# ReME: A Data-Centric Framework
# for Training-Free Open-Vocabulary Segmentation

Xiwei Xuan, Ziquan Deng, and Kwan-Liu Ma

University of California, Davis

{xwxuan, ziqdeng, klma}@ucdavis.edu

## Abstract

*Training-free open-vocabulary semantic segmentation (OVS) aims to segment images given a set of arbitrary textual categories without costly model fine-tuning. Existing solutions often explore attention mechanisms of pre-trained models, such as CLIP, or generate synthetic data and design complex retrieval processes to perform OVS. However, their performance is limited by the capability of reliant models or the suboptimal quality of reference sets. In this work, we investigate the largely overlooked data quality problem for this challenging dense scene understanding task, and identify that a high-quality reference set can significantly benefit training-free OVS. With this observation, we introduce a data-quality-oriented framework, comprising a data pipeline to construct a reference set with well-paired segment-text embeddings and a simple similarity-based retrieval to unveil the essential effect of data. Remarkably, extensive evaluations on ten benchmark datasets demonstrate that our method outperforms all existing training-free OVS approaches, highlighting the importance of data-centric design for advancing OVS without training. Our code is available here.*

## 1. Introduction

Open-vocabulary semantic segmentation (OVS) aims to segment images according to a set of arbitrary textual categories. While conventional approaches rely on training or fine-tuning large models with segment-text [13, 16, 22, 36, 39, 44, 49, 75] or image-text [11, 12, 63, 68, 71–73, 82] datasets, these methods incur significant annotation and computational costs. Recent advancements in vision-language models (VLMs), such as CLIP [51], have spurred interest in training-free OVS methods that leverage these pre-trained models without additional fine-tuning.

Some approaches assume VLMs have sufficient classification capabilities to recognize semantics, and attempt to enhance their pixel-level localization by modifying the at-
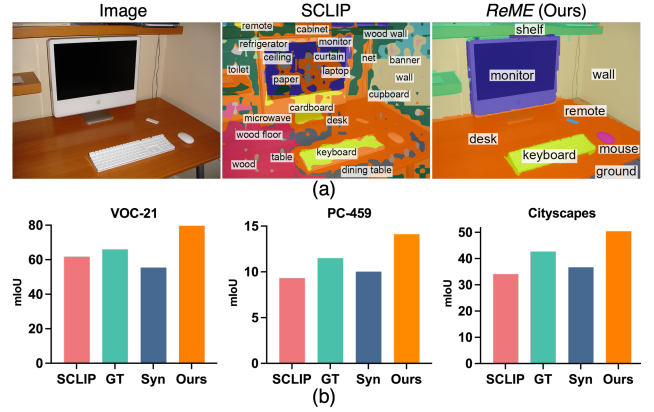


Figure 1. **(a)** Qualitatively compare *ReME* with SCLIP [62] (a representative CLIP-attention-based method) on an example from ADE20K [86] with 150 classes. Both methods have no mask post-processing. *ReME* produces higher-quality masks with a simple segmentation algorithm, yet SCLIP generates noisy masks paired with many irrelevant classes, such as *"curtain"* and *"microwave"*, etc. **(b)** Performance comparison between SCLIP and retrieval-based OVS using different reference sets. **GT**: ground-truth segment-label of COCO Stuff [8]; **Syn**: synthetic reference set from FreeDA [6]; **Ours**: our reference set from real images with quality refinement.

tention mechanisms [32, 33, 57, 60, 62, 64, 87]. However, VLMs are trained under weak image-text supervision [35], and merely tweaking their attention modules post-hoc cannot break through the inherent ceiling imposed by coarse training signals [24]. As a result, such approaches often struggle with low-quality segmentation masks and noisy labels. As shown by an example in Fig. 1 (a), methods relying on CLIP attention often fail to even match the segmentation quality of an algorithm-based segmenter [20], which, despite being class-agnostic, can propose perceptually meaningful image regions. (Refer to Fig. 5 for more examples).

Confronted with such limitations of VLMs, the retrieval paradigm offers a key inspiration: correcting the vulnerabilities of VLMs by reference substances in an external knowl-

edge base. We investigate this from the data-centric perspective – assuming there exists a high-quality reference set for OVS, can we unlock the potential of retrieval for better results? Specifically, we experiment with ground truth (GT) segment-text pairs from COCO Stuff [8], viewing it as a "high-quality" reference set. Using a class-agnostic segmentation algorithm [20] for test images, we apply a simple similarity-based strategy (refer to Sec. 3.3) for retrieving and aggregating mask labels. No models or complex retrieving algorithms are involved to unveil the fundamental data capabilities. As shown in Fig. 1 (b), retrieval from GT annotations consistently surpasses SCLIP [62], which is a well-performed method using CLIP attention. This observation indicates the significant potential of data for enhancing OVS with retrieval.

However, in the absence of labor-intensive GT annotations, despite the anticipated performance gain, existing retrieval-based methods remain modest compared to attention-based approaches. To fill the data gap, they often resort to diffusion models (DMs) [54] to generate synthetic images, utilizing text-corresponding attention masks to construct reference sets [5, 6, 27, 66]. However, synthetic data often lacks the realism and richness of real images, making them a suboptimal retrieval resource [18, 59, 79]. Fig. 1 (b) compares retrieval results from the COCO Stuff GT and a state-of-the-art synthetic reference set from FreeDA [6], constructed from COCO captions. The comparison indicates that the synthetic reference set falls significantly short of its GT counterpart, highlighting its limitations compared to real images. Some methods employ a quick fix of data bottleneck by designing sophisticated retrieval strategies to tolerate suboptimal data [66], which is less preferred as the central issue remains unaddressed [78].

In this work, we aim to construct **a well-aligned, rich, and contextually relevant set of segment-text pairs from real images**, investigating the capabilities of such data to benefit training-free OVS. Specifically, we introduce *ReME*, a data-centric framework that Refines Multi-modal Embeddings for retrieval-based, training-free OVS.

*ReME* harnesses the intrinsic capabilities of VLMs with the aggregation effect – while they may not yield satisfactory results for individual instances, VLMs demonstrate a certain degree of effectiveness when analyzing data at a collective level [15, 46, 50, 55, 56, 85, 88]. Using only images as input, the data pipeline of *ReME* starts by constructing a semantically rich base set with VLM-paired image segments and textual labels. With the observation of superior discriminativeness in same-modal data features, *ReME* leverages their collective patterns to recalibrate cross-modal misalignments introduced by the VLM, producing a reference set with enhanced quality. In the inference, we apply a basic class-agnostic segmentation algorithm to test images, then align and aggregate the proposed

masks into final predictions by referring to our reference set with a simple similarity-based strategy.

As shown in Fig. 1 (b), compared to using GT annotations for reference, our reference set indicates even better OVS capacity. Specifically, governed by our data quality consideration in semantic richness and alignment correctness, *ReME* brings significantly further performance gains for challenging benchmarks, such as PC-459 [43], and domain-specific scenarios, such as Cityscapes [14].

Extensive experiments on ten OVS benchmarks demonstrate the superior performance of *ReME* over existing training-free methods. Overall, our main contributions are summarized as follows:

- We introduce *ReME*, a data-centric framework for training-free OVS by constructing high-quality segment-text data from real images, with a simple retrieval process to perform OVS without any model fine-tuning.
- We identify the critical yet underexplored role of data quality, and provide a data pipeline for refining multi-modal data collectively from an intra-modal perspective, which unveils the data potential to regularize or complement pre-trained models.
- Extensive experiments on ten benchmark datasets demonstrate the consistent performance gains of *ReME* over 14 training-free OVS baselines.

## 2. Related Work

**Training-Free OVS.** Many training-free OVS methods analyze attention mechanisms [61, 76, 77] of VLMs or diffusion models to produce segmentation results [7, 32, 40, 53, 57, 64, 69, 87]. However, their reliant models trained without pixel-level supervision may produce inferior results, constraining the upper bound of these approaches. To improve the situation, methods such as CaR [60] leverage multiple CLIP models to iteratively refine segmentation results, which are less practical due to significant inference costs. Another line of methods [5, 6, 23, 27, 58, 66] constructs a reference set with extensive segment-text pairs, which serves as reference knowledge to support better alignment between test segments and text. Most of such methods either use synthetic data [5, 6, 27, 66] that falls short of diversity and realism [18, 59], or opt for sophisticated retrieval strategy or shortcut like adding GT classes to bypass the data quality issues [45, 66]. To fill the gap, our framework investigates the importance of data for this challenging task.

**Multi-Modal Datasets Curation.** High-quality datasets have driven remarkable progress in multi-modal models [11, 19, 41, 48, 65, 70, 80]. Recent techniques of image-text data curation typically involve two branches: *data filtering* [9, 21] that removes noisy data, and *data improvement* [31, 88] that refines multi-modal alignment. Despite training better VLMs is often the objective of data curation, these models themselves are also used for curating better
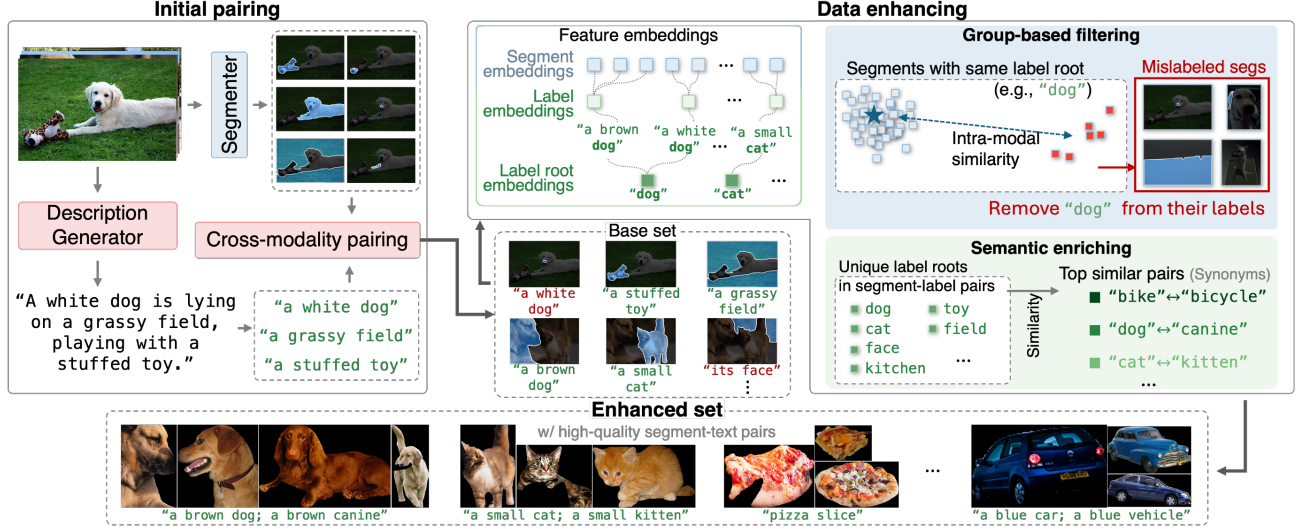
Figure 2. **The data pipeline of *ReME* to construct a high-quality reference set.** The pipeline includes two phases: the ***initial pairing*** that produces pre-matched segment-text pairs using images as input; and the ***data enhancing*** that performs group-based filtering according to more discriminative intra-modal similarity among segment embeddings and semantic enriching with similar labels.

data. For instance, CLIP is frequently used for data filtering with CLIP score (i.e., similarity between visual and text features) [50, 55, 56], and some methods use BLIP-2 to reduce semantic noise [46, 88]. As both *data filtering* and *data improvement* address essential aspects of data quality, we consider both of them in our approach for more comprehensive data enhancement. In addition, multi-modality data curation in fine-grain such as segment-text pairs remains underexplored, requiring further research to support more precise, context-aware applications. Beyond OVS, our proposed data cleaning methodology has the potential for broader applicability across other contexts.

**Pre-Training Models.** Pre-training models often have high reusability and adaptability, which can be easily reused in various approaches without further fine-tuning [27, 40, 58, 81, 83]. For instance, *Vision-Language Models (VLMs)* like CLIP [51] and ALIGN [26] link text with images for zero-shot image classification. *Multimodal Large Language Models (MLLMs)*, such as LLaVA [38], BLIP-2 [34], and GPT-Vision [1], support diverse tasks such as captioning, reasoning, and visual chat. Vision models like DINO [10] and DINOv2 [47] provide universal visual feature representations. Despite their success, such models are not almighty solutions, having various limitations for specified tasks [30]. In this work, we study how to effectively harness the value of data to compensate for the shortcomings of such models in a training-free manner.

## 3. Our Approach

The goal of an OVS method is to segment an input image given an arbitrary set of textual labels. Our training-free framework *ReME* consists of a data pipeline with two key phases, ***initial pairing*** and ***data enhancing***, to construct a high-quality reference set of segment-text pairs (Fig. 2); followed by a simple ***similarity based retrieval*** process to unveil the performance gains enabled by our data (Fig. 4).

### 3.1. Initial Pairing

Our data pipeline starts by constructing a diverse base set with segment-text pairs, where we aim to ensure broad coverage of various objects and contextual elements.

