



# Semantic Library Adaptation: LoRA Retrieval and Fusion for Open-Vocabulary Semantic Segmentation

Reza Qorbani<sup>3\*</sup> Gianluca Villani<sup>1,2\*</sup> Theodoros Panagiotakopoulos<sup>1,5</sup>  
Marc Botet Colomer<sup>1</sup> Linus Härenstam-Nielsen<sup>6,7</sup> Mattia Segu<sup>1,9</sup> Pier Luigi Dovesi<sup>1,4†</sup>  
Jussi Karlgren<sup>1,4</sup> Daniel Cremers<sup>6,7</sup> Federico Tombari<sup>6,8</sup> Matteo Poggi<sup>10</sup>

<sup>1</sup>The Good AI Lab <sup>2</sup>University of Toronto <sup>3</sup>KTH <sup>4</sup>AMD Silo AI  
<sup>5</sup>King <sup>6</sup>Technical University of Munich <sup>7</sup>Munich Center for Machine Learning  
<sup>8</sup>Google <sup>9</sup>ETH Zurich <sup>10</sup>University of Bologna

<https://thegoodailab.org/semLa>

## Abstract

Open-vocabulary semantic segmentation models associate vision and text to label pixels from an undefined set of classes using textual queries, providing versatile performance on novel datasets. However, large shifts between training and test domains degrade their performance, requiring fine-tuning for effective real-world applications. We introduce **Semantic Library Adaptation (SemLA)**, a novel framework for training-free, test-time domain adaptation. SemLA leverages a library of LoRA-based adapters indexed with CLIP embeddings, dynamically merging the most relevant adapters based on proximity to the target domain in the embedding space. This approach constructs an ad-hoc model tailored to each specific input without additional training. Our method scales efficiently, enhances explainability by tracking adapter contributions, and inherently protects data privacy, making it ideal for sensitive applications. Comprehensive experiments on a 20-domain benchmark built over 10 standard datasets demonstrate SemLA’s superior adaptability and performance across diverse settings, establishing a new standard in domain adaptation for open-vocabulary semantic segmentation.

## 1. Introduction

Semantic segmentation aims to classify images at the pixel level, assigning a label to each pixel in an image. Traditionally, models are trained to recognize a fixed set of classes, but recent advances have led to *open-vocabulary* (OV) semantic segmentation, where models identify categories from an undefined and unbounded set using textual queries. This task leverages the synergy between text and

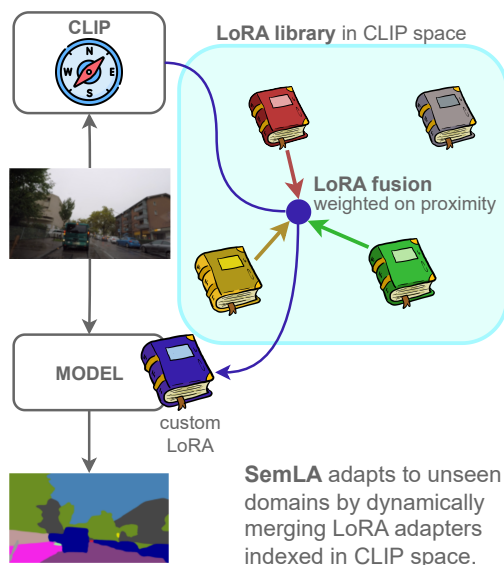


Figure 1. **Overview of SemLA.** During test-time, SemLA uses CLIP as a domain navigator, to retrieve and fuse relevant adapters, to get a LoRA tailored to the target domain.

visual embeddings [51], enabling flexible and dynamic labeling [10, 71]. However, like all models, OV semantic segmentation models are vulnerable to *domain shifts* – situations where the model encounters data distributions different from those it was trained on, leading to degraded performance. This issue is significant because it impacts both segmentation accuracy and open-vocabulary generalization capabilities, which are crucial in real-world applications.

Data is plentiful but heterogeneous – varying in labels, annotation styles, and other factors – and often contains sensitive information, limiting its accessibility. Consequently, we find ourselves in a paradox where an abundance of data does not necessarily translate into robust models capable

\* Joint first authorship

† Project Lead

of handling domain shifts [67]. Although several domain adaptation methods have been proposed, they have not been applied to OV semantic segmentation settings.

Classic domain adaptation approaches are limited: they usually focus on a single target domain, require access to source data during adaptation, involve slow and expensive processes, and often deteriorate performance on the source dataset [21, 23, 64, 68, 87]. Test-time and online adaptation address some issues but are usually too slow for real-time applications and lack explainability [3, 49]. Drawing inspiration from recent developments in Large Language Models (LLMs) and the use of *Parameter-Efficient Fine-Tuning* (PEFT) techniques enriched with metadata – such as adapters in Hugging Face and Civit-ai – we aim to transfer this idea to the field of OV semantic segmentation, by integrating external information beyond what was found in the static training data through retrieval mechanisms [84].

We propose a training-free test-time adaptation approach for OV semantic segmentation, *Semantic Library Adaptation* (**SemLA**), achieved through adapter retrieval and merging. Our approach begins with creating a set of adapters using Low-Rank Adaptation (LoRA) [26, 42]. However, simply having a collection of adapters is insufficient, as the model still can’t discern which adapter to select. Just as a collection of books requires an indexing system to become a functional library, we introduce an indexing mechanism to guide retrieval. Specifically, we employ Contrastive Language-Image Pretraining (CLIP) [51] to represent each LoRA adapter as the centroid of the CLIP embeddings of its training data. At test time, we can then select the most relevant adapters based on the proximity of the target image embedding to these centroids.

Real-world domains are, however, far too numerous to rely on single adapters for every possible scenario due to the inherent combinatorial complexity of the domain space. To address this, we aggregate the most relevant adapters based on their proximity to the target domain. Just as in a library where the exact book you seek may not be available, the necessary knowledge can still be gathered by consulting multiple related sources. By fusing several adapters based on their relevance to the target image, we effectively construct an ad-hoc model fine-tuned for each situation. Figure 1 presents an overview of SemLA at test-time.

This is training-free, scales well to any library size, and accommodates vastly different domains and label sets, with adaptation happening almost instantaneously through CLIP inference and LoRA merging. Moreover, our approach enhances explainability by identifying which adapters are most useful for certain target domains, allowing us to understand the model’s reasoning and capabilities.

Our main contributions are summarized as follows:

- **Training-free test-time adaptation for OV semantic segmentation:** We introduce the first method for adap-

tation in OV segmentation. It requires no training at test time, enabling adaptive responses to diverse input images. Our framework is simple, scalable with large adapter libraries, and backbone-agnostic.

- **Novel benchmark for OV domain adaptation:** We provide comprehensive experiments with a 20-domain benchmark, built over 10 popular datasets, showcasing performances over vastly different datasets and superior performance compared to zero-shot and naive adapter merging [35] approaches.
- **Explainability and LoRA contribution analysis:** Our approach is inherently transparent and controllable even at test time, easily scaling to new targets. The adaptation phase occurs without accessing source data, thereby protecting data privacy. We present analyses showing how adapters synergize and how we can measure their contributions and influence over the adaptation process.

## 2. Related Work

Our work intersects the following research areas:

**Semantic Segmentation.** Deep learning techniques have led to increasingly effective semantic segmentation models. Fully Convolutional Networks (FCN) [41] and SegNet [1] extended convolutional neural networks with upsampling layers (deconvolution) to produce pixel-wise predictions. Subsequent works improved both speed [46, 76] and accuracy [7–9]. Enhancements in accuracy involved enlarging the receptive field [7–9, 74, 83], designing refinement modules [20, 22, 86], introducing boundary cues [6, 14, 59], and exploiting attention mechanisms [19, 36, 65] and Vision Transformers [72, 73, 78].

With the introduction of vision-language models like CLIP [51], *open-vocabulary* (OV) semantic segmentation [10, 71] has emerged as a new approach to segmentation, where the set of classes can be arbitrarily defined at any time through textual queries. This flexibility, however, introduces additional sources of domain shifts, such as vocabulary misalignment across different domains.

**Domain Adaptation.** Domain adaptation focuses on adapting a model pre-trained on a source domain to perform well on a new, unlabeled target domain [82]. Adaptation can be performed either offline or during deployment. In the offline case, usually known as *unsupervised domain adaptation* (UDA), early approaches used style transfer strategies [16, 23, 37, 68, 75, 87] or self-training to adapt the model, for example by carefully retrieving reliable pseudo-labels for the target domain by exploiting the confidence of the model itself [43, 89, 90], class balancing [24, 88], or prototypes [5, 79, 80]. While the offline scenario usually focuses on domain shifts occurring only once, such as moving from synthetic [52] to real [12] images, the deployment approach aims to handle *continuous* adaptation to avoid catastrophic forgetting of the source domain. Some methods as-

sume availability of data from the training domain and involve replay buffers [2, 31, 32], style transfer [17, 18, 69], contrastive learning [57, 62], or adversarial learning [70]. Others consider a more constrained case, *source-free* or *test-time* adaptation, where no data from the source domain is available [56], and involve pseudo-source data generation [40], partial freezing of the original model [38], feature alignment [39], batch norm retraining through entropy minimization [63], or prototype adaptation [29]. Recently, *on-line* adaptation emerged to tackle multiple domain shifts occurring unpredictably during deployment [3, 49, 54, 58, 61]. Although these strategies are flexible for unseen domains, they introduce significant overhead at deployment. Finally, *training-free* adaptation [35, 55] has been proposed to address domain shifts by aggregating parameters of different single-domain experts. This approach relies on linear mode connectivity to uniformly merge different models. We demonstrate that a careful selection of a few meaningful experts allows for better results.

