | 0.460 | 0.598 | 1.362 | 0.649 | 661 | 36.40 |
| ProxyCLIP [20] | 1.297 | 0.494 | 0.451 | 1.28 | 3.168 | 1.338 | 717 | 43.21 |
| Trident [36] | 1.426 | 0.900 | 1.123 | 1.444 | 2.365 | 1.451 | 1349 | 46.90 |
| TextRegion [42] | 0.732 | 0.670 | 0.711 | 0.744 | 2.476 | 1.067 | 857 | 47.50 |
| **RADSeg** | **0.617** | **0.202** | **0.457** | **0.594** | **1.363** | **0.647** | **664** | **46.92** |
| **RADSeg+** | **1.172** | **0.734** | **1.074** | **1.341** | **2.128** | **1.290** | **684** | **49.05** |

Figure A.3. Additional visualizations of 2D semantic segmentation results generated by **RADSeg**, **RADSeg+** and the most competitive baselines on images from all the benchmarks. Along with achieving SOTA mIoU for 2D OVSS, **RADSeg** and **RADSeg+** produce more precise segmentation maps and sharper object boundaries for both single-object as well as multi-object complex scenes.
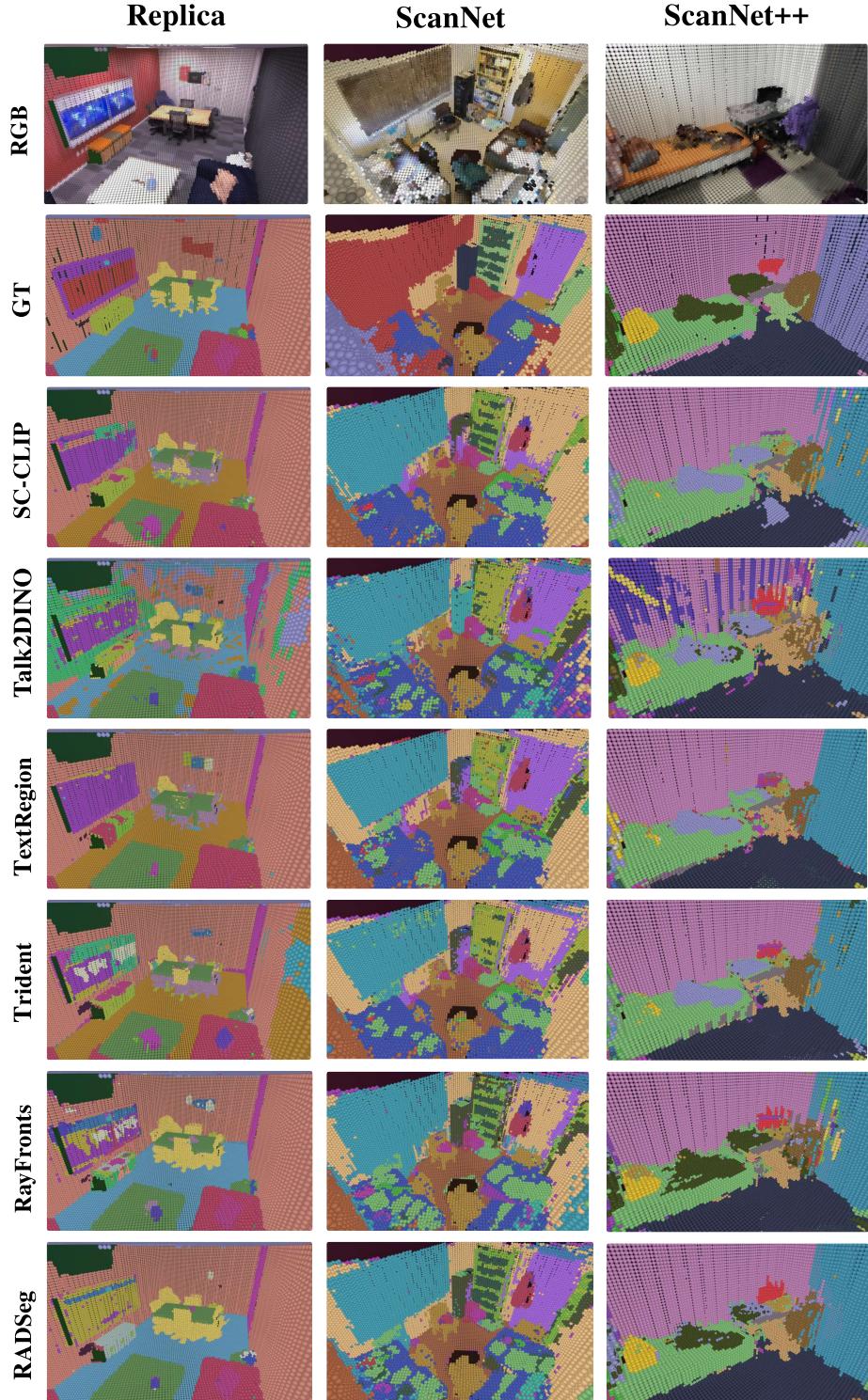
Figure A.4. Sample visualizations of 3D semantic segmentation results generated by **RADSeg** and all the baselines for scenes from Replica (scene-`office2`), ScanNet (scene-`0378`), and ScanNet++ (scene-`ea42cd27e6`. "RGB" and "GT" refer to the RGB scene reconstruction and Ground Truth semantics for each corresponding scene. **RADSeg** achieves SOTA mIoU for 3D OVSS and produces cleaner and more accurate semantic voxels.