As shown in Fig. 2, for each input image, a segmenter generates class-agnostic segment masks, and an image description generator produces semantically rich descriptions. Following [48, 70], we extract noun phrases (i.e., nouns with descriptive modifiers such as *"a lovely white rabbit"*) from the descriptions to produce candidate textual labels, which reserve richer semantics than pure nouns. Then, we pair the segments and labels using a capable VLM. Despite the noises, CLIP still provides a certain degree of correctness for segment labeling, as evidenced by CLIP-based OVS methods [23, 32, 57, 60, 87]. We thus leverage CLIP to perform initial pairing. Given an image, we extract CLIP embeddings of its segments and noun phrases, $S \in \mathbb{R}^{m \times d}$ and $L \in \mathbb{R}^{n \times d}$, respectively, where $m$, $n$ is the number of segments and labels, and $d$ is the embedding dimension. Note that all embeddings across the paper are L2 normalized so that $X \cdot Y^T$ is equivalent to their cosine similarity $\langle X, Y \rangle$. We then compute $sim = S \cdot L^T \in \mathbb{R}^{m \times n}$ to pair each label with its closest segment, where $sim_{ij}$ is the similarity between $S_i$ and $L_j$. The top-matching $S_{i^*}$ of $L_j$ is identified by: $i^* = \arg \max_i sim_{ij}$. The unpaired segments are dropped to reduce the redundancy in class-agnostic image segments.

Despite including diverse segments paired with noun

Figure 3. **The superiority of intra-modality over cross-modality for data issue detection.** The plot provides the UMAP projection of segment embeddings in the base set labeled as "dog", colored by cross-modal similarity scores (CLIP scores). **Blue boxes** highlight misalignments detected by our filtering; **orange boxes** are those detected by low CLIP scores, which remove correct pairings while leaving many misalignments unaddressed.

phrase labels, the base set still has unaddressed quality concerns, which is our major focus in the following phase.

## 3.2. Data Enhancing

As shown in Fig. 2, through initial pairing, we obtain a base set with paired segments and labels. Various factors negatively affect its quality: (1) irrelevant textual labels caused by object hallucinations of the description generator, (2) meaningless or partial segments from over-segmentation, and (3) segment-text pairing errors. These challenges are essentially introduced by the ambiguity in cross-modality. We opt for navigating them with intra-modality operations—which we find to be more discriminative to effectively clean and enrich such data.

Unlike the common practice of dropping segment-label pairs with <u>low cross-modal similarity scores (i.e., CLIP scores)</u> [50, 55, 56], we reexamine data features and test whether CLIP scores can favorably clean noisy data. As shown in Fig. 3, our observation brings two key insights: (1) CLIP scores tend to mistakenly remove correct pairs while leaving many misalignments unaddressed; (2) Visual features, i.e., embeddings of image segments, have favorable discriminative power for data issue detection. (Please refer to Sec. 4.3 for extensive discussions.)

**Group-Based Filtering.** Guided by our observation, in this phase, we detect misalignments by leveraging the discriminativeness of intra-modal features. Specifically, we aim to leverage the visual features of segments to identify outliers in their own modality. First, for our labels in the format of noun phrases, we view those with the same "root" noun as an identical label. Next, we group each single segment-label pair according to identical labels. For example, segments labeled by *"a small dog"* or *"an adorable dog"* belong to the same group denoted by the root *"dog."* Note that we reserve all descriptive modifiers in the phrases to keep se-

mantic richness; "roots" are just used for grouping.

This is an important step for detecting data misalignments. Since each group of segments corresponds to an identical label, their visual features should be inherently consistent, and misaligned data would be automatically highlighted as outliers. As shown in Fig. 3, the majority of segments marked by ❶ accurately depict dogs, while others are misalignments located as outliers in the projection, indicating their distinctive features.

In detail, we define the segment center $S_{center}$ of each group, computed by the median of all visual features. $S_{center}$ denotes the representative features for correctly labeled segments. We then compute the intra-modal similarity between each segment in the current group and $S_{center}$, indicating the possibility of a segment being mistakenly labeled. $\delta_{filter}$ percent of segments with the lowest scores are considered mistakenly paired, where we remove the corresponding noun phrases from their labels. By performing this filtering guided by intra-modality feature similarity, we effectively reduce misalignment in the base set.

**Semantic Enriching.** By clearing outliers in each group, we have significantly mitigated segment-text misalignments. However, data enhancement remains incomplete in terms of semantic diversity – While statistical analysis of the entire label corpus of our base set reveals the textual modality includes different words for the same concept (e.g., *"bike"* and *"bicycle"*), such alternatives are missing in individual segments. This results in a lack of semantic richness, as we expect our method to recognize concepts in various representations. Consistently, we design an intra-modality-based approach to diversify our labels. As shown in Fig. 2, we collect embeddings of identical labels, i.e., the "root" noun of our noun phrase labels, denoted as $\{L_1, \ldots, L_n\}$. Then, we compute their pair-wise cosine similarity $\langle L_i, l_j \rangle$, and identify $k$ top-similar pairs to be treated as synonyms. Next, label semantics are enriched by adding these synonyms. For instance, *"cat"*-*"kitten"* is a pair of synonyms, so, if a segment has a label *"a small cat"*, we add *"a small kitten"* to become its label as well. Those synonym-enriched noun phrases are our finalized labels. This process further diversifies our textual descriptions by leveraging the global semantic richness present within the data. Note that, an LLM could produce even better semantic enriching, while we opt for a lightweight-yet-effective solution to avoid the additional overhead of large models.

## 3.3. Similarity-Based Retrieval

Following Sec. 3.2, we have constructed a reference set with well-paired segments and labels and stored their embeddings. At inference time, given a test image and a set of target classes, our goal is to assign pixel-level labels to the test image by referring to the reference set. Inspired by Tip-Adapter [84], we perform a simple retrieval strategy
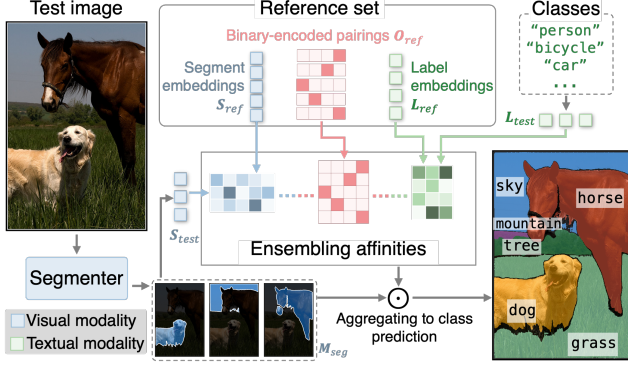
Figure 4. The process of similarity-based retrieval from the reference set given a test image and classes. Note that the same modality embeddings (i.e., reference & test segments; reference & test labels) are encoded by the same feature encoder, respectively.

grounded in feature similarity, where labels of test segments are estimated by their closest matches in the reference set.

As shown in Fig. 4, the reference set includes embeddings of $m$ segments $S_{\text{ref}} \in \mathbb{R}^{m \times d_1}$ and $n$ labels $L_{\text{ref}} \in \mathbb{R}^{n \times d_2}$. The label assignments are encoded by a binary matrix $O_{\text{ref}} \in \mathbb{R}^{m \times n}$, with an entry 0 indicating the corresponding label is absent in the segment and 1 otherwise.

Given a test image, we segment it into $k$ class-agnostic masks, denoted by $M_{\text{seg}} \in \mathbb{R}^{k \times h \times w}$, where $h \times w$ is the image size. We extract the segment and test class embeddings $S_{\text{test}} \in \mathbb{R}^{k \times d_1}$ and $L_{\text{test}} \in \mathbb{R}^{c \times d_2}$, using the same encoders applied to the reference set, where $d_1$ and $d_2$ are the embedding dimensions of the image and text modalities, respectively. The affinity between test segments and reference set labels is computed as:

$$A_1 = \text{Softmax}(S_{\text{test}} \cdot S_{\text{ref}}^T) \cdot O_{\text{ref}}, \qquad (1)$$

where $A_1 \in \mathbb{R}^{k \times n}$, and Softmax normalizes the affinities to formulate the probability distribution of labels for each test segment. Similarly, the affinity matrix between the reference labels and test classes is represented as:

$$A_2 = \text{Softmax}(L_{\text{ref}} \cdot L_{\text{test}}^T), \qquad (2)$$

where $A_2 \in \mathbb{R}^{n \times c}$. Ensembling the affinities, for $k$ test segments, the predicted logit is defined as: $P_{\text{seg}} = A_1 \cdot A_2, P_{\text{seg}} \in \mathbb{R}^{k \times c}$. Then, we calculate pixel-wise label probabilities $P_{\text{test}} \in \mathbb{R}^{h \times w \times c}$ by aggregating ensembled affinities according to the segment masks $M_{seg}$, where the probability of class $j$ at coordinates $(x, y)$ is:

$$P_{\text{test}}^{(x,y,j)} = \sum_{i=1}^{k} P_{\text{seg}}^{(i,j)} \cdot M_{\text{seg}}^{(i,x,y)}. \qquad (3)$$

Lastly, the predicted mask $\hat{l} \in \mathbb{R}^{h \times w}$ is determined by: $\hat{l}^{(x,y)} = \arg\max_j P_{\text{test}}^{(x,y,j)}$, where $\hat{l}^{(x,y)} \in [0, c-1]$ indicates class predictions, which conclude our retrieval phase.

# 4. Experiments

## 4.1. Experimental Setup

**Datasets for standard benchmarks.** We conducted experiments on 10 widely-used OVS benchmarks to evaluate *ReME*, including validation splits of Pascal VOC (VOC) [17], Pascal Context (PC) [43], COCO Object (Object) [37], COCO Stuff (Stuff) [8], Cityscapes (City) [14], and ADE20K [86]. Specifically, VOC has 20 object classes (VOC-20). PC involves PC-459 with 459 classes and PC-59 with 59 frequent classes. Object and Stuff provide COCO-2017 image annotations of 81 and 171 classes, respectively. City captures urban street scenes with 19 classes. ADE20K includes A-150 and A-849 with 150 and 847 classes. For VOC-20 and PC-59, we consider pixels not belonging to any class as *"background"*, represented by VOC-21 and PC-60, respectively. We use the standard mean Intersection-over-Union (mIoU) to measure OVS performance.

**Implementation.** We obtain textual descriptions of input images with LLaVA-1.5 [38]. For class-agnostic segmenter, we deploy Felzenszwalb's algorithm by default, which is a lightweight superpixel-based algorithm, following the same settings as [6]. The initial pairing is performed by CLIP [51] with ViT-B. In the data-enhancing and retrieval phases, we use CLIP-encoded text features, and employ DINOv2 with ViT-L as the default visual feature encoder. Hyperparameters $\delta_{filter}, k_{sim}$ are both set to 30. To compare with SAM-involved OVS baselines, we optionally use SAM as the segmenter, where we prompt a $32 \times 32$ point grid to obtain masks. More implementation details are in supplementary.

Our reference set construction requires only images as input—no GT captions, classes, or masks, and we default to using COCO-2017 [8, 37] images, which depict everyday scenes with objects in their natural contexts.

## 4.2. Comparison to the state-of-the-arts

**Baselines.** We compare *ReME* with 14 training-free OVS approaches, including ReCo [58], MaskCLIP [87], SCLIP [62], NACLIP [25], GEM [7], PnP [40], FreeDA [6], RIM [66], OVDiff [27], CLIPtrase [57], Diff-Segmenter [64], CaR [60], ProxyCLIP [32], and CorrCLIP [83]. Since post-processing techniques for mask refinements, such as DenseCRF [29], enhance OVS performance but introduce additional computational overhead, we indicate whether each method applies such refinements. Moreover, as SAM is widely regarded as a strong backbone for segment-related tasks, we conduct our comparison across both SAM-free and SAM-involved approaches.