**Parameter-Efficient Fine-Tuning.** The advent of large language models (LLMs) and vision-language models (VLMs), with billions of parameters trained on vast datasets, unveiled the need for new fine-tuning paradigms to adapt them to specific use cases. Low-Rank Adaptation (LoRA) [25] was proposed to fine-tune LLMs by optimizing new sets of low-rank weights to learn residuals over the original parameters, preserving the knowledge derived from large-scale datasets while reducing computational load during fine-tuning. Advanced strategies allow dynamically setting the rank of LoRA parameters [81]. Recently, the possibility of merging different LoRAs learned for different tasks has gained popularity [27], yielding various policies for combining multiple LoRAs [42].

**Model Merging.** Model merging has emerged as a strategy to fuse knowledge from multiple sources directly in the weight space, effectively representing an alternative to ensemble and federated learning [34]. In [67], merging parameters of zero-shot and fine-tuned models improves out-of-distribution performance. Similarly, [28, 30] use linear interpolation for multi-task learning. While these methods interpolate all weights of the models, our method merges LoRA adapters, which are orders of magnitude smaller in size. Inspired by Model Soups [66], AdapterSoup [11] averages adapters in parameter space, automatically retrieving and uniformly merging them for adaptation to downstream language tasks. LoraRetriever [84] introduced a similar pipeline but instead uses instruction fine-tuning to train a retriever for selecting the most relevant adapters. Compared to these methods, SemLA performs adapter retrieval and fusion simply relying on the powerful semantic indexing capabilities of CLIP [51], which effectively serves as a *domain navigator*.

**Positioning Our Work.** Our framework falls into the

category of *domain adaptation* methods, as it actively adapts to new, unseen domains by modifying the model weights based on the input. However, we diverge from classical UDA, as our source dataset isn’t fixed but can be changed at test time. It is represented by the data used to train the model backbone and all the datasets employed in the training of the adapters composing our LoRA library. These effectively represent the points in the LoRA weight-space over which we interpolate to generate new models for every unseen target image.

### 3. Method and Framework

In this section, we introduce SemLA. It consists of two main stages: (1) *Construction of the LoRA Adapters Library*, and (2) *Dynamic Test-Time Adaptation*. We begin by detailing the backbone architecture and the integration of Low-Rank Adaptation (LoRA) into it, followed by the description of our library creation and the adaptation process.

#### 3.1. Backbone Architecture with LoRA Integration

Our framework builds on the **CAT-Seg** architecture [10], a state-of-the-art model designed for OV semantic segmentation. It aligns visual and textual information by mapping images and textual queries into a shared semantic space using CLIP embeddings [51]. This allows the model to perform segmentation without being constrained to a predefined set of classes, enabling flexible and dynamic labeling. SemLA can be deployed with any model enjoying these properties – thus not only CAT-Seg, as we will discuss later.

To facilitate efficient and adaptable domain-specific tuning, we integrate **Low-Rank Adaptation (LoRA)** [26] into CAT-Seg. LoRA introduces learnable low-rank matrices to each linear layer, allowing for parameter-efficient fine-tuning without altering the original weights. Specifically, for each linear layer with weights  $\mathbf{W} \in \mathbb{R}^{d \times k}$ , we augment it with low-rank matrices  $\mathbf{B} \in \mathbb{R}^{d \times r}$  and  $\mathbf{A} \in \mathbb{R}^{r \times k}$ , where  $r \ll \min(d, k)$ . The adapted weights become:

$$\mathbf{W}' = \mathbf{W} + \Delta\mathbf{W}, \quad \Delta\mathbf{W} = \mathbf{B}\mathbf{A}. \quad (1)$$

By training only the LoRA parameters ( $\mathbf{B}$ ,  $\mathbf{A}$ ) for each domain, we create lightweight domain-specific adapters without the need to retrain the entire model. For consistency and to simplify the merging process, we use the same rank  $r$  for all LoRA adapters.

Moreover, we denote a LoRA adapter with  $\Delta\mathcal{W}_i = (\mathcal{B}_i, \mathcal{A}_i)$ , where  $\mathcal{B}_i = \{B_{1,i}, \dots, B_{m,i}\}$  and  $\mathcal{A}_i = \{A_{1,i}, \dots, A_{m,i}\}$  denote all the LoRA parameters applied across all  $m$  linear layers of the CAT-Seg model. Given a LoRA adapter  $\Delta\mathcal{W}_i$ , the corresponding adapted CAT-Seg model is denoted by  $\text{CAT-Seg}(\Delta\mathcal{W}_i)$ , or  $\text{CAT-Seg}(\Delta\mathcal{W}_i, \mathbf{x}_t)$  if we wish to specify the dependence on the input image  $\mathbf{x}_t$ . In the following, we explain the pro-

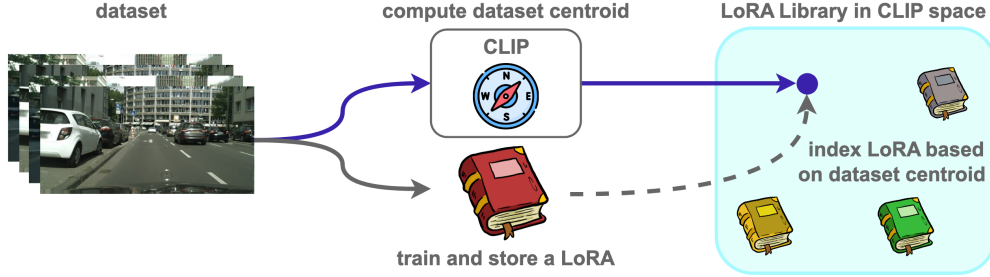


Figure 2. **Construction and Expansion of the LoRA Adapter Library.** Each LoRA adapter is created by fine-tuning on a specific dataset and subsequently added to the library. The library index for each adapter is represented by the CLIP centroid of its training data.

cedure for a single linear layer; updates for all other linear layers are performed analogously.

### 3.2. Construction of the LoRA Adapters Library

The first stage involves building a library of domain-specific LoRA adapters, each associated with a representative embedding in the CLIP space. This library serves as the foundation for our dynamic adaptation at test time. This stage, depicted in Figure 2, is performed offline and can be repeated for each newly labeled domain that becomes available for training.

#### 3.2.1. Domain Embeddings from CLIP

For each training dataset (domain)  $\mathcal{D}_i$ , we compute a centroid embedding  $\mathbf{c}_i$  to represent the domain in the CLIP embedding space as follows:

- (a) **Embedding Computation:** For each image  $\mathbf{x}_j \in \mathcal{D}_i$ , we obtain its CLIP embedding  $\mathbf{e}_j$ :

$$\mathbf{e}_j = \text{CLIP}_{\text{image}}(\mathbf{x}_j). \quad (2)$$

- (b) **Centroid Calculation:** We calculate the centroid of the embeddings from the  $N_i$  images in  $\mathcal{D}_i$ :

$$\mathbf{c}_i = \frac{1}{N_i} \sum_{j=1}^{N_i} \mathbf{e}_j, \quad (3)$$

The centroid  $\mathbf{c}_i$  captures the semantic essence of  $\mathcal{D}_i$ , facilitating efficient similarity comparisons during adaptation.

#### 3.2.2. Training Domain-Specific LoRA Adapters

For each domain  $\mathcal{D}_i$ , we fine-tune the LoRA parameters  $\Delta\mathcal{W}_i = (\mathcal{B}_i, \mathcal{A}_i)$  on CAT-Seg using  $\mathcal{D}_i$ , while keeping the original model weights  $\mathbf{W}$  frozen. This results in a domain-specific adapter  $\text{LoRA}_i = \Delta\mathcal{W}_i$ , which, when combined with the centroid  $\mathbf{c}_i$ , forms a tuple  $(\mathbf{c}_i, \Delta\mathcal{W}_i)$  representing domain  $\mathcal{D}_i$  in our library.

#### 3.2.3. Assembly of the Adapters Library

By collecting all domain-specific adapters and their centroids, we construct the **LoRA Adapters Library**:

$$\mathcal{L} = \{(\mathbf{c}_i, \Delta\mathcal{W}_i) \mid i = 1, 2, \dots, M\}, \quad (4)$$

where  $M$  is the total number of training domains. This library enables the model to dynamically adapt to various domains by selecting and merging relevant adapters.

#### 3.2.4. Extending the Library

Our approach is extremely scalable. In any point in the future, as soon as new annotated data become available for a new domain  $\mathcal{D}_*$ , we can perform steps 3.2.1 and 3.2.2 to obtain  $(\mathbf{c}_*, \Delta\mathcal{W}_*)$  and append it to  $\mathcal{L}$  to enrich it

$$\mathcal{L} = \mathcal{L} + (\mathbf{c}_*, \Delta\mathcal{W}_*) \quad (5)$$

### 3.3. Dynamic Test-Time Adaptation

The second stage focuses on adapting the model for each input image at test time by selecting and merging the most relevant adapters from the library, guided by CLIP.

#### 3.3.1. Computing the Test Image Embedding

Given a test image  $\mathbf{x}_t$ , we compute its CLIP embedding  $\mathbf{e}_t$ :

$$\mathbf{e}_t = \text{CLIP}_{\text{image}}(\mathbf{x}_t). \quad (6)$$

This embedding serves as a query to identify the most relevant adapters in the library.