**Comparison.** Table 1 presents the quantitative comparison results, where we observe that *ReME* consistently outperforms training-free OVS baselines across ten evaluated benchmarks. Notably, *ReME* achieves superior performance in challenging benchmarks with complex scenarios
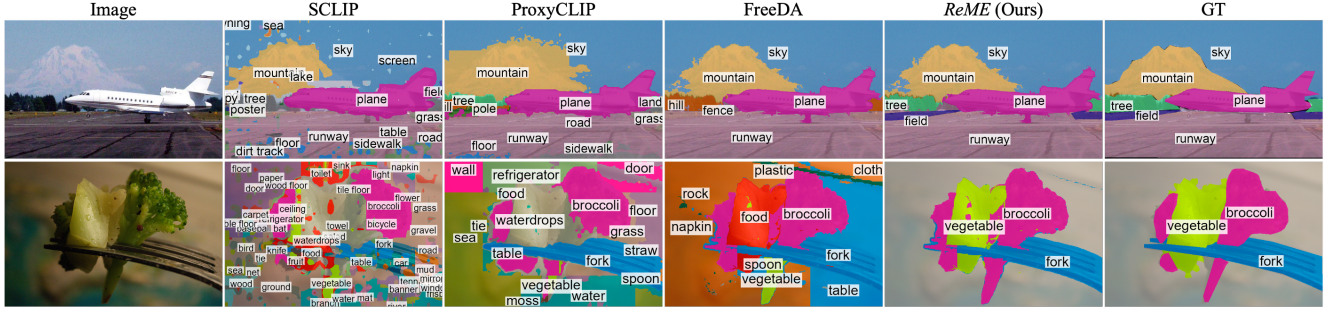
Figure 5. **Qualitative results of *ReME* in comparison with other training-free OVS methods.** Examples are from ADE20K [86] (w/ 150 classes) and COCO Stuff [8] (w/ 171 classes), respectively. SCLIP is based on CLIP attention; ProxyCLIP enhances CLIP attention with DINO features; FreeDA and *ReME* are retrieval-based methods, adopting the same superpixel-algorithm [20] for class-agnostic segmentation. We observe increasing quality of OVS results from left to right, with less noise in both masks and assigned labels.

| Methods | Post-processing | mIoU | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | VOC-20 | VOC-21 | City | PC-59 | PC-60 | Object | Stuff | A-150 | PC-459 | A-847 |
| *Training-free Methods without SAM* | | | | | | | | | | | |
| GEM [7] | ✗ | 46.2 | 24.7 | - | 32.6 | 21.2 | - | 15.1 | 10.1 | 4.6 | 3.7 |
| MaskCLIP [87] | ✓ | 74.9 | 38.8 | 12.6 | 25.5 | 23.6 | 20.6 | 14.6 | 9.8 | - | - |
| ReCo [58] | ✓ | 62.4 | 27.2 | 23.2 | 24.7 | 21.9 | 17.3 | 16.3 | 12.4 | - | - |
| SCLIP [62] | ✓ | 83.5 | 61.7 | 34.1 | 36.1 | 31.5 | 32.1 | 23.9 | 17.8 | 9.3 | 6.1 |
| CaR [60] | ✓ | <u>91.4</u> | 67.6 | 15.1 | 39.5 | 30.5 | 36.6 | 11.2 | 17.7 | 11.5 | 5.0 |
| NACLIP [25] | ✓ | 83.0 | 64.1 | 38.3 | 38.4 | 35.0 | 36.2 | 25.7 | 19.1 | 9.0 | 6.5 |
| CLIPtrase [57] | ✓ | 81.2 | 53.0 | 21.1 | 34.9 | 30.8 | <u>39.6</u> | 24.1 | 17.0 | 9.9 | 5.9 |
| PnP [40] | ✓ | 79.1 | 51.3 | 19.3 | 31.0 | 28.0 | 36.2 | 17.9 | 14.2 | 5.5 | 4.2 |
| FreeDA [6] | ✓ | 87.9 | 55.4 | 36.7 | <u>43.5</u> | <u>38.3</u> | 37.4 | <u>28.8</u> | 22.4 | 10.2 | 5.3 |
| ProxyCLIP [32] | ✗ | 83.2 | 60.6 | <u>40.1</u> | 37.7 | 34.5 | 39.2 | 25.6 | <u>22.6</u> | 11.2 | <u>6.7</u> |
| DiffSegmenter [64] | ✓ | 71.4 | 60.1 | - | 27.5 | 25.1 | 37.9 | - | - | - | - |
| OVDiff [27] | ✓ | 80.9 | <u>68.4</u> | 23.4 | 32.9 | 31.2 | 36.2 | 20.3 | 14.1 | <u>12.0</u> | 6.6 |
| ***ReME* (Ours)** | ✗ | **92.3** | **79.6** | **50.4** | **44.9** | **41.6** | **45.5** | **33.1** | **26.1** | **14.1** | **8.4** |
| *ReME* (Ours - VOC) | ✗ | 84.7 | **75.0** | 43.9 | 40.9 | **38.7** | 40.8 | 22.6 | **25.2** | 12.8 | 8.3 |
| *ReME* (Ours - ADE) | ✗ | 84.3 | **72.3** | 42.1 | 44.0 | **39.7** | 35.8 | 27.0 | **26.0** | 13.2 | 8.6 |
| *Training-free Methods with SAM* | | | | | | | | | | | |
| RIM [66] | ✗ | 77.8 | - | - | 34.3 | - | <u>44.5</u> | - | - | - | - |
| CaR w/ SAM [60] | ✗ | - | 70.2 | 16.9 | 40.5 | 31.1 | 37.6 | 12.4 | 17.9 | <u>11.8</u> | 5.7 |
| CLIPtrase w/ SAM [57] | ✗ | 82.3 | 57.1 | - | 36.4 | 32.0 | 44.2 | 24.8 | 17.2 | 10.6 | <u>6.0</u> |
| ProxyCLIP w/ SAM [32] | ✗ | 80.4 | 59.3 | 37.0 | 37.0 | 33.6 | 35.4 | 25.0 | <u>19.1</u> | 6.9 | 4.8 |
| CorrCLIP [83] | ✗ | <u>91.6</u> | <u>74.1</u> | <u>47.7</u> | <u>45.5</u> | <u>40.3</u> | 43.6 | <u>30.6</u> | - | - | - |
| ***ReME* w/ SAM (Ours)** | ✗ | **93.2** | **82.2** | **59.0** | **53.1** | **44.6** | **48.2** | **33.3** | **28.2** | **15.8** | **8.8** |

Table 1. **Comparison to state-of-the-art training-free OVS approaches.** The best results are **bolded**, with the second-best <u>underlined</u>. We also analyze data robustness by varying the image resources from the default COCO-2017 to VOC and ADE, respectively (w/o SAM version), where leading performances over SAM-free baselines are **bolded**.

and large numbers of classes, including A-150, PC-459, and A-847. On Cityscapes, an increase of 8.4 mIoU points (11.3 mIoU increase for SAM-involved comparison) further indicates that our approach excels in domain-specific scenarios. We also provide qualitative results of *ReME* in Fig. 5.

**Data robustness.** To evaluate the data robustness of *ReME*, we vary image resources of our reference set, replacing the default COCO-2017 images with those from VOC and ADE. As shown in Table 1, SAM-free *ReME* maintains

strong performance over baselines. Across 10 benchmarks, *ReME* achieves the highest mIoU in 6 benchmarks with VOC and 7 with ADE. These results highlight the effectiveness and flexibility of our approach in diverse settings.

### 4.3. Analysis, Ablation, and Discussion

**Hyperparameters.** Our data-enhancement process involves two key hyperparameters: *the drop ratio per group* ($\delta_{filter}$) in *group-based filtering* (i), and *the number of top-similar label pairs* ($k_{sim}$) in *semantic enriching* (ii). To de-
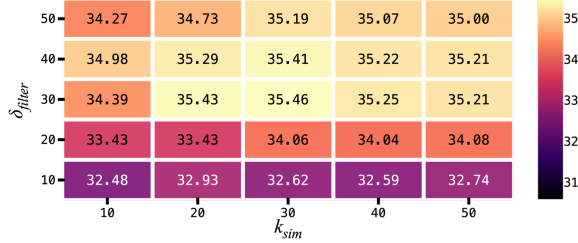
Figure 6. Hyperparameter analysis. $\delta_{filter}$ is the drop ratio in *group-based filtering*; $k_{sim}$ is the number of top-similar pairs in *semantic enriching*.

| Components | mIoU | | | |
|---|---|---|---|---|
| | VOC-20 | PC-59 | Object | A-150 |
| Base set (no enhancement) | 70.03 | 35.42 | 39.38 | 22.03 |
| w/ (i) Group-based filtering | 91.10 | 40.66 | 42.48 | 24.09 |
| w/ (ii) Semantic enriching | 79.50 | 36.69 | 39.92 | 23.41 |
| w/ (i) and (ii) | **92.34** | **44.89** | **45.50** | **26.13** |

Table 2. Impact of data enhancement components.

termine their values, we perform a grid search on $1k$ (1%) randomly sampled images from the training split of COCO Stuff (with no overlap with evaluation data), as shown in Fig. 6. Initially, increasing both parameters improves performance, as more misaligned data is filtered out and label diversity is enhanced. However, beyond a certain point, further increases cause a decline in performance: a high drop ratio removes correct pairings and reduces data diversity, while an excessive $k_{sim}$ introduces noise by including less-similar labels. We use 30 as the default value for both parameters in our implementation.

**Data enhancing component analysis.** We analyze the effects of our data enhancement components: (i) group-based filtering and (ii) semantic enriching. As shown in Table 2, the base set without enhancement yields relatively low performance, highlighting the critical role of data quality. We also observe higher individual effectiveness of (i), while enriching labels alone provides modest gains due to data noises. Notably, the performance improves significantly when both are integrated, demonstrating the complementary benefits of our data enhancement components.

**Analysis of different data-cleaning approaches.** We apply group-based filtering to remove $\delta_{filter}\%$ of outliers from each label-identified group (refer to Sec. 3.2). This experiment examines the impact of different data-cleaning strategies by comparing three methods: (a) the common approach of dropping segment-label pairs with the globally lowest cross-modal similarity scores (CLIP scores); (b) a variant that also uses CLIP scores, removing segment-label pairs within each constructed groups; (c) our group-based filtering based on intra-modal feature similarity. By only changing this module, we evaluate OVS performance as

| Data cleaning | mIoU | | | |
|---|---|---|---|---|
| | VOC-20 | PC-59 | Object | A-150 |
| Global CLIP score[a] | 79.34 | 39.79 | 40.79 | 21.18 |
| Group-based CLIP score[b] | 80.05 | 41.84 | 41.81 | 22.79 |
| Ours [c] | **92.34** | **44.89** | **45.50** | **26.13** |

Table 3. Analysis of different data filtering approaches.

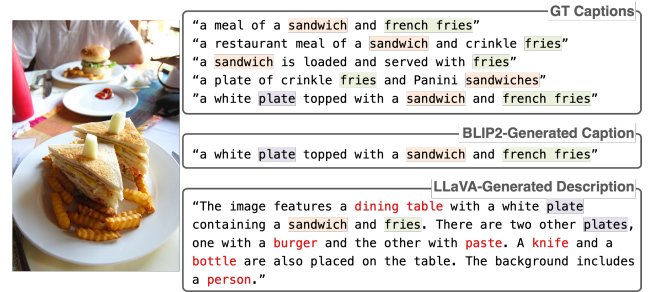| Feature encoder | mIoU | | | |
|---|---|---|---|---|
| | VOC-20 | PC-59 | Object | A-150 |
| $CLIP_B$ | 91.61 | 36.51 | 39.82 | 24.72 |
| $CLIP_L$ | 92.16 | 37.63 | 40.52 | 25.56 |
| $DINOv2_B$ | 91.72 | 43.65 | 44.71 | 25.29 |
| $DINOv2_L$ | **92.34** | **44.89** | **45.50** | **26.13** |

Table 4. Analysis of visual encoder variations.



Figure 7. Image descriptions from different resources. Red text highlights concepts uniquely present in the LLaVA description.

shown in Table 3, where we keep the drop ratio fixed at 30%. The results indicate that (a) and (b) perform much worse than (c), suggesting that filtering based on CLIP scores is less effective, likely due to its weaker alignment with actual pairing quality.

**Analysis of different visual encoders.** We investigate the impact of different visual feature encoders, including CLIP and DINOv2 with ViT-B and ViT-L architectures, denoted by $CLIP_B$, $CLIP_L$, $DINOv2_B$, and $DINOv2_L$, respectively. As shown in Table 4, upgrading from the base to large variants within the same encoder family results in slight performance improvements. When comparing CLIP and DINOv2 at the same scale, DINOv2 achieves better performance as expected, given its specialization in visual representation learning. Notably, by using CLIP alone and removing DINOv2 from our framework, we still achieve performance that is comparable to or better than baselines using similar backbones. This suggests that our approach is robust to variations in visual encoders.

**Analysis of the description generator.** Providing semantically rich image descriptions is important for our data construction. However, not all descriptions are equally informative. Fig. 7 presents an image from the COCO-2017 [8]

| Description generator | mIoU | | | |
|---|---|---|---|---|
| | VOC-20 | PC-59 | Object | A-150 |
| LLaVA [38] | **92.34** | **44.89** | **45.50** | **26.13** |
| BLIP-2 [34] | 89.41 | 40.85 | 37.64 | 24.31 |
| GT Caption | 89.02 | 40.15 | 37.77 | 24.09 |

Table 5. Ablation study of the image description generator.

| Segmenter | mIoU | | | |
|---|---|---|---|---|
| | VOC-20 | PC-59 | Object | A-150 |
| Superpixel [20] | 92.34 | 44.89 | 45.50 | 26.13 |
| SAM [28] | 93.15 | **53.10** | 48.21 | **28.21** |
| SAM2 [52] | **93.18** | 52.04 | **48.40** | **28.21** |

Table 6. Ablation study of the segmenter.

| Methods | Input data | | Reference set | | |
|---|---|---|---|---|---|
| | Format | Size | #labels | #pairs | Storage |
| Ours | Images | 118k | 41k | 1,023k | 4GB |
| FreeDA [6] | Captions | 591k | 19k | 2,167k | 17GB |

Table 7. Comparing data effectiveness with FreeDA, a retrieval-based OVS baseline that uses a synthetic reference set.

| Methods | PC-59 | A-150 | PC-459 | A-847 |
|---|---|---|---|---|
| ProxyCLIP [32] | 0.17 | 0.27 | 0.21 | 0.32 |
| SCLIP [62] | 0.20 | 0.34 | 0.57 | 0.53 |
| CaR [60] | 5.06 | 14.67 | 33.46 | 75.09 |
| FreeDA [6] | 0.89 | 0.89 | 4.89 | 5.01 |
| Ours | 0.29 | 0.34 | 0.31 | 0.34 |

Table 8. Inference time comparison (seconds/image).

dataset alongside its descriptions from three resources: (1) GT captions, (2) BLIP-2 [34], and (3) LLaVA [38], where we highlight the detected visual concepts. We can observe common concepts presented in all three, "sandwich", "fries", and "plate". However, concepts such as "dining table", "burger", "knife", "bottle", and "person" are missing from both (1) and (2). To further study the impact of using BLIP-2 and LLaVA as description generators, we conduct a quantitative evaluation. As shown in Table 5, LLaVA demonstrates superior performance, achieving significantly higher mIoU compared to BLIP-2, where BLIP-2 captions yield only marginal improvements over GT captions. These results highlight LLaVA's ability to generate richer, more detailed descriptions and quantitatively confirm its effectiveness in enhancing training-free OVS with retrieval.

**Analysis of the segmenter.** We also investigate the impact of different segmenters in Table 6 by replacing our default superpixel algorithm-based segmenter with advanced segmentation models, SAM [28] and SAM2 [52]. While a more advanced segmentation model leads to slightly better performance, our data-quality-enhancement pipeline enables even a simple segmentation algorithm to achieve strong results compared to baseline methods.

**Analysis of data effectiveness.** We evaluate data effectiveness by comparing *ReME* with FreeDA [6], a representative retrieval-based method that also utilizes COCO-2017 dataset [8, 37]. Unlike our approach, which takes images as input and generates textual descriptions, FreeDA uses GT captions to generate synthetic images for constructing its reference set. As shown in Table 7, despite FreeDA incorporating five times more input samples, its reference set contains fewer unique labels than *ReME*, due to the limited semantic diversity of GT captions (refer to Fig. 7). Our data-enhancement strategy improves data quality while reducing dataset size, resulting in a reference set with over a million fewer segment-text pairs than FreeDA. Using this considerably smaller reference set, *ReME* achieves supe-

rior performance even with a much simpler retrieval mechanism, highlighting the effectiveness of high-quality, semantically rich data over quantity.

**Analysis of inference time.** Inference time for training-free OVS varies due to factors such as the number of classes, image resolution, and inference strategies. To evaluate this, we compare our method with four representative approaches: (1) ProxyCLIP [32], which augments CLIP attention with DINO features, (2) SCLIP [62], which exploits CLIP self-attention followed by PAMR post-processing [67], (3) CaR [60], which iteratively queries two CLIP models for mask proposal and classification, and (4) FreeDA [6], a retrieval-based method. All experiments are conducted on two NVIDIA 4090 GPUs, measuring total inference time per dataset to compute the average. As shown in Table 8, (3) is the most time-consuming due to its iterative processing. Our method surpasses the retrieval-based baseline (4), benefiting from a smaller reference set and streamlined retrieval design, while remaining competitive with (1) and (2), which do not rely on reference sets.

## 5. Conclusion

In this work, we introduce *ReME*, a data-centric framework for training-free OVS by refining multi-modal embeddings. We observe the overlooked challenge of data quality, and demonstrate its critical impact on this dense scene understanding task. Following our data pipeline, we produce a reference set with well-aligned, rich, and contextually relevant segment-text pairs. Extensive experimental results highlight that enhancing data quality can be more beneficial than relying on complex retrieval mechanisms or model-specific adaptations. We hope this work can inspire future research for exploring data-centric strategies to further improve training-free OVS.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 3

[2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022. 3

[3] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023. 3

[4] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 3

[5] Luca Barsellotti, Roberto Amoroso, Lorenzo Baraldi, and Rita Cucchiara. Fossil: Free open-vocabulary semantic segmentation through synthetic references retrieval. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1464–1473, 2024. 2

[6] Luca Barsellotti, Roberto Amoroso, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. Training-free open-vocabulary segmentation with offline diffusion-augmented prototype generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3689–3698, 2024. 1, 2, 5, 6, 8, 3, 4, 7

[7] Walid Bousselham, Felix Petersen, Vittorio Ferrari, and Hilde Kuehne. Grounding everything: Emerging localization properties in vision-language transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3828–3837, 2024. 2, 5, 6, 7, 8

[8] 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. 1, 2, 5, 6, 7, 8, 4, 12

[9] Liangliang Cao, Bowen Zhang, Chen Chen, Yinfei Yang, Xianzhi Du, Wencong Zhang, Zhiyun Lu, and Yantao Zheng. Less is more: Removing text-regions improves clip training efficiency and robustness. *arXiv preprint arXiv:2305.05095*, 2023. 2

[10] 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

[11] Junbum Cha, Jonghwan Mun, and Byungseok Roh. Learning to generate text-grounded mask for open-world semantic segmentation from only image-text pairs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11165–11174, 2023. 1, 2, 4, 8

[12] Jun Chen, Deyao Zhu, Guocheng Qian, Bernard Ghanem, Zhicheng Yan, Chenchen Zhu, Fanyi Xiao, Sean Chang Culatana, and Mohamed Elhoseiny. Exploring open-vocabulary semantic segmentation from clip vision encoder distillation only. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 699–710, 2023. 1, 4, 8

[13] Seokju Cho, Heeseong Shin, Sunghwan Hong, Anurag Arnab, Paul Hongsuck Seo, and Seungryong Kim. Catseg: Cost aggregation for open-vocabulary semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4113–4123, 2024. 1, 3, 4, 8

[14] 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

[15] Ziquan Deng, Xiwei Xuan, Kwan-Liu Ma, and Zhaodan Kong. A reliable framework for human-in-the-loop anomaly detection in time series. *arXiv preprint arXiv:2405.03234*, 2024. 2

[16] Jian Ding, Nan Xue, Gui-Song Xia, and Dengxin Dai. Decoupling zero-shot semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11583–11592, 2022. 1, 4, 8

[17] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88:303–338, 2010. 5

[18] Lijie Fan, Kaifeng Chen, Dilip Krishnan, Dina Katabi, Phillip Isola, and Yonglong Tian. Scaling laws of synthetic images for model training... for now. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7382–7392, 2024. 2

[19] Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander T Toshev, and Vaishaal Shankar. Data filtering networks. In *The Twelfth International Conference on Learning Representations*, 2024. 2

[20] Pedro F Felzenszwalb and Daniel P Huttenlocher. Efficient graph-based image segmentation. *International journal of computer vision*, 59:167–181, 2004. 1, 2, 6, 8, 9

[21] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. *Advances in Neural Information Processing Systems*, 36, 2024. 2

[22] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Scaling open-vocabulary image segmentation with image-level labels. In *European Conference on Computer Vision*, pages 540–557. Springer, 2022. 1, 4, 8

[23] Zhongrui Gui, Shuyang Sun, Runjia Li, Jianhao Yuan, Zhaochong An, Karsten Roth, Ameya Prabhu, and Philip Torr.

kNN-CLIP: Retrieval enables training-free segmentation on continually expanding large vocabularies. *Transactions on Machine Learning Research*, 2024. 2, 3

[24] Yong Guo, David Stutz, and Bernt Schiele. Robustifying token attention for vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17557–17568, 2023. 1

[25] Sina Hajimiri, Ismail Ben Ayed, and Jose Dolz. Pay attention to your neighbours: Training-free open-vocabulary semantic segmentation. *arXiv preprint arXiv:2404.08181*, 2024. 5, 6, 7, 8

[26] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pages 4904–4916. PMLR, 2021. 3

[27] Laurynas Karazija, Iro Laina, Andrea Vedaldi, and Christian Rupprecht. Diffusion models for open-vocabulary segmentation. In *European Conference on Computer Vision*, pages 299–317. Springer, 2025. 2, 3, 5, 6, 7, 8

[28] 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. 8, 6

[29] Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. *Advances in neural information processing systems*, 24, 2011. 5

[30] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. LISA: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589, 2024. 3

[31] Zhengfeng Lai, Haotian Zhang, Bowen Zhang, Wentao Wu, Haoping Bai, Aleksei Timofeev, Xianzhi Du, Zhe Gan, Jiulong Shan, Chen-Nee Chuah, et al. Veclip: Improving clip training via visual-enriched captions. In *European Conference on Computer Vision*, pages 111–127. Springer, 2025. 2

[32] Mengcheng Lan, Chaofeng Chen, Yiping Ke, Xinjiang Wang, Litong Feng, and Wayne Zhang. Proxyclip: Proxy attention improves clip for open-vocabulary segmentation. In *ECCV*, 2024. 1, 2, 3, 5, 6, 8, 7

[33] Jingyao Li, Pengguang Chen, Shengju Qian, and Jiaya Jia. Tagclip: Improving discrimination ability of open-vocabulary semantic segmentation. *arXiv preprint arXiv:2304.07547*, 2023. 1

[34] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 3, 8, 6

[35] Yi Li, Hualiang Wang, Yiqun Duan, and Xiaomeng Li. Clip surgery for better explainability with enhancement in open-vocabulary tasks. *arXiv preprint arXiv:2304.05653*, 2023. 1

[36] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7061–7070, 2023. 1, 4, 8

[37] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 5, 8

[38] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 3, 5, 8, 6

[39] Huaishao Luo, Junwei Bao, Youzheng Wu, Xiaodong He, and Tianrui Li. Segclip: Patch aggregation with learnable centers for open-vocabulary semantic segmentation. In *International Conference on Machine Learning*, pages 23033–23044. PMLR, 2023. 1, 4, 8

[40] Jiayun Luo, Siddhesh Khandelwal, Leonid Sigal, and Boyang Li. Emergent open-vocabulary semantic segmentation from off-the-shelf vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4029–4040, 2024. 2, 3, 5, 6, 7, 8

[41] Anas Mahmoud, Mostafa Elhoushi, Amro Abbas, Yu Yang, Newsha Ardalani, Hugh Leather, and Ari S Morcos. Sieve: Multimodal dataset pruning using image captioning models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22423–22432, 2024. 2

[42] Divyanshu Malik, Xiwei Xuan, and Kwan-Liu Ma. Towards interactive 3d surgical scene reconstruction: An incremental training and monitoring framework. In *2025 IEEE 22nd International Symposium on Biomedical Imaging (ISBI)*, pages 1–4. IEEE, 2025. 5

[43] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 891–898, 2014. 2, 5, 4, 11

[44] Jishnu Mukhoti, Tsung-Yu Lin, Omid Poursaeed, Rui Wang, Ashish Shah, Philip HS Torr, and Ser-Nam Lim. Open vocabulary semantic segmentation with patch aligned contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19413–19423, 2023. 1

[45] Quang Nguyen, Truong Vu, Anh Tran, and Khoi Nguyen. Dataset diffusion: Diffusion-based synthetic data generation for pixel-level semantic segmentation. *Advances in Neural Information Processing Systems*, 36, 2024. 2

[46] Thao Nguyen, Samir Yitzhak Gadre, Gabriel Ilharco, Sewoong Oh, and Ludwig Schmidt. Improving multimodal datasets with image captioning. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 3

[47] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez,

Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2024. 3

[48] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023. 2, 3

[49] Yu Q. and et al. Convolutions die hard: Ovs with single frozen convolutional clip. In *NeurIPS*, 2023. 1

[50] Filip Radenovic, Abhimanyu Dubey, Abhishek Kadian, Todor Mihaylov, Simon Vandenhende, Yash Patel, Yi Wen, Vignesh Ramanathan, and Dhruv Mahajan. Filtering, distillation, and hard negatives for vision-language pre-training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6967–6977, 2023. 2, 3, 4

[51] 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, 3, 5

[52] 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 preprint arXiv:2408.00714*, 2024. 8, 6

[53] Pitchaporn Rewatbowornwong, Nattanat Chatthee, Ekapol Chuangsuwanich, and Supasorn Suwajanakorn. Zero-guidance segmentation using zero segment labels. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1162–1172, 2023. 2

[54] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2

[55] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 2, 3, 4

[56] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5b: An open large-scale dataset for training next generation image-text models. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. 2, 3, 4

[57] Tong Shao, Zhuotao Tian, Hang Zhao, and Jingyong Su. Explore the potential of clip for training-free open vocabulary semantic segmentation. *arXiv preprint arXiv:2407.08268*, 2024. 1, 2, 3, 5, 6, 7, 8

[58] Gyungin Shin, Weidi Xie, and Samuel Albanie. Reco: Retrieve and co-segment for zero-shot transfer. *Advances in Neural Information Processing Systems*, 35:33754–33767, 2022. 2, 3, 5, 6, 7, 8

[59] Krishnakant Singh, Thanush Navaratnam, Jannik Holmer, Simone Schaub-Meyer, and Stefan Roth. Is synthetic data all we need? benchmarking the robustness of models trained with synthetic images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2505–2515, 2024. 2

[60] Shuyang Sun, Runjia Li, Philip Torr, Xiuye Gu, and Siyang Li. Clip as rnn: Segment countless visual concepts without training endeavor. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13171–13182, 2024. 1, 2, 3, 5, 6, 8, 7

[61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 2

[62] 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. 1, 2, 5, 6, 8, 7

[63] Haoxiang Wang, Pavan Kumar Anasosalu Vasu, Fartash Faghri, Raviteja Vemulapalli, Mehrdad Farajtabar, Sachin Mehta, Mohammad Rastegari, Oncel Tuzel, and Hadi Pouransari. Sam-clip: Merging vision foundation models towards semantic and spatial understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3635–3647, 2024. 1, 4, 8

[64] Jinglong Wang, Xiawei Li, Jing Zhang, Qingyuan Xu, Qin Zhou, Qian Yu, Lu Sheng, and Dong Xu. Diffusion model is secretly a training-free open vocabulary semantic segmenter. *arXiv preprint arXiv:2309.02773*, 2023. 1, 2, 5, 6, 7, 8