#### 3.3.2. Selecting Relevant Adapters

We measure the similarity between the test image embedding and each domain centroid, e.g. as Euclidean distance:

$$d_i = \|\mathbf{e}_t - \mathbf{c}_i\|_2, \quad \forall (\Delta\mathcal{W}_i, \mathbf{c}_i) \in \mathcal{L}. \quad (7)$$

We then select the top- $K$  adapters with the smallest distances, denoted by the index set  $\mathcal{K} = \{i_1, i_2, \dots, i_K\}$ .

#### 3.3.3. Computing Adapter Weights

To quantify the relevance of each selected adapter, we apply a softmax function with temperature  $\tau$  to the proximities, computed as distance reciprocals:

$$w_i = \frac{\exp\left(\frac{1}{d_i \cdot \tau}\right)}{\sum_{k \in \mathcal{K}} \exp\left(\frac{1}{d_k \cdot \tau}\right)}, \quad \forall i \in \mathcal{K}. \quad (8)$$

These weights  $w_i$  determine the contribution of each adapter in the fusion process.

### 3.3.4. Merging Adapters via Concatenation

We specify here the merging procedure for a given linear layer adapter [42]. Since all adapters share the same rank  $r$ , the merging is straightforward:

$$\mathbf{A}'_i = w_i \mathbf{A}_i, \quad \forall i \in \mathcal{K}. \quad (9)$$

$$\mathbf{A}_{\text{fused}} = [\mathbf{A}'_{i_1}, \mathbf{A}'_{i_2}, \dots, \mathbf{A}'_{i_K}]^\top \in \mathbb{R}^{rK \times d}. \quad (10)$$

$$\mathbf{B}_{\text{fused}} = [\mathbf{B}_{i_1}, \mathbf{B}_{i_2}, \dots, \mathbf{B}_{i_K}] \in \mathbb{R}^{k \times rK}. \quad (11)$$

$$\Delta \mathbf{W}_{\text{fused}} = \mathbf{B}_{\text{fused}} \mathbf{A}_{\text{fused}}. \quad (12)$$

This fused update  $\Delta \mathbf{W}_{\text{fused}}$  integrates knowledge from multiple domains, weighted by their relevance to the test image.

### 3.3.5. Updating the Model and Prediction

We update the model weights for each linear layer:

$$\mathbf{W}' = \mathbf{W} + \Delta \mathbf{W}_{\text{fused}}. \quad (13)$$

The adapted model is then used to segment the test image:

$$\hat{\mathbf{y}}_t = \text{CAT-Seg}(\Delta \mathbf{W}_{\text{fused}}, \mathbf{x}_t). \quad (14)$$

This process is executed for each test image individually, enabling real-time, training-free adaptation.

## 4. Experiment Setup and Results

To evaluate SemLA’s performance in OV semantic segmentation under domain shifts, we design a benchmark with a diverse set of datasets covering substantial data and vocabulary variations. Unlike most of the adaptation literature, often using limited domains or fixed label spaces, our benchmark stress tests the OV capability of the network across varied scenarios, label sets, and weather conditions.

Scalability is essential for practical deployment, so we constructed a large library of LoRA adapters to show that our method scales effectively even with several adapters. This setup mirrors real-world scenarios where models must handle numerous domains and conditions without compromising performance.

### 4.1. Datasets

Our benchmark consists of **20 diverse semantic segmentation domains**, derived from 10 datasets, selected to represent a wide range of domains and to provide large data and vocabulary shifts. This includes:

- **Driving Datasets:** Commonly used in adaptation scenarios, featuring complex urban scenes with varying conditions – Cityscapes (CS) [12], BDD100K (BDD) [77], Mapillary Vistas (MV) [47], and Indian Driving Dataset (IDD) [15].
- **Weather-Specific Datasets:** Small but challenging datasets focusing on extreme weather conditions like fog, rain, and night-time scenarios – ACDC [53], MUSES [4].

- **General-Purpose Datasets:** Datasets with a vast number of classes, covering a broad spectrum of scenes and objects beyond driving scenarios – ADE20K [85] (ADE150), PC59 [44], NYU [45], and COCONutL [13].

This diverse collection ensures that our benchmark captures significant variations in both data distribution and label spaces, providing a rigorous evaluation for OV semantic segmentation.

### 4.2. Evaluation Protocol

To adhere to common conventions from the adaptation literature, we adopt a **leave-one-out** evaluation strategies

1. For each dataset, we train an adapter on its training set.
2. When evaluating on a particular dataset, we **remove** its corresponding adapter from the library, to prevent the model from leveraging any direct knowledge of it.
3. We repeat the process for each dataset, ensuring that the model is always tested on unseen data, as it never accesses to domain-specific adapters.

#### 4.2.1. Implementation Details

For our SemLA method, we select the top  $K = 7$  closest LoRA adapters based on CLIP embedding proximity, using a softmax temperature  $\tau = 0.01$  to weight the adapters during fusion. This configuration balances performance and computational efficiency. We refer to supplementary material for further details.

### 4.3. Comparative Methods

We compare the following methods:

- **Zero-Shot CAT-Seg [10]:** The baseline model without any adaptation or LoRA integration.
- **Uniform LoRA Merging:** In this approach, all LoRA adapters (excluding the one corresponding to the test dataset) are uniformly averaged without considering their relevance to the target domain – which can be seen as an adapters-variant of [35] revised to our setting with supervised fine-tuning of the adapters.
- **SemLA (Ours):** Our proposed method, which performs adaptive LoRA merging based on domain proximity.
- **Uniform (Late Fusion):** Outputs from all LoRA adapters (excluding the one corresponding to the test dataset) are uniformly averaged at the softmax output level.
- **SemLA (Late Fusion):** Similar to SemLA, weights are computed based on domain proximity but applied to the softmax outputs. This method has higher computational complexity than SemLA (one forward pass per adapter).

Furthermore, we report the performance achieved by each of the adapters trained on the specific dataset where we evaluate, although such adapters had the unfair advantage of having access to training data from the target domain – accordingly, we refer to them as **Oracles**.

Method	ACDC				MUSES								CS	BDD	MV	A150	IDD	PC59	NYU	COCONuL*	h-mean
	rain	snow	fog	night	clear (d)	clear (n)	rain (d)	rain (n)	fog (d)	fog (n)	snow (d)	snow (n)									
Zero-shot [10]	46.53	48.04	47.09	37.93	44.43	39.29	38.95	27.78	53.73	25.35	43.56	33.29	47.11	47.95	25.69	37.83	35.39	63.33	49.38	(68.26)	39.39
Oracles	70.94	69.22	69.98	51.55	69.36	57.09	54.28	52.11	75.85	61.26	66.25	54.35	67.47	60.06	49.57	53.99	64.34	68.68	61.90	(70.44)	61.05
Uniform ( <i>Late Fusion</i> )	64.58	64.76	69.64	48.51	60.53	51.88	50.62	37.30	79.00	35.68	62.39	44.04	61.82	57.08	30.34	36.04	37.43	63.24	47.39	(66.40)	49.26
SemLA ( <i>Late Fusion</i> )	66.09	66.77	71.01	50.79	62.69	55.70	53.27	45.71	68.86	45.65	65.20	49.71	62.86	57.60	29.68	37.28	39.64	63.87	50.64	(65.84)	52.09
Uniform [35]	67.40	66.35	69.71	49.98	58.28	55.78	54.70	45.09	73.75	45.16	61.02	49.08	62.18	<b>58.19</b>	<b>31.51</b>	37.25	38.83	63.06	48.93	(67.62)	51.89
SemLA (ours)	<b>67.71</b>	<b>68.95</b>	<b>71.92</b>	<b>51.73</b>	<b>61.09</b>	<b>60.06</b>	<b>57.60</b>	<b>47.35</b>	72.97	<b>52.38</b>	<b>67.28</b>	<b>55.92</b>	<b>63.91</b>	57.30	31.12	<b>38.18</b>	<b>40.16</b>	<b>64.75</b>	<b>51.35</b>	(67.26)	<b>54.16</b>
	+0.31	+2.60	+2.21	+1.75	+2.81	+4.28	+2.90	+2.26	-0.78	+7.22	+6.26	+6.84	+1.73	-0.89	-0.39	+0.93	+1.33	+1.69	+2.42	-0.36	+2.27

Table 1. **Adaptation for OV semantic segmentation – CAT-Seg [10]**. Performance comparison across our 20-domain benchmark, leave-one-out setting. On MUSES, (d) and (n) stand for *day* and *night*. ( ) means excluded from h-mean. SemLA with  $\tau = 0.05$ , and  $K = 5$ .

#### 4.4. Main Results and Discussion

Table 1 presents the mean Intersection over Union (mIoU) scores for each method across all datasets. Our method, SemLA, consistently outperforms both the zero-shot baseline and the uniform merging approach [35] (improvements over the latter are shown in the last row), reaching close to the performance of the single adapters – and, sometimes, even outperforming them (e.g., on ACDC fog and night).

One particular dataset, *CoconuL* (labeled with \*), largely overlaps with the pretraining data of CAT-Seg. Evaluating our performance on it – even in the leave-one-out protocol – would not fit the usual setting established for the adaptation task. Therefore, we exclude it from the final harmonic mean calculations. However, we provide its results in brackets for completeness, which show (unsurprisingly) an exceptional performance of the Zero-shot model.