[65] Xiaoqi Wang, Wenbin He, Xiwei Xuan, Clint Sebastian, Jorge Piazentin Ono, Xin Li, Sima Behpour, Thang Doan, Liang Gou, Han-Wei Shen, et al. Use: Universal segment embeddings for open-vocabulary image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4187–4196, 2024. 2

[66] Yuan Wang, Rui Sun, Naisong Luo, Yuwen Pan, and Tianzhu Zhang. Image-to-image matching via foundation models: A new perspective for open-vocabulary semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3952–3963, 2024. 2, 5, 6, 8

[67] Anne S Wannenwetsch and Stefan Roth. Probabilistic pixel-adaptive refinement networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11642–11651, 2020. 8

[68] Ji-Jia Wu, Andy Chia-Hao Chang, Chieh-Yu Chuang, Chun-Pei Chen, Yu-Lun Liu, Min-Hung Chen, Hou-Ning Hu, Yung-Yu Chuang, and Yen-Yu Lin. Image-text co-decomposition for text-supervised semantic segmentation. In

*Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26794–26803, 2024. 1, 4, 8

[69] 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

[70] Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan. Florence-2: Advancing a unified representation for a variety of vision tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4818–4829, 2024. 2, 3

[71] Yun Xing, Jian Kang, Aoran Xiao, Jiahao Nie, Ling Shao, and Shijian Lu. Rewrite caption semantics: Bridging semantic gaps for language-supervised semantic segmentation. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 4, 8

[72] 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*, pages 18134–18144, 2022. 4, 8

[73] Jilan Xu, Junlin Hou, Yuejie Zhang, Rui Feng, Yi Wang, Yu Qiao, and Weidi Xie. Learning open-vocabulary semantic segmentation models from natural language supervision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2935–2944, 2023. 1, 4, 8

[74] Mengde Xu, Zheng Zhang, Fangyun Wei, Yutong Lin, Yue Cao, Han Hu, and Xiang Bai. A simple baseline for open-vocabulary semantic segmentation with pre-trained vision-language model. In *European Conference on Computer Vision*, pages 736–753. Springer, 2022. 4, 8

[75] 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. 1, 4, 8

[76] Xiwei Xuan, Xiaoyu Zhang, Oh-Hyun Kwon, and Kwan-Liu Ma. VAC-CNN: A visual analytics system for comparative studies of deep convolutional neural networks. *IEEE Transactions on Visualization and Computer Graphics*, 28 (6):2326–2337, 2022. 2

[77] Xiwei Xuan, Ziquan Deng, Hsuan-Tien Lin, Zhaodan Kong, and Kwan-Liu Ma. Suny: A visual interpretation framework for convolutional neural networks from a necessary and sufficient perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8371–8376, 2024. 2

[78] Xiwei Xuan, Ziquan Deng, Hsuan-Tien Lin, and Kwan-Liu Ma. SLIM: Spuriousness mitigation with minimal human annotations. In *European Conference on Computer Vision*, pages 215–231. Springer, 2024. 2

[79] Xiwei Xuan, Jorge Piazentin Ono, Liang Gou, Kwan-Liu Ma, and Liu Ren. AttributionScanner: A visual analytics

system for model validation with metadata-free slice finding. *IEEE Transactions on Visualization and Computer Graphics*, pages 1–12, 2025. 2

[80] Xiwei Xuan, Xiaoqi Wang, Wenbin He, Jorge Piazentin Ono, Liang Gou, Kwan-Liu Ma, and Liu Ren. VISTA: A visual analytics framework to enhance foundation model-generated data labels. *IEEE Transactions on Visualization and Computer Graphics*, 2025. 2

[81] Xinyuan Yan, Xiwei Xuan, Jorge Piazentin Ono, Jiajing Guo, Vikram Mohanty, Shekar Arvind Kumar, Liang Gou, Bei Wang, and Liu Ren. Vislix: An xai framework for validating vision models with slice discovery and analysis. In *Computer Graphics Forum*, page e70125. Wiley Online Library, 2025. 3

[82] Muyang Yi, Quan Cui, Hao Wu, Cheng Yang, Osamu Yoshie, and Hongtao Lu. A simple framework for text-supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7071–7080, 2023. 1, 4, 8

[83] Dengke Zhang, Fagui Liu, and Quan Tang. Corrclip: Reconstructing correlations in clip with off-the-shelf foundation models for open-vocabulary semantic segmentation. *arXiv preprint arXiv:2411.10086*, 2024. 3, 5, 6, 8

[84] Renrui Zhang, Wei Zhang, Rongyao Fang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. Tip-adapter: Training-free adaption of clip for few-shot classification. In *European conference on computer vision*, pages 493–510. Springer, 2022. 4

[85] Xiaoyu Zhang, Xiwei Xuan, Alden Dima, Thurston Sexton, and Kwan-Liu Ma. Labelvizier: Interactive validation and relabeling for technical text annotations. In *2023 IEEE 16th Pacific Visualization Symposium (PacificVis)*, pages 167–176. IEEE, 2023. 2

[86] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 127:302–321, 2019. 1, 5, 6, 4, 10

[87] 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. 1, 2, 3, 5, 6, 7, 8

[88] Deyao Zhu, Jun Chen, Kilichbek Haydarov, Xiaoqian Shen, Wenxuan Zhang, and Mohamed Elhoseiny. Chatgpt asks, blip-2 answers: Automatic questioning towards enriched visual descriptions. *arXiv preprint arXiv:2303.06594*, 2023. 2, 3

# ReME: A Data-Centric Framework
# for Training-Free Open-Vocabulary Segmentation

## Supplementary Material

Our supplementary material is organized as follows:

## A. Problem Definition

To complement the definition of open-vocabulary segmentation (OVS), we formulate it mathematically as follows.

Given an input image $I$ and a candidate set of class labels $\mathcal{L} = \{\mathcal{L}_n\}_{n=1}^N$, the objective of OVS is to assign a class label $\mathcal{L}_n \in \mathcal{L}$ to each pixel in $I$. Each $\mathcal{L}_n$ represents the $n$-th class described by free-form text, where $N$ denotes the total number of candidate classes. Unlike traditional semantic segmentation, where the category set is fixed and predefined during training ($\mathcal{L} = \mathcal{L}_{\text{train}}$), OVS allows for segmentation of arbitrary and unseen categories, operating under a zero-shot setting. This flexibility facilitates adaptive and robust dense scene understanding in dynamic real-world scenarios.

## B. Approach and Implementation

### B.1. Prompt for Generating Image Descriptions

We design a specific prompt to obtain semantically enriched image descriptions using LLaVA. Our prompt is:

*"Describe this image in detail. Mention all visible objects, their parts, contexts, and characteristics like size, color, and texture. Also, describe the background/foreground context, including any natural scene or man-made structures, such as wall, ceiling, sky, and cloud. FOCUS ONLY on visible objects or contexts. Avoid speculation or guesses."*

### B.2. Filtering Ambiguous Labels

**Object Hallucination in MLLM Outputs.** Multi-modal large language models (MLLMs) such as LLaVA often suffer from object hallucination. This includes generating descriptions of tangible objects not present in the input image, which we address through our group-based filtering phase.
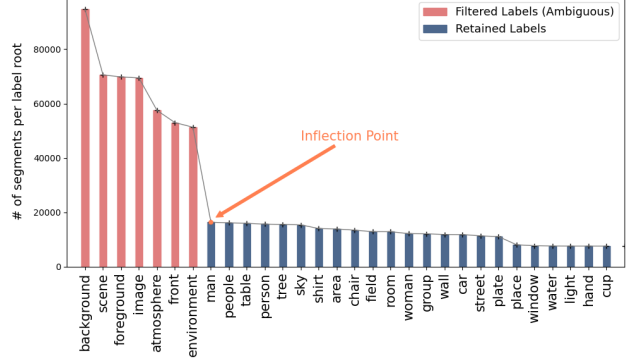


Figure A1. **The number of corresponding segments for each unique label root.** The knee of the distribution curve, *Inflection Point*, indicates the threshold for filtering out ambiguous labels.
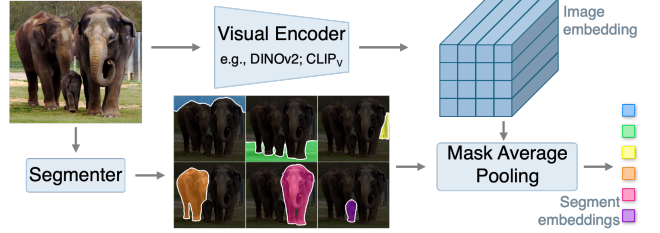


Figure A2. Visual feature encoding for image segments.

Additionally, MLLMs frequently produce ambiguous outputs reflecting abstract or subjective concepts evoked by the image. For instance, descriptions like "*The room has a cozy atmosphere*" lead to ambiguous labels such as "*atmosphere*," which are ungrounded in observable entities and irrelevant to segmentation tasks.

**Fast Filtering of Ambiguous Labels.** To address this issue, we propose a fast and effective approach to eliminate ambiguous labels arising from evoked descriptions. Due to their abstract nature, these labels appear frequently across MLLM-generated descriptions, often corresponding to an unusually large number of segments in the dataset. This observation forms the basis of our aggregation-based analysis. As described in Sec. 3.2, we group segment-text pairs by consistent label roots. For each group represented by a unique label root, we compute the total number of corresponding segments (i.e., the group size) and plot the distribution of group sizes. By identifying the knee of the curve—referred to as the inflection point (see Fig. A1)—we filter out labels exceeding this point, such as "background,"
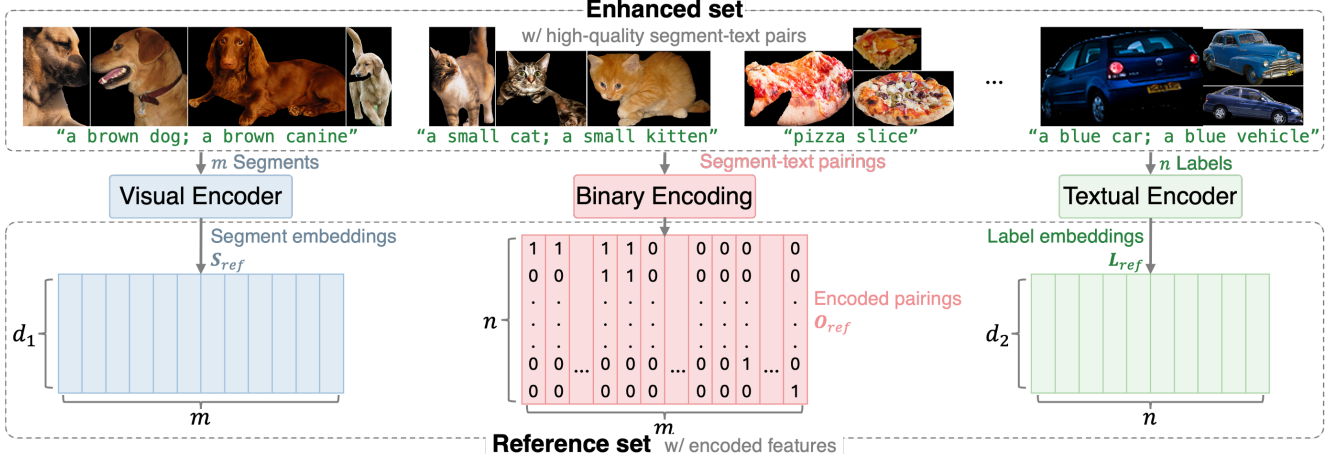
Figure A3. **Illustration of reference set encoding for similarity-based retrieval.** Image segments, textual labels, and their relationships are encoded as $\boldsymbol{S}_{\text{ref}}$, $\boldsymbol{L}_{\text{ref}}$, and $\boldsymbol{O}_{\text{ref}}$, respectively. These embeddings collectively form the reference set, enabling efficient retrieval.

"scene," "image," and "atmosphere." These labels dominate the dataset and detract from meaningful segmentation labels. Removing them ensures the dataset remains focused on concrete and observable objects, improving its relevance and usability for segmentation tasks.

### B.3. Feature Encoding

**Visual feature encoding.** We compute segment embeddings following common practices [6, 27, 65, 66]. According to the requirements, a visual encoder such as DINOv2 or $\text{CLIP}_V$ is used, denoted as $\varphi$. As shown in Fig. A2, given an input image $\boldsymbol{I}$ and its $K$ corresponding segment masks $\mathcal{M} = \{\boldsymbol{M}_k\}_{k=1}^K$, the visual encoder processes the image to obtain its embedding. To align the segment masks with the encoder's output resolution, the masks are resized using a downscaling function, $\zeta$. Lastly, we apply mask average pooling (MAP) to produce embedding for each segment $\boldsymbol{S}_k$. This process is represented as:

$$\boldsymbol{S}_k = \text{MAP}(\varphi(\boldsymbol{I}), \zeta(\boldsymbol{M}_k)). \tag{A1}$$

**Textual feature encoding.** We generate text embeddings using a textual encoder $\text{CLIP}_T$, denoted by $\phi$. For a given label $\mathcal{L}$, we deploy four templates to prompt the encoder: *"A photo of {}," "This is a photo of {}," "There is {} in the scene,"* and *"A photo of {} in the scene."* The text encoder processes each prompted input, and the resulting embeddings are averaged to form the final label embedding $\boldsymbol{L}$. This process is expressed as:

$$\boldsymbol{L} = \frac{1}{P} \sum_{p=1}^{P} \phi(\psi_p(\mathcal{L})), \tag{A2}$$

where $P$ is the number of templates, and $\psi_p(\mathcal{L})$ represents applying $p$-th template to label $\mathcal{L}$.

All encoded features, regardless of modality, are L2-normalized to facilitate our cosine similarity computation.

### B.4. Reference Set Construction

Following our intra-modality data enhancement phase (refer to Sec. 3.2), we have obtained a high-quality set of segment-text pairs. Fig. A3 depicts how we obtain specific embeddings to construct the reference set for streamlined retrieval. The visual encoder processes image segments to extract $d_1$-dimensional segment embeddings ($\boldsymbol{S}_{\text{ref}}$), while a textual encoder generates $d_2$-dimensional label embeddings ($\boldsymbol{L}_{\text{ref}}$). To represent the relationships between segments and their associated labels, we utilize binary encoding to formulate $\boldsymbol{O}_{\text{ref}} \in \mathbb{R}^{m \times n}$, where $m$ and $n$ are the numbers of unique segments and labels, respectively. Each row of $\boldsymbol{O}_{\text{ref}}$ corresponds to a segment, and a column entry of '1' indicates an association with a specific label and '0' otherwise. The resulting reference set is defined by $\{\boldsymbol{S}_{\text{ref}}, \boldsymbol{O}_{\text{ref}}, \boldsymbol{L}_{\text{ref}}\}$, combining visual, textual, and relational encodings. This structured representation enables efficient similarity-based retrieval in the subsequent phase.

### B.5. Pseudocode for Similarity-Based Retrieval

To complement Sec. 3.3, we provide a Python-style pseudocode in Alg. A1 to detail the similarity-based retrieval process. The variable names are consistent with those in Sec. 3.3 for ease of reference, and comments within the pseudocode indicate the steps corresponding to the equations discussed in the main paper.

## C. More Experimental Results and Discussion

### C.1. Additional Ablation Study Results

In this section, we provide comprehensive results and additional examples to supplement the findings presented in the

**Algorithm A1** Pseudocode for similarity-based retrieval

```
# Inputs:
#   S_ref [m,d_1] - Segment embeddings in the reference set
#   O_ref [m,n] - Binary encoding of segment-label relationships
#   L_ref [n,d_2] - Label embeddings in the reference set
#   I_test [h,w] - Test image
#   L_test - c test classes in text
# Outputs:
#   l_pred [h,w] - Predicted label mask for the test image

# Segment the test image into k class-agnostic masks
M_seg = segmenter(I_test) # [k,h,w]

# Same encoder as used for S_ref and L_ref
S_test = visual_encoder(I_test, M_seg) # [k,d_1]
L_test = textual_encoder(L_test) # [c,d_2]

# Intra-modality similarities
sim_seg = np.dot(S_test, S_ref.T) # [k,m]
sim_text = np.dot(L_ref, L_test.T) # [n,c]

# Compute and ensemble affinities # Eqns.(1-3)
A1 = np.dot(softmax(sim_seg, axis=1), O_ref) # [k,n]
A2 = softmax(sim_text, axis=1) # [n,c]
P_seg = np.dot(A1, A2) # [k,c]

# Aggregate segment-class probabilities # Eqn.(4)
P_test = np.einsum('ij,ihw->hwj', P_seg, M_seg) # [h,w,c]

# Compute the predicted label mask # Eqn.(5)
l_pred = np.argmax(P_test, axis=2) # [h,w]
```

main paper. The supplementary tables and figures expand on the quantitative and qualitative analyses in Sec. 4.3, offering a more complete view of our ablation studies. We cross-reference the corresponding tables/figures in the main paper for clarity and context.

- **Data enhancement component analysis.** Full quantitative results for analyzing contributions of individual components in our data enhancement pipeline are presented in Table A4 (supplementing Table 2 in the main paper).
- **Analysis of different data filtering approaches.** A comprehensive comparison of different data filtering approaches is provided in Table A5, extending the analysis from Table 3 in the main paper. We include a variant of our group-based filtering, noted as (d). Compared to our default approach that use the same drop ratio for all groups, (d) adapts each group's drop ratio to its segment consistency, ranging from 0 to 50%, with weights $w = \frac{1}{n} \sum_{i=1}^{n} (1 - \langle S_i, S_{center} \rangle)$, allowing for more drops in sparser groups. We can observe that this variant brings further performance gain.

  Additionally, we provide more examples to showcase the superiority of intra-modality over cross-modality in Fig. A4, to complement Fig. 3 in the main paper.
- **Feature encoder backbones.** Full results of using different feature encoder backbones are detailed in Table A6 (bottom), supplementing Table 4 in the main paper.
- **Analysis of the description generator.** Full results on the impact of the description generator are shown in Table A7, supplementing Table 5. In addition, to further evidence the semantic richness of LLaVA-generated descriptions as discussed in Fig. 7 in the main paper, we provide more examples in Fig. A5.
- **Analysis of the segmenter.** Additional results for the

| Method | VOC-20 | PC-59 | A-150 | PC-459 |
|---|---|---|---|---|
| LLaVA [38] as Classifier | 72.65 | 35.50 | 20.03 | 7.22 |
| Qwen [4] as Classifier | 70.67 | 36.03 | 21.21 | 6.36 |
| LLaVA [38] as Filter | 73.06 | 37.55 | 22.05 | 7.60 |
| Qwen [4] as Filter | 72.18 | 38.10 | 22.81 | 8.15 |
| *ReME* | 92.34 | 44.89 | 26.13 | 14.12 |
| *ReME* (OpenFlamingo [3]) | 92.54 | 44.77 | 25.95 | 13.39 |

Table A1. **Top**: Analysis of large MLLM capabilities. We use LLaVA-1.5 [38] and Qwen-2.5 VL [4] to (1) classify segmentation masks without any references, and (2) perform label filtering for data enhancing, respectively. They all perform significantly worse than *ReME*. This demonstrates that challenging tasks such as OVS require strategic adaptation rather than direct use. **Bottom**: Analysis of our performance gain from inherent segment-text pretraining. We replace LLaVA-1.5 [38] with OpenFlamingo [2] trained purely on image-text data. The performance remains comparable, indicating *ReME*'s effectiveness without dense annotations.

| Method († w/ seg-text training) | VOC-20 | PC-59 | A-150 | PC-459 |
|---|---|---|---|---|
| CAT-Seg (GT COCO) [13]† | 94.57 | 57.45 | 31.81 | 19.04 |
| CAT-Seg (*ReME*)† | 94.60 | 59.76 | 32.24 | 22.03 |
| FreeDA [6] | 87.91 | 43.49 | 22.43 | 10.24 |
| FreeDA (*ReME*) | 92.35 | 44.80 | 24.91 | 13.89 |
| *ReME* | 92.34 | 44.89 | 26.13 | 14.12 |

Table A2. **Data transferability.** We apply *ReME* data to two representative methods by replacing their training/reference data: (1) training-based CAT-Seg [13], and (2) retrieval-based FreeDA [6]. The results demonstrate the strong utility of our data across both training-based and training-free OVS.

impact of various segmenters are presented in Table A8, which complements Table 6 in the main paper.

- **Analysis of the large MLLM capabilities.** To analyze the capabilities of large MLLM compared to our data enhancement framework, we perform two experiments. (1) We directly leverage advanced MLLMs, including LLaVA-1.5 [38] and Qwen-2.5 VL [4], to assign class labels to class-agnostic segmentation masks, without using any data as references. (2) We perform data filtering with each MLLM, rather than using our group-based data filtering. The results are shown in Table A1, marked as "* as Classifier", and "* as Filter", respectively. They perform significantly worse than ReME. This observation aligns with widely discussed challenges in directly using VLMs for fine-grained data matching—they tend to hallucinate object labels and produce noisy predictions. These results highlight: while pre-trained models present potential, challenging tasks like reasoning segmentation [30] or OVS require strategic adaptation rather than direct use. For instance, LISA [30] fine-tunes vLLM+SAM backbones, while *ReME* studies data-centricity—they contribute in complementary ways.
- **Data transferability.** We apply *ReME* data to two representative methods by replacing their training/reference

data: (1) training-based CAT-Seg [13], and (2) retrieval-based FreeDA [6]. As shown in Table A2, CAT-Seg (*ReME*) even surpasses the version trained on COCO ground-truth, and FreeDA (*ReME*) also outperforms the original version with its default reference set. The results demonstrate the strong utility of *ReME* data across both training-based and training-free OVS settings.

## C.2. Additional Qualitative Results

We perform additional qualitative comparisons with other training-free baselines. The results are shown in Fig. A6. In addition, we present qualitative results of *ReME*-SAM on datasets with a large number of categories. Specifically, we include ADE20K [86] with 847 categories (Fig. A8), Pascal Context [43] with 459 categories (Fig. A9), and COCO Stuff [8] with 171 categories (Fig. A10).

## C.3. Backbone Usage for Training-Free Methods

Table A10 presents the backbone usage across various training-free methods. As shown, earlier approaches predominantly relied on a single CLIP backbone, but their overall performance falls short compared to more recent methods that leverage multiple backbones. Compared to these multi-backbone methods, our approach (1) remains entirely off-the-shelf, avoiding structural modifications to the backbone as implemented in ProxyCLIP, and (2) achieves the best performance while maintaining controlled backbone usage.

Additionally, existing methods employ different backbone variants, such as ViT-B/16 and ViT-L/14, with some supporting even larger models like ViT-H/14. In our comparisons, we use ViT-L/14 by default. However, if a method performs better with ViT-B/16, we report the superior result.

## C.4. Free-form Queries and In-the-wild Results

**Generalizability evaluation. Quantitative.** We evaluate generalizability using free-form text. To ensure a fair comparison, we use the same superpixel segmenter as FreeDA. We prompt GPT4o three times independently to generate diverse free-form class variations (e.g.,"cat"→"small domestic feline") and then perform retrieval. Results across three runs are summarized in Table A9. Shifting to free-from text, FreeDA and ProxyCLIP experience significant performance drops, whereas *ReME* consistently outperforms them. **Qualitative.** Following FreeDA, we collect in-the-wild text and qualitatively evaluate out method. The results are shown in Fig.A7.

## C.5. Data Usage for Training-Required Methods

For training-required OVS methods using image-text pairs, they often demand extensive training. Table A3 provides the training data size for such methods, where we can observe

| Methods | Training or Fine-tuning dataset | Size |
|---|---|---|
| GroupViT[72] | CC12M+YFCC | 26 million |
| SimSeg[82] | CC15M | 15 million |
| TCL[11] | CC15M | 15 million |
| CoCu[71] | CC15M+YFCC | 29 million |
| ZeroSeg[12] | CC3M+COCO | 3.4 million |
| OVSegmentor[73] | CC4M | 4 million |
| SegCLIP[39] | CC3M+COCO | 3.4 million |
| CoDe[68] | CC15M | 15 million |
| SAM-CLIP[63] | CC15M+YFCC+IN21k | 41 million |

Table A3. Data usage for training-required OVS methods.

that millions of image-text pairs from diverse datasets are leveraged, indicating their higher computational cost.

## C.6. Comparison with Training-required Methods

Although it falls beyond our primary scope of comparison, we also evaluate our approach against training-required methods, as shown in Table A11. Our method **outperforms all approaches fine-tuned with image-text data**. When compared to methods fine-tuned with segment-text, our approach surpasses LSeg+ [22], ZegFormer [16], and ZSseg [74], but falls short compared to OVSeg [36], SAN [75], and CATSeg [13]. This performance gap is commonly observed across all training-free methods when compared to models that demand fine-tuning on segment-text.

However, it is important to note that training-free methods have significantly fewer resources: (1) no training is performed, and (2) no labor-intensive pixel-level annotations. i.e., segment-text data, are required. As a training-free method, we achieve the smallest performance gap compared to these segment-text fine-tuned models.

To sum up, our contributions remain distinct: **A. *ReME* achieves state-of-the-art performance among all training-free methods while also surpassing models trained on millions of image-text pairs**, demonstrating reduced dependence on large-scale training. **B.** Our framework provides a novel perspective on multi-modal data quality, offering contributions that extend beyond OVS.

## D. Limitation

One limitation of our framework is the decision to drop misaligned pairs in the base set rather than correcting them by reassigning appropriate labels. For instance, in Fig. 3 of the main paper, misaligned pairs where "dog" is associated with segments not depicting dogs are simply filtered out. A more sophisticated approach could involve identifying the correct segments for those labels and reassigning appropriate labels to the affected segments. This refinement would increase the diversity of the final reference set and further enhance the quality of the resulting segment-text embeddings. However, given the diversity and scale of our image

resource, COCO-2017 [8], we opt for a simpler and more efficient data enhancement phase.