**Performance Analysis.** Our results demonstrate that SemLA significantly outperforms the zero-shot baseline across all datasets, highlighting the effectiveness of our domain adaptation strategy. Uniform merging [35], despite its simplicity, already shows notable improvements over the zero-shot model, primarily due to linear mode connectivity, showing that LoRA merging is effective when one has access to a LoRA library without any target awareness – e.g., in domain generalization settings – as already proven in Li et al. [35]. However, the introduction of careful adapters weighting by SemLA leads to substantial performance gains. This underscores that selecting relevant adapters based on domain proximity is more important than merely scaling up the number of adapters available in the library. Moreover, we observe that both Uniform and SemLA (*Late Fusion*) perform worse than their LoRA merging counterparts, despite their higher computational complexity. These likely underperform because operating solely at the output level and cannot adjust internal model representations, unlike parameter-level fusion.

**Adapter Contribution Analysis.** To understand how the adapters contribute across datasets, we analyze the selection patterns shown in Figure 3, where a heatmap indicates the frequency and weight of each adapter’s selection per test dataset. Adapter support is spread across the library, showing that the model benefits from multiple adapters rather than relying on a few only. On the one hand, we can appreciate high weights for strongly related domains –

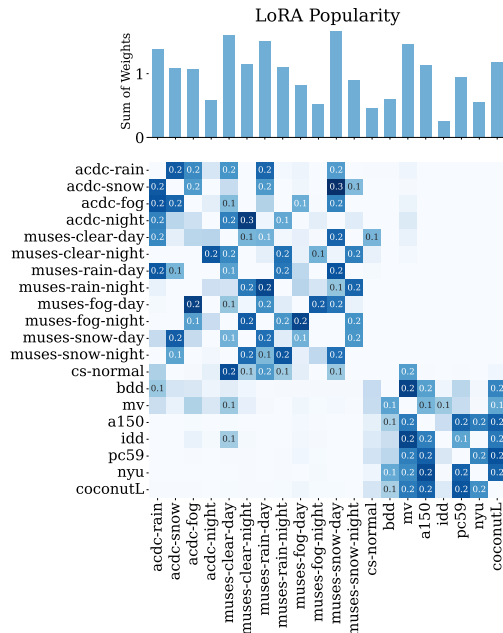


Figure 3. **Adapter contribution heatmap.** Rows represent individual test datasets, and columns correspond to specific LoRA adapters. The color intensity of each cell indicates the frequency and weight of selection (with values below 0.1 omitted). The diagonal is empty due to the leave-one-out strategy.

e.g., MUSES-snow-day heavily relies on ACDC-snow and vice-versa; in general, we can observe polarized selection for driving datasets (top left) or generic ones (bottom right). Nonetheless, even adapters from unrelated domains are often selected, suggesting the existence of shared useful features, though each category still predominantly influences related datasets. By examining columns (and aggregating them in a histogram), certain adapters emerge as “popular” choices, highlighting their stronger influence.

#### 4.5. Analysing CLIP as a Domain Navigator

A key assumption of our framework is that CLIP embeddings provide a meaningful mapping of the domains and that proximity in the CLIP space correlates with segmentation performance. To verify this assumption, we examine i) how our method navigates complex domains by combining relevant adapters; and ii) whether proximity in the CLIP space is a good heuristic for predicting performance.

Method	ACDC				MUSES								CS	BDD	MV	A150	IDD	PC59	NYU	COCOnuT*	h-mean
	rain	snow	fog	night	clear (d)	clear (n)	rain (d)	rain (n)	fog (d)	fog (n)	snow (d)	snow (n)									
Oracles	70.94	69.22	69.98	51.55	69.36	57.09	54.28	52.11	75.85	61.26	66.25	54.35	67.47	60.06	<b>49.57</b>	<b>53.99</b>	<b>64.34</b>	<b>68.68</b>	61.90	(70.44)	61.05
Uniform	67.76	66.69	69.91	50.44	58.60	56.30	55.73	45.96	73.33	46.29	61.36	50.27	62.77	58.42	32.87	38.48	39.89	63.77	50.12	(68.13)	52.78
SemLA $\tau = 0.005$	<b>70.97</b>	69.40	70.94	51.68	<b>69.51</b>	57.95	54.25	52.14	<b>75.88</b>	61.26	66.16	51.22	<b>67.52</b>	<b>60.20</b>	48.70	49.42	64.27	68.27	<b>62.68</b>	(68.83)	60.57
SemLA $\tau = 0.05$	68.82	<b>70.43</b>	<b>73.38</b>	<b>53.06</b>	67.19	<b>59.93</b>	<b>59.30</b>	<b>52.21</b>	72.45	<b>62.10</b>	<b>69.33</b>	<b>59.38</b>	67.49	58.18	38.06	44.20	49.49	67.45	57.93	(68.96)	58.79

Table 2. **Adaptation for OV semantic segmentation – CAT-Seg [10], all-inclusive setting.** Performance comparison across our 20-domain benchmark, when *all* domains are available in the LoRA Library, hence to *source*. On MUSES, (d) and (n) stand for *day* and *night*. ( ) means excluded from h-mean. SemLA with  $K = 5$ .

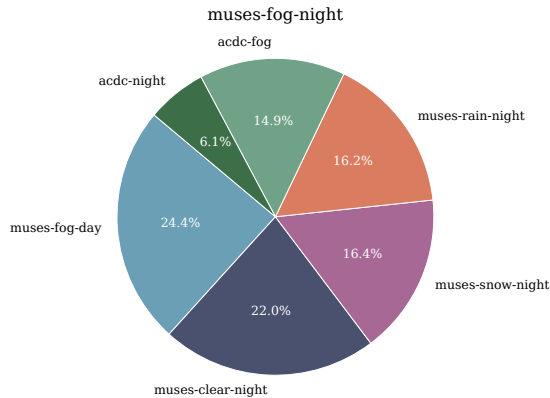


Figure 4. **Adapter weight distribution for MUSES-Fog-Night.** The fused adapter combines knowledge from foggy and night-time conditions by weighting relevant adapters. Adapters with a weight lower than 5% are not included.

**Domain Composition Analysis.** We observe that CLIP effectively navigates the domain space, as shown in the adapter selection for **MUSES-Fog-Night**. This dataset, comprising foggy night images, lacks a direct match in our library. However, adapters trained on night and fog conditions are available. Examining LoRA contributions for MUSES-Fog-Night in Figure 4, we see that fog and night adapters dominate with around 40% and 60% contribution respectively, with higher weights for MUSES adapters, around 80%, likely due to shared camera characteristics. This illustrates our ability to combine relevant knowledge even when an exact match is absent and showcases the easy interpretability of our approach.

**CLIP Distance - Performance Correlation.** Having observed that CLIP embeddings effectively guide adapter selection, we proceed to the second question: *Is image-LoRA proximity in the CLIP space a good heuristic for predicting segmentation performance?* Figure 5 indicates that this is indeed the case. The negative slope of the regression line indicates an inverse relationship between CLIP distance and mIoU – i.e., smaller CLIP distances (higher proximity) generally correspond to better segmentation. This empirical observation confirms that CLIP embedding distance is a good heuristic for predicting adapter effectiveness.

#### 4.6. All-Inclusive Settings

We evaluate the – unlikely – scenario where all the LoRA adapters corresponding to the benchmark datasets are in-

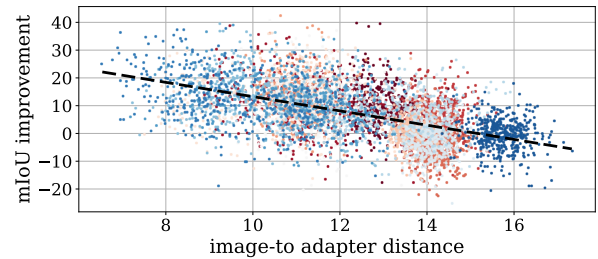


Figure 5. **CLIP-guidance effectiveness for LoRA selection on ACDC.** Each point represents an image-adapter combination, with adapters separated by color. x-axis: distance from the corresponding image embedding to the adapter embedding. y-axis: improvement in mIoU when using the adapter relative to the zero-shot base network. The linear regression curve (dashed line) indicates that embedding similarity correlates with higher mIoU. We show the full adapter library, excluding those trained on ACDC.

cluded in the library at test-time (i.e., in contrast with common adaptation settings). Table 2 reports SemLA tested with two  $\tau$  and Uniform [35] settings. The results reveal a two-fold effect: when the system is directed to prioritize the target-domain LoRA adapter, the model’s performance almost matches the Oracle. However, when the temperature  $\tau$  is higher, hence LoRA fusion is promoted, the model – while slightly lower on average – overcomes the Oracles on 14 domains over 20, primarily in highly related datasets. This finding underscores that LoRA synergies can extend beyond individual adapters, achieving performance that, in some cases, surpasses even the Oracle.

#### 4.7. Results with a Different Backbone: SED

Our method is inherently backbone-agnostic, as we only operate on LoRA retrieval and adaptation, relying on simple and widespread methods. To empirically prove this claim, we generate another library based on **SED** [71], another popular open-vocabulary semantic segmentation model. In Table 3, we collect the outcome of this new evaluation, carried out once again on our 20-domain benchmark. All our achievements with CAT-Seg are confirmed when switching to SED, with SemLA clearly outperforming both the zero-shot model and the uniform merging strategy [35], thus confirming the flexibility of our framework and the possibility of tailoring it around different OV semantic segmentation backbones.

Method	ACDC				MUSES								CS	BDD	MV	A150	IDD	PC59	NYU	COCONutL*	h-mean
	rain	snow	fog	night	clear (d)	clear (n)	rain (d)	rain (n)	fog (d)	fog (n)	snow (d)	snow (n)									
Zero-shot [10]	42.64	43.97	46.88	35.06	38.05	35.01	35.18	27.87	48.09	19.53	42.63	23.53	41.90	43.41	25.05	35.27	35.20	60.87	47.20	(66.36)	35.56
Oracles	61.70	68.25	69.33	48.39	62.04	55.55	53.61	44.82	69.40	53.20	65.60	55.85	68.48	58.61	47.28	50.66	63.15	66.55	60.60	(68.01)	58.05
Uniform [35]	59.59	61.07	68.33	47.27	58.79	53.87	50.65	37.12	<b>69.82</b>	35.52	61.46	44.08	61.85	<b>54.78</b>	<b>30.37</b>	34.92	39.59	60.82	48.42	(64.05)	48.60
<b>SemLA (ours)</b>	<b>60.01</b>	<b>67.00</b>	<b>69.96</b>	<b>49.70</b>	<b>60.98</b>	<b>54.79</b>	<b>55.15</b>	<b>38.89</b>	69.46	<b>38.25</b>	<b>66.80</b>	<b>48.89</b>	<b>64.35</b>	54.40	28.99	<b>36.91</b>	<b>41.08</b>	<b>62.23</b>	<b>51.82</b>	<b>(65.46)</b>	<b>50.57</b>
	+0.42	+5.93	+1.63	+2.43	+2.19	+0.92	+4.50	+1.77	-0.38	+2.70	+5.34	+4.81	+2.50	-0.38	-1.38	+1.99	+1.59	+1.69	+3.40	+1.41	+1.97

Table 3. **Adaptation for OV semantic segmentation – SED [71]**. Performance comparison across 20-domain benchmark, leave-one-out setting. On MUSES, (d) and (n) stand for day and night. (.) means excluded from h-mean. SemLA with  $\tau = 0.05$ , and  $K = 5$ .

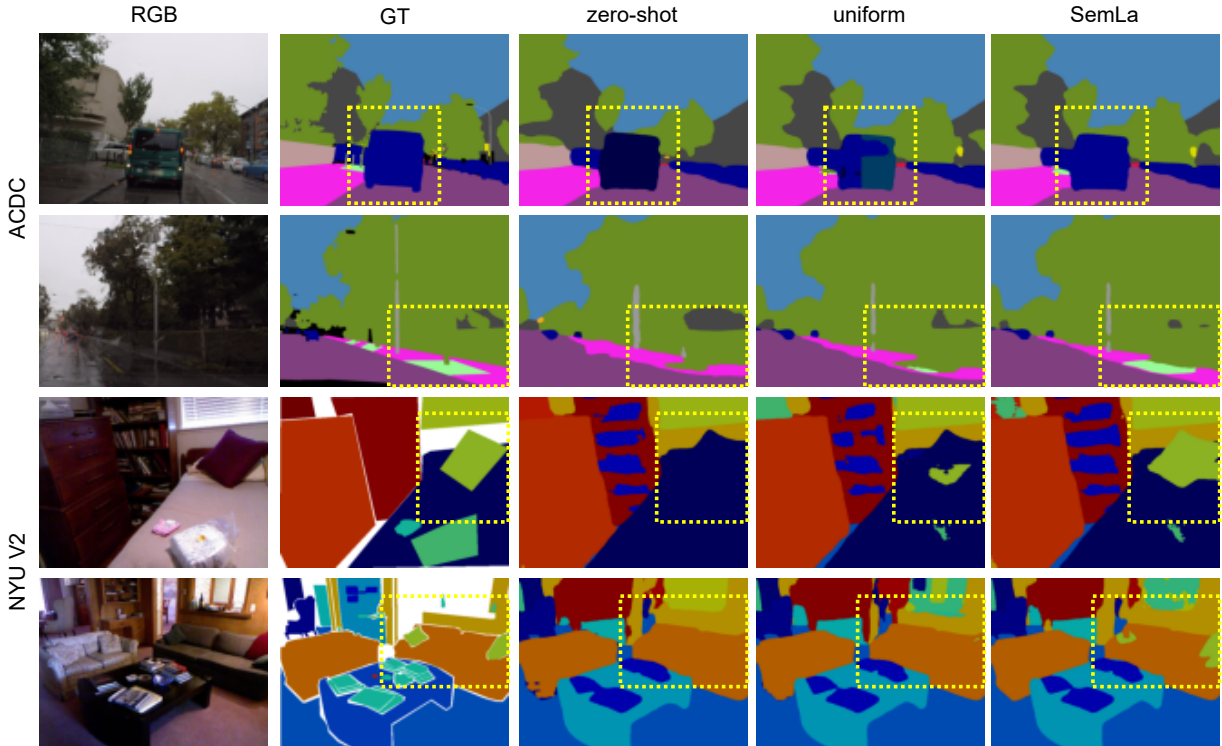


Figure 6. **Qualitative results on ACDC and NYU**. From left to right: Input Image (RGB), ground-truth semantic masks (GT) and predictions by the Zero-Shot model [10], Uniform Merging [35], and SemLA.

#### 4.8. Qualitative Results

Figure 6 presents qualitative comparisons between the predictions by different methods – zero-shot CAT-Seg, uniform merged model [35] and ours. SemLA produces more accurate and detailed segmentation masks, demonstrating its superiority over zero-shot and uniform merging models.

#### 4.9. Limitations and Future Work

While our method demonstrates strong performance, several areas remain for future exploration. First, we do not explicitly address vocabulary alignment between source and target domains; incorporating vocabulary matching could enhance performance by ensuring compatibility with the target label space. Second, exploring automated domain discovery through CLIP clustering and unsupervised training of new LoRA adapters could lead to a self-augmenting library, enhancing adaptability without manual intervention. Scaling SemLA to thousands of adapters introduces challenges in recognizing and addressing library weaknesses and gaps; developing automated strategies to identify and

mitigate them will be essential. Lastly, applying our approach to multi-tasks settings, *e.g.* panoptic segmentation or depth estimation [50], could validate its versatility.

### 5. Conclusion

We have presented SemLA, a framework for training-free test-time adaptation in open-vocabulary semantic segmentation. By constructing a library of LoRA adapters indexed with CLIP embeddings and dynamically merging them based on domain proximity, our method effectively mitigates domain shifts. Extensive experiments show that SemLA outperforms zero-shot baselines, uniform adapter merging, and sometimes even fine-tuned models across diverse datasets. We confirmed that CLIP embeddings reliably guide adapters selection and that proximity in the embedding space is a good heuristic for performance. With enhanced explainability, scalability, and data privacy, SemLA is suitable for practical applications. Future work will address vocabulary alignment, automated domain discovery, and extending the framework to other tasks.



**Acknowledgments.** The authors thank Hedvig Kjellström for the helpful discussions and guidance, and acknowledge The European High Performance Computing Joint Undertaking (EuroHPC JU), EuroCC National Competence Center Sweden (ENCCS) and the CINECA award under the IS CRA initiative for the availability of high-performance computing resources and support.

This study was funded by the European Union – Next Generation EU within the framework of the National Recovery and Resilience Plan NRRP – Mission 4 “Education and Research” – Component 2 - Investment 1.1 “National Research Program and Projects of Significant National Interest Fund (PRIN)” (Call D.D. MUR n. 104/2022) – PRIN2022 – Project reference: “River-Watch: a citizen-science approach to river pollution monitoring” (ID: 2022MMBA8X, CUP: J53D23002260006).

## References

- [1] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481–2495, 2017. 2
- [2] Andreea Bobu, Judy Hoffman, Eric Tzeng, and Trevor Darrell. Adapting to continuously shifting domains. In *ICLR 2018 Workshop Program Chairs*, 2018. 00000. 3
- [3] Marc Botet Colomer, Pier Luigi Dovesi, Theodoros Panagiotakopoulos, Joao Frederico Carvalho, Linus Härenstam-Nielsen, Hossein Azizpour, Hedvig Kjellström, Daniel Cremers, and Matteo Poggi. To adapt or not to adapt? real-time adaptation for semantic segmentation. In *IEEE International Conference on Computer Vision*, 2023. ICCV. 2, 3, 6
- [4] Tim Brödermann, David Bruggemann, Christos Sakaridis, Kevin Ta, Odysseas Liagouris, Jason Corkill, and Luc Van Gool. Muses: The multi-sensor semantic perception dataset for driving under uncertainty. *arXiv preprint arXiv:2401.12761*, 2024. 5
- [5] Chaoqi Chen, Weiping Xie, Wenbing Huang, Yu Rong, Xinghao Ding, Yue Huang, Tingyang Xu, and Junzhou Huang. Progressive feature alignment for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 627–636, 2019. 2
- [6] Liang-Chieh Chen, Jonathan T Barron, George Papandreou, Kevin Murphy, and Alan L Yuille. Semantic image segmentation with task-specific edge detection using cnns and a discriminatively trained domain transform. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4545–4554, 2016. 2
- [7] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. 2
- [8] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(4):834–848, 2018.
- [9] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018. 2
- [10] Seokju Cho, Heeseong Shin, Sunghwan Hong, Anurag Arnab, Paul Hongsuck Seo, and Seungryong Kim. Cat-seg: Cost aggregation for open-vocabulary semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4113–4123, 2024. 1, 2, 3, 5, 6, 7, 8
- [11] Alexandra Chronopoulou, Matthew Peters, Alexander Fraser, and Jesse Dodge. AdapterSoup: Weight averaging to improve generalization of pretrained language models. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 2054–2063, Dubrovnik, Croatia, 2023. Association for Computational Linguistics. 3
- [12] 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 *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 5
- [13] Xueqing Deng, Qihang Yu, Peng Wang, Xiaohui Shen, and Liang-Chieh Chen. Coconut: Modernizing coco segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 5
- [14] Henghui Ding, Xudong Jiang, Ai Qun Liu, Nadia Magnenat Thalmann, and Gang Wang. Boundary-aware feature propagation for scene segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6819–6829, 2019. 2
- [15] Shubham Dokania, AH Hafez, Anbumani Subramanian, Manmohan Chandraker, and CV Jawahar. Idd-3d: Indian driving dataset for 3d unstructured road scenes. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 4482–4491, 2023. 5
- [16] A. Dundar, M. Y. Liu, Z. Yu, T. C. Wang, J. Zedlewski, and J. Kautz. Domain stylization: A fast covariance matching framework towards domain adaptation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2020. 2
- [17] Qi Fan, Mattia Segu, Yu-Wing Tai, Fisher Yu, Chi-Keung Tang, Bernt Schiele, and Dengxin Dai. Normalization perturbation: A simple domain generalization method for real-world domain shifts. *arXiv preprint arXiv:2211.04393*, 2022. 3
- [18] Qi Fan, Mattia Segu, Yu-Wing Tai, Fisher Yu, Chi-Keung Tang, Bernt Schiele, and Dengxin Dai. Towards robust object detection invariant to real-world domain shifts. In *The Eleventh International Conference on Learning Representations (ICLR 2023)*. OpenReview, 2023. 3
- [19] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3146–3154, 2019. 2