In domains with limited data availability and constrained diversity [42], this limitation could be addressed easily through a plug-in component. After group-based filtering, this component could leverage intra-modality similarity to identify the closest neighbors for each element in mis-aligned pairs, enabling the estimation of correct matches with minimal computational overhead.

| Components | mIoU | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | VOC-20 | VOC-21 | City | PC-59 | PC-60 | Object | Stuff | A-150 | PC-459 | A-847 | AVG$^{10}$ |
| Base set (no enhancement) | 70.03 | 62.30 | 30.94 | 35.42 | 30.46 | 39.38 | 27.01 | 22.03 | 9.14 | 6.19 | 33.29 |
| w/ (i) Synonym-guided enriching | 79.50 | 66.91 | 33.47 | 36.69 | 34.81 | 39.92 | 28.02 | 23.41 | 9.56 | 6.22 | 35.85 |
| w/ (ii) Group-based filtering | 91.10 | 76.41 | 47.36 | 40.66 | 38.52 | 42.48 | 31.80 | 24.09 | 12.96 | 7.13 | 41.25 |
| w/ Both (i) and (ii) | **92.34** | **79.63** | **50.42** | **44.89** | **41.64** | **45.50** | **33.12** | **26.13** | **14.12** | **8.43** | **43.62** |

Table A4. Impact of data enhancement components.

| Data filtering alternatives | mIoU | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | VOC-20 | VOC-21 | City | PC-59 | PC-60 | Object | Stuff | A-150 | PC-459 | A-847 | AVG$^{10}$ |
| Global filtering*[a] | 79.34 | 71.19 | 41.37 | 39.79 | 37.25 | 40.79 | 31.16 | 21.18 | 11.71 | 7.88 | 38.47 |
| *Group-based filtering (with cross-modality CLIP score)[b] | 80.05 | 72.92 | 43.06 | 41.84 | 39.44 | 41.81 | 31.88 | 22.79 | 12.19 | 8.33 | 39.63 |
| *Group-based filtering (with intra-modality similarity score)[c] | **92.34** | **79.63** | **50.42** | 44.89 | 41.64 | 45.50 | 33.12 | 26.13 | 14.12 | 8.43 | 43.62 |
| *Group-based filtering (with intra-modality similarity score; weighted ratio)[d] | 92.26 | 79.61 | 50.38 | **44.97** | **41.88** | **45.60** | **33.17** | **26.51** | **14.74** | **8.58** | **43.77** |

Table A5. Analysis of different data filtering approaches.

| Feature encoder | mIoU | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | VOC-20 | VOC-21 | City | PC-59 | PC-60 | Object | Stuff | A-150 | PC-459 | A-847 | AVG$^{10}$ |
| CLIP | 91.61 | 68.77 | 38.53 | 36.51 | 35.08 | 39.82 | 26.85 | 24.72 | 13.76 | 7.51 | 37.81 |
| DINOv2$_B$ | 91.72 | 79.13 | 50.20 | 43.65 | 41.37 | 44.71 | 32.58 | 25.29 | 13.79 | 7.68 | 43.01 |
| DINOv2$_L$ | **92.34** | **79.63** | **50.42** | **44.89** | **41.64** | **45.50** | **33.12** | **26.13** | **14.12** | **8.43** | **43.62** |

Table A6. Analysis of feature encoder variations.

| Captioners | mIoU | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | VOC-20 | VOC-21 | City | PC-59 | PC-60 | Object | Stuff | A-150 | PC-459 | A-847 | AVG$^{10}$ |
| LLaVA [38] | **92.34** | **79.63** | **50.42** | **44.89** | **41.64** | **45.50** | **33.12** | **26.13** | **14.12** | **8.43** | **43.62** |
| BLIP-2 [34] | 89.41 | 56.32 | 40.06 | 40.85 | 38.42 | 37.64 | 30.76 | 24.31 | 12.42 | 6.47 | 37.67 |
| GT Caption | 89.02 | 55.57 | 40.19 | 40.15 | 38.37 | 37.77 | 29.68 | 24.09 | 11.64 | 5.37 | 37.18 |

Table A7. Ablation study of the image description generator.

| Segmenters | mIoU | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | VOC-20 | VOC-21 | City | PC-59 | PC-60 | Object | Stuff | A-150 | A-847 | PC-459 | AVG$^{10}$ |
| Superpixel [20] | 92.34 | 79.63 | 50.42 | 44.89 | 41.64 | 45.50 | 33.12 | 26.13 | 14.12 | 8.43 | 43.62 |
| SAM [28] | 93.15 | 82.20 | 59.04 | **53.10** | **44.58** | 48.21 | 33.32 | 28.21 | 15.82 | 8.80 | 46.64 |
| SAM2 [52] | **93.18** | **82.26** | **61.19** | 52.03 | 43.42 | **48.40** | **33.36** | **28.21** | 8.83 | 15.97 | **46.69** |

Table A8. Ablation study of the segmenter.

| Methods | mIoU | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PC-59 | PC-59* | Δ(%) | A-150 | A-150* | Δ(%) | PC-459 | PC-459* | Δ(%) | A-847 | A-847* | Δ(%) |
| Ours | **44.89** | **42.89**±0.9 | ↓ **4.46** | **26.13** | **26.12**±0.3 | ↓ **0.04** | **14.12** | **13.14**±0.1 | ↓ **6.94** | **8.43** | **7.35**±0.1 | ↓ 12.81 |
| FreeDA [6] | 43.50 | 36.18±0.8 | ↓ 16.83 | 22.4 | 16.27±1.0 | ↓ 27.37 | 10.20 | 7.16±0.2 | ↓ 29.80 | 5.30 | 2.09±0.1 | ↓ 54.52 |
| ProxyCLIP [32] | 37.7 | 33.15±1.2 | ↓ 12.05 | 22.6 | 17.12±0.3 | ↓ 24.26 | 11.20 | 8.41±0.3 | ↓ 24.84 | 6.70 | 6.39±0.2 | ↓ **4.63** |

Table A9. Generalizability evaluation with free-form queries.

| Methods | Backbone | Post-proc | mIoU | | | | | | | | | |
|---------|----------|-----------|------|------|------|-------|-------|--------|-------|-------|--------|-------|
| | | | VOC-20 | VOC-21 | City | PC-59 | PC-60 | Object | Stuff | A-150 | PC-459 | A-847 |
| GEM [7] | CLIP | ✗ | 46.2 | 24.7 | - | 32.6 | 21.2 | - | 15.1 | 10.1 | 4.6 | 3.7 |
| MaskCLIP [87] | CLIP, DeepLabV2 | ✓ | 74.9 | 38.8 | 12.6 | 25.5 | 23.6 | 20.6 | 14.6 | 9.8 | - | - |
| ReCo [58] | CLIP, DenseCLIP | ✓ | 62.4 | 27.2 | 23.2 | 24.7 | 21.9 | 17.3 | 16.3 | 12.4 | - | - |
| SCLIP [62] | CLIP | ✓ | 83.5 | 61.7 | 34.1 | 36.1 | 31.5 | 32.1 | 23.9 | 17.8 | 9.3 | 6.1 |
| CaR [60] | CLIP | ✓ | 91.4 | 67.6 | 15.1 | 39.5 | 30.5 | 36.6 | 11.2 | 17.7 | 11.5 | 5.0 |
| NACLIP [25] | CLIP | ✓ | 83.0 | 64.1 | 38.3 | 38.4 | 35.0 | 36.2 | 25.7 | 19.1 | 9.0 | 6.5 |
| CLIPtrase [57] | CLIP | ✓ | 81.2 | 53.0 | 21.1 | 34.9 | 30.8 | 39.6 | 24.1 | 17.0 | 9.9 | 5.9 |
| PnP [40] | CLIP, GPT4om, BLIP | ✓ | 79.1 | 51.3 | 19.3 | 31.0 | 28.0 | 36.2 | 17.9 | 14.2 | 5.5 | 4.2 |
| FreeDA [6] | CLIP, Stable Diffusion, DINO | ✓ | 87.9 | 55.4 | 36.7 | 43.5 | 38.3 | 37.4 | 28.8 | 22.4 | 10.2 | 5.3 |
| ProxyCLIP [32] | CLIP, DINO | ✗ | 83.2 | 60.6 | 40.1 | 37.7 | 34.5 | 39.2 | 25.6 | 22.6 | 11.2 | 6.7 |
| DiffSegmenter [64] | Stable Diffusion, BLIP, U-Net, DeepLabV2 | ✓ | 71.4 | 60.1 | - | 27.5 | 25.1 | 37.9 | - | - | - | - |
| OVDiff [27] | CLIP, Stable Diffusion, GPT, CutLER | ✓ | 80.9 | 68.4 | 23.4 | 32.9 | 31.2 | 36.2 | 20.3 | 14.1 | 12.0 | 6.6 |
| *ReME* (Ours) | CLIP, LLaVA, DINO | ✗ | **92.3** | **79.6** | **50.4** | **44.9** | **41.6** | **45.5** | **33.1** | **26.1** | **14.1** | **8.4** |

Table A10. **Comparison to training-free methods without SAM.** The best overall results are **bolded**, with the second-best results underlined.
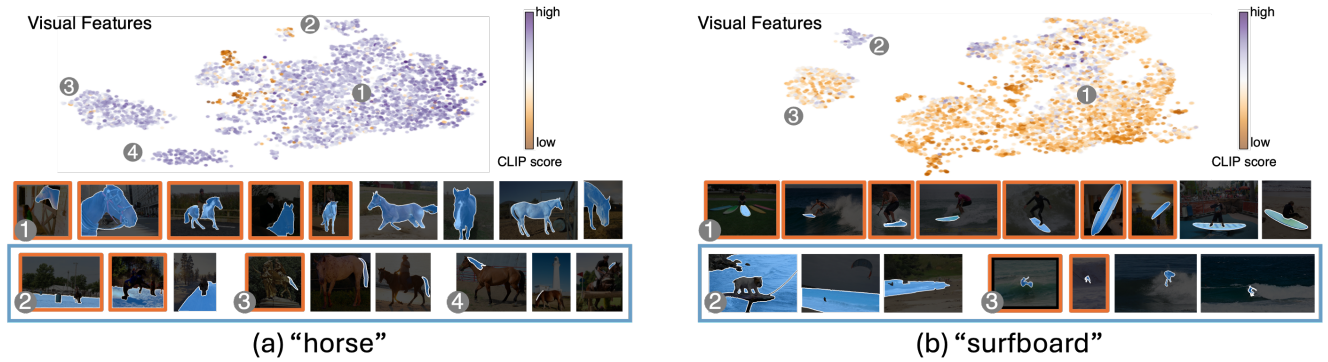


(a) "horse"   (b) "surfboard"

Figure A4. **The superiority of intra-modality over cross-modality for data issue detection.** Each figure provides a UMAP projection of segment embeddings labeled as "horse" or "surfboard", respectively, colored by cross-modal similarity scores (CLIP scores) between the segment and its corresponding label. Individual segments are shown below. Blue boxes highlight misalignments detected by our filtering; orange boxes are those detected by low CLIP scores, which remove correct pairings while leaving many misalignments unaddressed.



Figure A5. Image descriptions from different resources. Red text highlights concepts uniquely present in the LLaVA description.