- [20] Jun Fu, Jing Liu, Yuhang Wang, Yong Li, Yongjun Bao, Jinhui Tang, and Hanqing Lu. Adaptive context network for scene parsing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6748–6757, 2019. 2
- [21] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016. 2
- [22] Junjun He, Zhongying Deng, Lei Zhou, Yali Wang, and Yu Qiao. Adaptive pyramid context network for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7519–7528, 2019. 2
- [23] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. CyCADA: Cycle-consistent adversarial domain adaptation. In *Proceedings of the 35th International Conference on Machine Learning*, pages 1989–1998, Stockholmsmässan, Stockholm Sweden, 2018. PMLR. 2
- [24] Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Daformer: Improving network architectures and training strategies for domain-adaptive semantic segmentation. *arXiv preprint arXiv:2111.14887*, 2021. 2
- [25] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 3
- [26] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. 2, 3
- [27] Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. Lorahub: Efficient cross-task generalization via dynamic lora composition. *arXiv preprint arXiv:2307.13269*, 2023. 3
- [28] Gabriel Ilharco, Mitchell Wortsman, Samir Yitzhak Gadre, Shuran Song, Hannaneh Hajishirzi, Simon Kornblith, Ali Farhadi, and Ludwig Schmidt. Patching open-vocabulary models by interpolating weights. *Advances in Neural Information Processing Systems*, 35:29262–29277, 2022. 3
- [29] Yusuke Iwasawa and Yutaka Matsuo. Test-time classifier adjustment module for model-agnostic domain generalization. In *Advances in Neural Information Processing Systems*, 2021. 3
- [30] Dong-Hwan Jang, Sangdoon Yun, and Dongyoon Han. Model stock: All we need is just a few fine-tuned models. In *European Conference on Computer Vision*, pages 207–223. Springer, 2025. 3
- [31] Yevhen Kuznietsov, Marc Proesmans, and Luc Van Gool. Towards unsupervised online domain adaptation for semantic segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 261–271, 2022. 3
- [32] Qicheng Lao, Xiang Jiang, Mohammad Havaei, and Yoshua Bengio. Continuous domain adaptation with variational domain-agnostic feature replay. *arXiv preprint arXiv:2003.04382*, 2020. 3
- [33] 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, 4
- [34] Weishi Li, Yong Peng, Miao Zhang, Liang Ding, Han Hu, and Li Shen. Deep model fusion: A survey, 2023. 3
- [35] Wenyi Li, Huan-ang Gao, Mingju Gao, Beiwen Tian, Rong Zhi, and Hao Zhao. Training-free model merging for multi-target domain adaptation. In *European Conference on Computer Vision*. Springer, 2024. 2, 3, 5, 6, 7, 8
- [36] Xia Li, Zhisheng Zhong, Jianlong Wu, Yibo Yang, Zhouchen Lin, and Hong Liu. Expectation-maximization attention networks for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9167–9176, 2019. 2
- [37] Yunsheng Li, Lu Yuan, and Nuno Vasconcelos. Bidirectional learning for domain adaptation of semantic segmentation. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [38] Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. *CoRR*, 2020. 3
- [39] Yuejiang Liu, Parth Kothari, Bastien Germain van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, and Alexandre Alahi. TTT++: When does self-supervised test-time training fail or thrive? In *Advances in Neural Information Processing Systems*, 2021. 3
- [40] Yuang Liu, Wei Zhang, and Jun Wang. Source-free domain adaptation for semantic segmentation. 2021. 3
- [41] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 2
- [42] Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>, 2022. 2, 3, 5
- [43] Ke Mei, Chuang Zhu, Jiaqi Zou, and Shanghang Zhang. Instance adaptive self-training for unsupervised domain adaptation. *Lecture Notes in Computer Science*, page 415–430, 2020. 2
- [44] 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. 5
- [45] Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *ECCV*, 2012. 5
- [46] Vladimir Nekrasov, Chunhua Shen, and Ian Reid. Lightweight refinenet for real-time semantic segmentation. In *British Conference on Computer Vision (BMVC)*, 2018. 2

- [47] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Bulo, and Peter Kotschieder. The mapillary vistas dataset for semantic understanding of street scenes. In *Proceedings of the IEEE international conference on computer vision*, pages 4990–4999, 2017. 5
- [48] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*. 3, 5, 6
- [49] Theodoros Panagiotakopoulos, Pier Luigi Dovesi, Linus Härenstam-Nielsen, and Matteo Poggi. Online domain adaptation for semantic segmentation in ever-changing conditions. In *European Conference on Computer Vision (ECCV)*, 2022. 2, 3
- [50] Matteo Poggi and Fabio Tosi. Federated online adaptation for deep stereo. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 8
- [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, 2, 3, 5, 6
- [52] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *European conference on computer vision*, pages 102–118. Springer, 2016. 2
- [53] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Acdc: The adverse conditions dataset with correspondences for semantic driving scene understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10765–10775, 2021. 5
- [54] Mattia Segu, Bernt Schiele, and Fisher Yu. Darth: Holistic test-time adaptation for multiple object tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9717–9727, 2023. 3
- [55] Mattia Segu, Alessio Tonioni, and Federico Tombari. Batch normalization embeddings for deep domain generalization. *Pattern Recognition*, 135:109115, 2023. 3
- [56] Serban Stan and Mohammad Rostami. Unsupervised model adaptation for continual semantic segmentation. In *AAAI*, 2021. 3
- [57] Peng Su, Shixiang Tang, Peng Gao, Di Qiu, Ni Zhao, and Xiaogang Wang. Gradient regularized contrastive learning for continual domain adaptation. 2020. 00000. 3
- [58] Tao Sun, Mattia Segu, Janis Postels, Yuxuan Wang, Luc Van Gool, Bernt Schiele, Federico Tombari, and Fisher Yu. SHIFT: a synthetic driving dataset for continuous multi-task domain adaptation. In *Computer Vision and Pattern Recognition*, 2022. 3
- [59] Towaki Takikawa, David Acuna, Varun Jampani, and Sanja Fidler. Gated-scnn: Gated shape cnns for semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5229–5238, 2019. 2
- [60] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9 (11), 2008. 3
- [61] Riccardo Volpi, Pau de Jorge, Diane Larlus, and Gabriela Csurka. On the road to online adaptation for semantic image segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [62] Vibashan VS, Poojan Oza, and Vishal M. Patel. Towards online domain adaptive object detection. *2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2023. 3
- [63] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. In *International Conference on Learning Representations*, 2021. 3
- [64] Haoran Wang, Tong Shen, Wei Zhang, Lingyu Duan, and Tao Mei. Classes matter: A fine-grained adversarial approach to cross-domain semantic segmentation. In *The European Conference on Computer Vision (ECCV)*, 2020. 2
- [65] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. 2
- [66] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig Schmidt. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *Proceedings of the 39th International Conference on Machine Learning*, pages 23965–23998. PMLR, 2022. 3
- [67] Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. Robust fine-tuning of zero-shot models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7959–7971, 2022. 2, 3, 6
- [68] Zuxuan Wu, Xintong Han, Yen-Liang Lin, Mustafa Gökhan Uzunbas, Tom Goldstein, Ser Nam Lim, and Larry S. Davis. Dcan: Dual channel-wise alignment networks for unsupervised scene adaptation. *Lecture Notes in Computer Science*, page 535–552, 2018. 2
- [69] Zuxuan Wu, Xin Wang, Joseph Gonzalez, Tom Goldstein, and Larry Davis. ACE: Adapting to changing environments for semantic segmentation. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2121–2130. IEEE, 2019. 3
- [70] Markus Wulfmeier, Alex Bewley, and Ingmar Posner. Incremental adversarial domain adaptation for continually changing environments. 2018. 00000. 3
- [71] Bin Xie, Jiale Cao, Jin Xie, Fahad Shahbaz Khan, and Yanwei Pang. Sed: A simple encoder-decoder for open-vocabulary semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 1, 2, 7, 8
- [72] Enze Xie, Wenjia Wang, Wenhai Wang, Peize Sun, Hang Xu,

- Ding Liang, and Ping Luo. Segmenting transparent objects in the wild with transformer. In *IJCAI*, 2021. [2](#)
- [73] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021. [2](#)
- [74] Maoge Yang, Kun Yu, Chi Zhang, Zhiwei Li, and Kuiyuan Yang. Denseaspp for semantic segmentation in street scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3684–3692, 2018. [2](#)
- [75] Yanchao Yang and Stefano Soatto. FDA: Fourier domain adaptation for semantic segmentation. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4084–4094. IEEE, 2020. [2](#)
- [76] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 325–341, 2018. [2](#)
- [77] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2636–2645, 2020. [5](#)
- [78] Yuhui Yuan, Xiaokang Chen, Xilin Chen, and Jingdong Wang. Segmentation transformer: Object-contextual representations for semantic segmentation. In *European Conference on Computer Vision (ECCV)*, 2020. [2](#)
- [79] Pan Zhang, Bo Zhang, Ting Zhang, Dong Chen, Yong Wang, and Fang Wen. Prototypical pseudo label denoising and target structure learning for domain adaptive semantic segmentation. 2021. [2](#)
- [80] Qiming Zhang, Jing Zhang, Wei Liu, and Dacheng Tao. Category anchor-guided unsupervised domain adaptation for semantic segmentation. *Advances in Neural Information Processing Systems*, 2019. [2](#)
- [81] Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In *The Eleventh International Conference on Learning Representations*, 2023. [3](#)
- [82] Youshan Zhang. A survey of unsupervised domain adaptation for visual recognition. *arXiv:2112.06745*, 2021. [2](#)
- [83] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *CVPR*, 2017. [2](#)
- [84] Ziyu Zhao, Leilei Gan, Guoyin Wang, Wangchunshu Zhou, Hongxia Yang, Kun Kuang, and Fei Wu. LoraRetriever: Input-aware LoRA retrieval and composition for mixed tasks in the wild. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 4447–4462, Bangkok, Thailand, 2024. Association for Computational Linguistics. [2](#), [3](#)
- [85] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017. [5](#)
- [86] Yizhou Zhou, Xiaoyan Sun, Zheng-Jun Zha, and Wenjun Zeng. Context-reinforced semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4046–4055, 2019. [2](#)
- [87] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017. [2](#)
- [88] Yang Zou, Zhiding Yu, BVK Kumar, and Jinsong Wang. Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In *Proceedings of the European conference on computer vision (ECCV)*, pages 289–305, 2018. [2](#)
- [89] Yang Zou, Zhiding Yu, Xiaofeng Liu, BVK Kumar, and Jinsong Wang. Confidence regularized self-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5982–5991, 2019. [2](#)
- [90] Yang Zou, Zhiding Yu, Xiaofeng Liu, B. V. K. Vijaya Kumar, and Jinsong Wang. Confidence regularized self-training. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. [2](#)



# Semantic Library Adaptation: LoRA Retrieval and Fusion for Open-Vocabulary Semantic Segmentation

## Supplementary Material



Figure 7. **t-SNE visualization of the LoRA Library and relative datasets in the CLIP embedding space.** Similar domains are clustered together, indicating areas with higher LoRA support and potentially stronger performance improvements. A maximum of 15,000 samples per dataset is used for the t-SNE fit and only 15% of the total datapoints are used for plotting.

## Introduction

In this supplementary document, we provide additional details and analyses to support and extend the findings presented in the main paper. Here we report additional details about the implementation specifics, auxiliary experimental results, extensive ablation studies, discuss practical considerations for real-world deployment, and offer additional qualitative examples that highlight the effectiveness and robustness of our proposed method, SemLA. The rest of this document is organized as follows:

- **Section A** details the implementation aspects of our method, including training procedures, hyperparameters, model architecture modifications, and code availability for replication purposes.
- **Section B** provides an in-depth analysis of the LoRA adapter library, including visualizations of the embeddings from training samples using t-SNE, labeling of the adapters with natural language using BLIP-2 [33] for the sake of interpretability, analysis of adapter contributions dataset by dataset, and support score analysis.
- **Section C** presents extensive ablation studies and performance analyses. We explore the effectiveness of fully fine-tuned models versus LoRA adapters, conduct hyperparameter sensitivity analysis, and assess alternative domain navigators such as DINOv2 [48] versus CLIP.
- **Section D** showcases additional qualitative results across various domains, further demonstrating the adaptability and efficacy of our approach.
- **Section E** discusses pragmatic considerations for real-world deployment of SemLA, including strategies to handle computational overhead, domain navigation in specialized domains, scalability concerns, and methods to ensure efficiency and reliability in production environments.

## A. Implementation Details

**LoRA Training Details.** We attach LoRAs to every  $mn$ .Linear layer in the CAT-Seg architecture, except for CLIPs token embedding layer, as this parameter was not trained in the original CAT-Seg implementation either. All LoRAs are trained with the same LoRA configuration – i.e., rank  $r=8$  and  $\alpha=16$ . The training hyperparameters are largely the same as CATSeg with minor modifications: compared to the original CAT-Seg implementation, we use a base learning rate of  $1e-4$ , weight decay of  $1e-5$ , and 1000 warm-up iterations. For ACDC and MUSES adapters we use a batch size of 2 and a warm-up factor of 0.01. For BDD and CS we increase the batch size to 4 while keeping other hyperparameters the same. For the remaining datasets, we increased the warm-up factor to 0.1 while keeping other hyperparameters the same. All the adapters were trained until convergence.

**Code and Models.** Full source code and documentation are available in our project page <https://thegoodailab.org/semLa>.

## B. Interpretability of the LoRA Library

### B.1. t-SNE Visualization of the LoRA Library

Figure 7 presents a t-SNE [60] visualization of the LoRA adapters’ centroids and their associated datasets in the CLIP embedding space. This visualization illustrates the distribution and relationships among different adapters and domains, highlighting similarities between them.

We observe that domains with similar visual characteristics are positioned closely, such as foggy conditions or nighttime scenes. This clustering validates the effectiveness of using CLIP embeddings for adapter selection.

### B.2. BLIP for Labeling LoRA Adapters

To highlight the transparency and interpretability of our system, we leverage the connection between the CLIP embedding space and natural language. By processing the centroids of our LoRA adapters with BLIP-2 [33], we obtain natural language captions describing the domain encapsulated in each adapter’s training set. This provides semantic insights into the content and characteristics of the training set used for each adapter.

Table 4 reports, for each dataset, the caption provided by BLIP, together with the answers to two simple questions “Where is this?” and “Describe the environment in two words”, when processing domain centroids. We can appreciate how both captions and answers are strongly related to the content in each dataset, even though retrieved information remains limited and coarse. Nevertheless, the possibility of extracting natural language captions of the LoRA centroids is an interesting feature, further motivating the use of CLIP as our domain navigator.

### B.3. Adapter Contributions

Figure 8 shows the adapter weight distribution for all datasets composing our benchmark involved in the leave-one-out experiments. The parameters used in this experiment are  $\tau = 0.01$  and top- $K = 7$ . The weights represent the relative contribution of each adapter to the fused model, highlighting their respective roles in the overall composition.

These pie charts provide insights into how different adapters contribute to the final model for specific target domains, demonstrating the effective combination of knowledge from relevant adapters.

### B.4. LoRA Support Score Analysis

We analyze the relationship between the LoRA support score and mIoU performance. We define LoRA support score for a test image  $\mathbf{x}_t$ :

$$\text{Support Score}(\mathbf{x}_t) = \sum_{i \in \mathcal{K}} \frac{w_i}{\|\mathbf{e}_t - \mathbf{c}_i\|_2}, \quad (15)$$

where  $w_i$  is the weight assigned to adapter  $i$ ,  $\mathbf{e}_t$  is the CLIP embedding of the test image,  $\mathbf{c}_i$  is the centroid of adapter  $i$ , and  $\mathcal{K}$  is the set of top- $K$  selected adapters.

We compute the support score for a sample of images and plot it against their corresponding mIoU scores. Figure 9 shows that images with higher support scores tend to have higher mIoU, confirming that the LoRA support score is a good predictor of segmentation performance. We also notice that for low values of support

Dataset	Caption:	Question: Where is this? Answer:	Question: Describe The environment in two words? Answer:
bdd	the view from the driver's seat of a car on a street in san francisco, ca, june 2018	the city of los angeles, california	The environment in two words is the environment in which it is located.
idd	road in kolkata, india, photo by person	a street in bangalore, india	city, road, traffic jam
nyu	a view of the kitchen in the house i'm renting in san francisco, ca, in summer 2008	the house i grew up in, in san francisco, california, usa	blue and white
acdc-rain	the rain is coming down hard, but the streets are dry, and the cars are moving along the road	berlin, germany, street view, rain	rain
acdc-fog	a view from the driver's seat of a car on a highway on a foggy morning in kiev, ukraine	the highway in bordeaux, france, on a foggy day in october 2018	foggy, rainy, cloudy, misty
muses-snow-day	the view from the driver's seat of a car on a city street with buildings in sight	the city of berlin, germany, on a rainy winter day	Rainy day in vienna, austria, with trees and buildings
coconutL	person	a small town in the middle of nowhere, nyc, usa	The environment is where the person lives, works, and plays.
acdc-night	street at night in kiev, ukraine, with traffic lights and a car on the road	austria	dark and light, city, traffic
a150	the blue house	a house in the middle of the woods	The environment is where the person lives, works, or plays.
Cityscapes	street view of berlin, germany	berlin, germany, in the year 2014	city, street, road, traffic light
pc59	person	the house of the person in the picture	blue sky, green grass
muses-fog-day	a view from the driver's seat of a car on a rainy day in bordeaux, france	berlin, germany, in the year 2014	city, road, street
muses-clear-day	the car driving on the street in the city	a foggy day in bordeaux, france, driving on the autoroute	foggy, rainy, misty
acdc-snow	street in krasnodar, russia, april 2019	berlin, germany, in the year 2040, a virtual reality simulation	city, road, street
muses-fog-night	the road at night, with car lights visible	austria	snow, winter
muses-clear-night	the car on the road at night, with city lights in the background	the road in the dark, in the middle of nowhere, at night	dark and light
mv	street view of kuala lumpur, malaysia, with the city's main road visible	austria	night, city, traffic, street lights
muses-snow-night	a view of the city from a car's windshield at night, with city lights and snow visible	australia	urban, city, cityscape
		a city in the uk, in wintertime, with a car driving on the road	snow, rain, night, city

Table 4. Text generation results using BLIP-2 [33]. For the image embedding, the average embedding across all images from each dataset was computed. Then different prompts were given to the model, as presented at the top of the table.