| Methods | Post-processing | mIoU | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | VOC-20 | VOC-21 | City | PC-59 | PC-60 | Object | Stuff | A-150 | PC-459 | A-847 |
| *Methods that require finetuning on segment-text data* | | | | | | | | | | | |
| LSeg+[22] | ✗ | - | 59.0 | - | 36.0 | - | - | - | 13.0 | 5.2 | 2.5 |
| ZegFormer [16] | ✗ | 86.2 | 62.7 | - | 42.8 | - | - | - | 16.9 | 9.1 | 4.9 |
| ZSseg [74] | ✗ | 88.4 | - | - | 44.7 | - | - | - | 20.5 | - | 7.0 |
| OVSeg [36] | ✗ | 94.5 | - | - | 55.7 | - | - | - | 29.6 | 12.4 | 9.0 |
| SAN [75] | ✗ | 94.6 | - | - | 57.7 | - | - | - | 32.1 | 15.7 | 12.4 |
| CAT-Seg [13] | ✗ | 97.0 | 82.5 | - | 63.3 | - | - | - | 37.9 | 23.8 | 16.0 |
| *Methods that require finetuning on image-text data* | | | | | | | | | | | |
| GroupViT [72] | ✗ | 74.1 | 52.3 | 11.1 | 23.4 | 22.4 | 24.3 | 15.3 | 10.6 | 4.9 | 4.3 |
| SimSeg [82] | ✓ | 57.4 | 35.2 | 10.8 | 26.2 | 23.4 | 29.7 | 18.5 | 11.4 | 5.0 | 4.7 |
| TCL [11] | ✓ | 83.2 | 55.0 | 23.1 | 33.9 | 30.4 | 31.6 | 19.6 | 17.1 | 8.7 | 6.3 |
| CoCu [71] | ✗ | 73.0 | 51.4 | 22.1 | 26.5 | 23.6 | 22.7 | 15.2 | 12.3 | 5.1 | 4.5 |
| ZeroSeg [12] | ✗ | - | 40.8 | - | 20.4 | - | 20.2 | - | - | - | - |
| OVSegmentor [73] | ✗ | - | 53.8 | - | - | 20.4 | 25.1 | - | 5.6 | - | - |
| SegCLIP [39] | ✗ | - | 52.6 | - | - | 24.7 | 26.5 | - | 8.7 | - | - |
| CoDe [68] | ✓ | 57.7 | - | 28.9 | 30.5 | - | 32.3 | 23.9 | 17.7 | - | - |
| SAM-CLIP [63] | ✗ | - | 60.6 | - | - | 29.2 | - | 31.5 | 17.1 | - | - |
| *Training-free Methods without SAM* | | | | | | | | | | | |
| GEM [7] | ✗ | 46.2 | 24.7 | - | 32.6 | 21.2 | - | 15.1 | 10.1 | 4.6 | 3.7 |
| MaskCLIP [87] | ✓ | 74.9 | 38.8 | 12.6 | 25.5 | 23.6 | 20.6 | 14.6 | 9.8 | - | - |
| ReCo [58] | ✓ | 62.4 | 27.2 | 23.2 | 24.7 | 21.9 | 17.3 | 16.3 | 12.4 | - | - |
| SCLIP [62] | ✓ | 83.5 | 61.7 | 34.1 | 36.1 | 31.5 | 32.1 | 23.9 | 17.8 | 9.3 | 6.1 |
| CaR [60] | ✓ | 91.4 | 67.6 | 15.1 | 39.5 | 30.5 | 36.6 | 11.2 | 17.7 | 11.5 | 5.0 |
| NACLIP [25] | ✓ | 83.0 | 64.1 | 38.3 | 38.4 | 35.0 | 36.2 | 25.7 | 19.1 | 9.0 | 6.5 |
| CLIPtrase [57] | ✓ | 81.2 | 53.0 | 21.1 | 34.9 | 30.8 | 39.6 | 24.1 | 17.0 | 9.9 | 5.9 |
| PnP [40] | ✓ | 79.1 | 51.3 | 19.3 | 31.0 | 28.0 | 36.2 | 17.9 | 14.2 | 5.5 | 4.2 |
| FreeDA [6] | ✓ | 87.9 | 55.4 | 36.7 | 43.5 | 38.3 | 37.4 | 28.8 | 22.4 | 10.2 | 5.3 |
| ProxyCLIP [32] | ✗ | 83.2 | 60.6 | 40.1 | 37.7 | 34.5 | 39.2 | 25.6 | 22.6 | 11.2 | 6.7 |
| DiffSegmenter [64] | ✓ | 71.4 | 60.1 | - | 27.5 | 25.1 | 37.9 | - | - | - | - |
| OVDiff [27] | ✓ | 80.9 | 68.4 | 23.4 | 32.9 | 31.2 | 36.2 | 20.3 | 14.1 | 12.0 | 6.6 |
| ***ReME (Ours)*** | ✗ | **92.3** | **79.6** | **50.4** | **44.9** | **41.6** | **45.5** | **33.1** | **26.1** | **14.1** | **8.4** |
| *Training-free Methods with SAM* | | | | | | | | | | | |
| RIM [66] | ✗ | 77.8 | - | - | 34.3 | - | 44.5 | - | - | - | - |
| CaR w/ SAM [60] | ✗ | - | 70.2 | 16.9 | 40.5 | 31.1 | 37.6 | 12.4 | 17.9 | 11.8 | 5.7 |
| CLIPtrase w/ SAM [57] | ✗ | 82.3 | 57.1 | - | 36.4 | 32.0 | 44.2 | 24.8 | 17.2 | 10.6 | 6.0 |
| ProxyCLIP w/ SAM [32] | ✗ | 80.4 | 59.3 | 37.0 | 37.0 | 33.6 | 35.4 | 25.0 | 19.1 | 6.9 | 4.8 |
| CorrCLIP [83] | ✗ | 91.6 | 74.1 | 47.7 | 45.5 | 40.3 | 43.6 | 30.6 | - | - | - |
| ***ReME (Ours)* w/ SAM** | ✗ | **93.2** | **82.2** | **59.0** | **53.1** | **44.6** | **48.2** | **33.3** | **28.2** | **15.8** | **8.8** |

Table A11. **Comparison to state-of-the-art OVS approaches.** The best overall results are **bolded**, with the second-best results underlined. We also analyze data robustness by varying the image resources of our reference set from the default COCO-2017 to VOC and ADE, respectively, where leading performances over baselines are **bolded**.
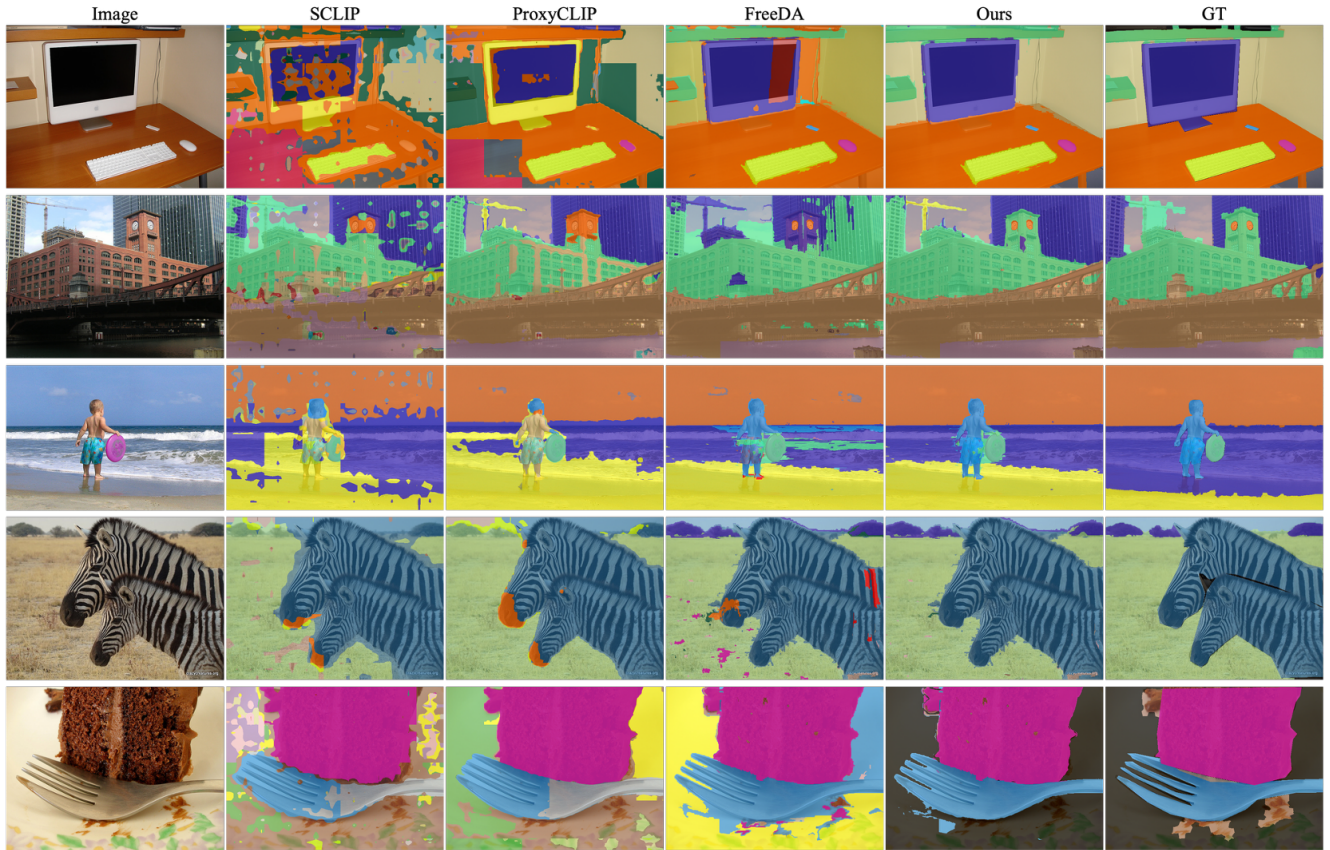
Figure A6. **Qualitative results of *ReME* in comparison with other training-free OVS methods.** SCLIP is based on CLIP attention; ProxyCLIP enhances CLIP attention with DINO features; FreeDA and *ReME* are retrieval-based methods, adopting the same superpixel-algorithm [20] for class-agnostic segmentation. We observe increasing quality of OVS results from left to right, with less noise in both masks and assigned labels.
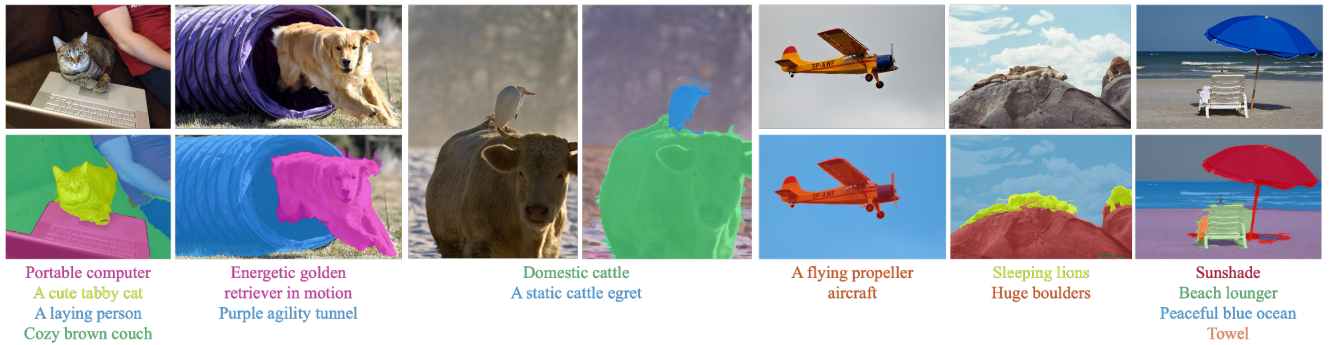


Figure A7. **In-the-wild segmentation results obtained by prompting *ReME* with diverse free-form textual inputs.**
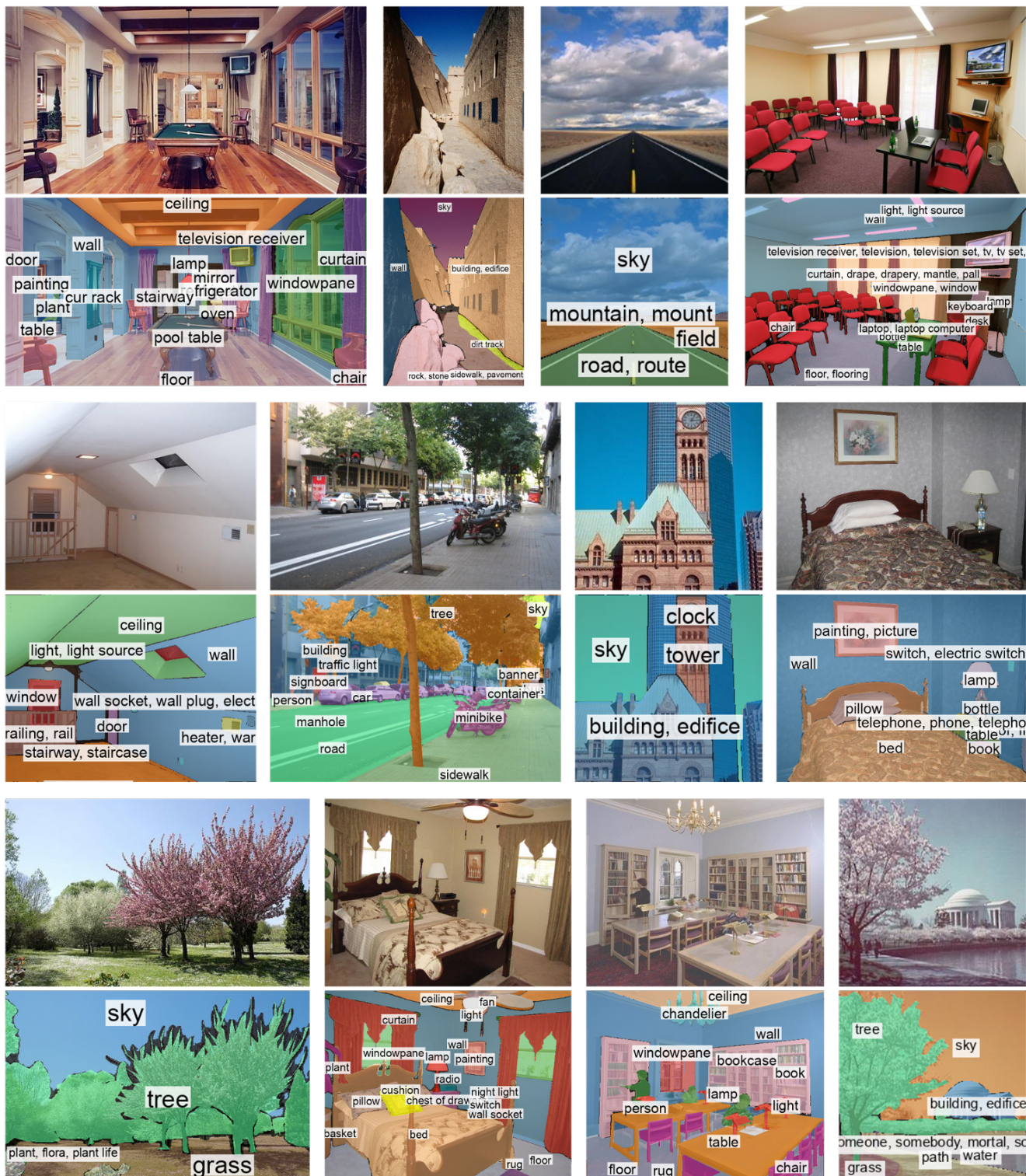
Figure A8. **Qualitative results on ADE20K [86] with 847 categories.**

Figure A9. **Qualitative results on Pascal Context [43] with 459 categories.**

Figure A10. **Qualitative results on COCO Stuff [8] with 171 categories.**