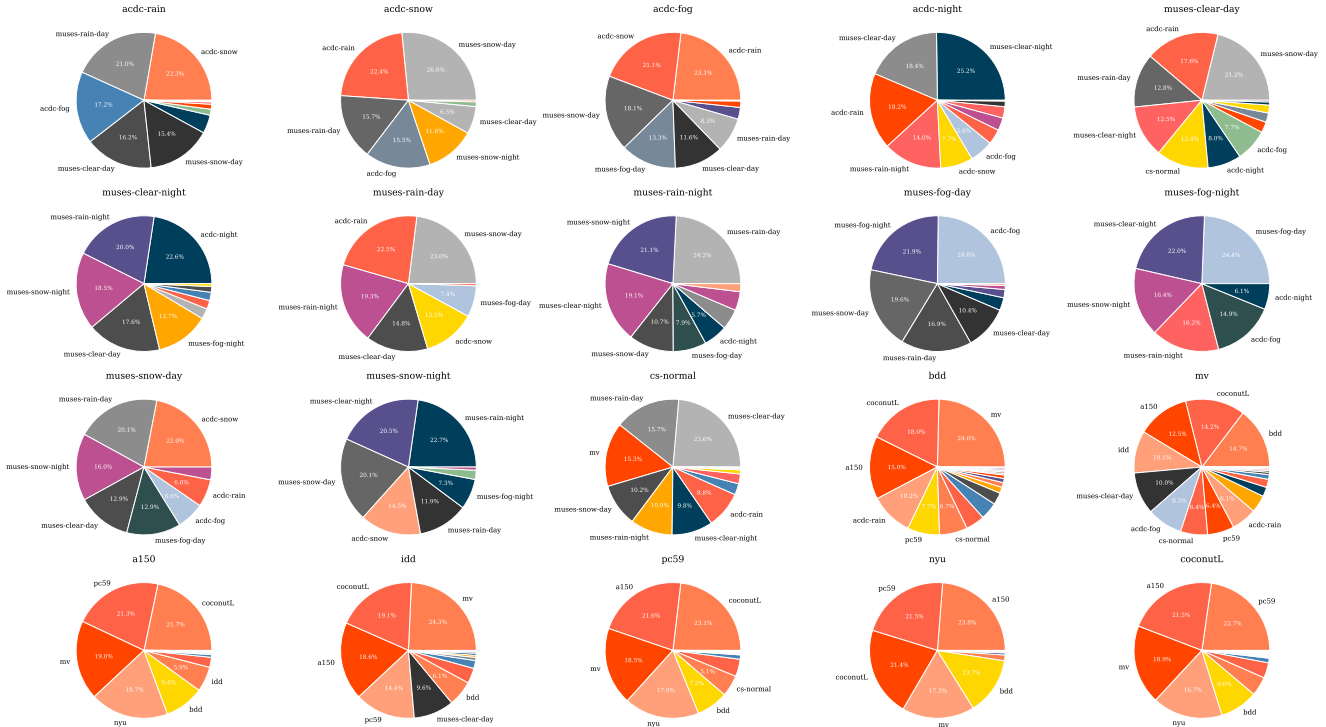


Figure 8. Adapters weight distribution for each benchmark dataset. Each pie chart is divided into sections proportional to the average contribution provided by each adapter based on CAT-Seg leave-one-out adaptation settings.

score (*e.g.* below 0.09), an improvement in support score does not strictly imply a stronger improvement in mIoU, showing that the underlying relation is likely not linear.

Overall, this analysis validates our assumption that proximity in the CLIP embedding space, combined with the weighting mechanism, is an effective heuristic for adapter selection.

## C. Ablations and Analysis

### C.1. Hyper-parameters Study

We conduct ablations over  $\tau$  and  $K$ , the two hyper-parameters controlling our system at test time.

- **Number of Adapters ( $K$ ):** Increasing  $K$  includes more adapters in the fusion, potentially providing more context but risking the introduction of unrelated knowledge while increasing the LoRA merging computational overhead.

- **Temperature ( $\tau$ ):** Regulates the weighting of adapters based on their distances. Lower  $\tau$  emphasizes closer adapters; higher  $\tau$  promotes a more uniform weighting.

Figure 10 shows a heatmap of overall performance across different values of  $K$  and  $\tau$ . Performance peaks at  $K = 7$  and  $\tau = 0.01$ , balancing relevance and diversity in adapter selection.

### C.2. Distance Metrics Comparison

We compare Euclidean distance (used in SemLA) against alternative distance measures, specifically cosine similarity and Mahalanobis distance. As shown in Table 5, cosine similarity – which would be a natural choice for CLIP embeddings – yields aligned performance with Euclidean distance, which is expected given CLIP embeddings exhibit almost uniform norms, making cosine similarity essentially a monotonic function of Euclidean distance. Conversely, Mahalanobis distance performs worse since covari-

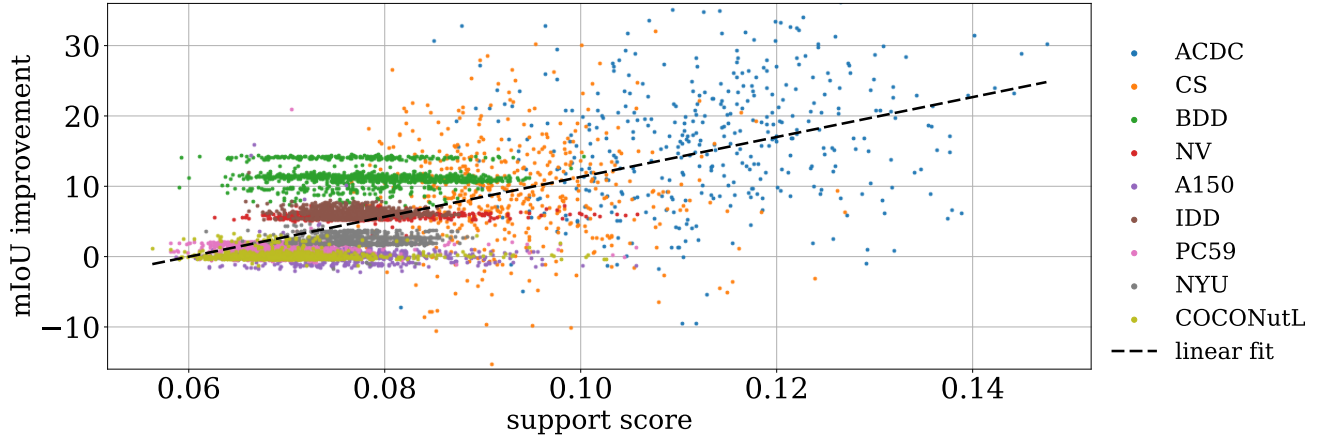


Figure 9. **Correlation between LoRA support score and mIoU.** Higher support scores correlate with better segmentation performance.

Method	ACDC				MUSES								CS	BDD	MV	A150	IDD	PC59	NYU	COCONutL*	h-mean
	rain	snow	fog	night	clear (d)	clear (n)	rain (d)	rain (n)	fog (d)	fog (n)	snow (d)	snow (n)									
Uniform [35]	67.40	66.35	69.71	49.98	58.28	55.78	54.70	45.09	73.75	45.16	61.02	49.08	62.18	58.19	31.51	37.25	38.83	63.06	48.93	(67.62)	51.89
SemLA with <b>Euclidean</b> (ours)	67.71	68.95	71.92	51.73	61.09	60.06	57.60	47.35	72.97	52.38	67.28	55.92	63.91	57.30	31.12	38.18	40.16	64.75	51.35	(67.26)	54.16
SemLA with <b>Cosine</b>	67.76	68.67	72.52	51.24	61.55	60.22	57.76	47.03	73.03	48.10	66.82	56.67	63.53	57.52	30.24	38.10	39.80	64.65	50.93	(67.31)	53.70
SemLA with <b>Mahalanobis</b> †	59.94	63.18	67.70	45.90	56.26	50.54	48.96	36.71	75.74	34.93	56.84	35.89	57.30	56.54	30.03	37.85	39.78	64.43	50.55	(67.69)	47.87

Table 5. **Ablation study – Distance Metrics Comparison (CAT-Seg [10]).** Comparing SemLA with alternative distances. On MUSES, (d) and (n) stand for day and night. ( ) means excluded from h-mean. † For Mahalanobis, source domains where the covariance cannot be computed are excluded from the library. The parameters for Cosine, Mahalanobis, and Late Fusions ( $\tau$  and  $K$ ) are tuned independently to achieve the best results with each variant.

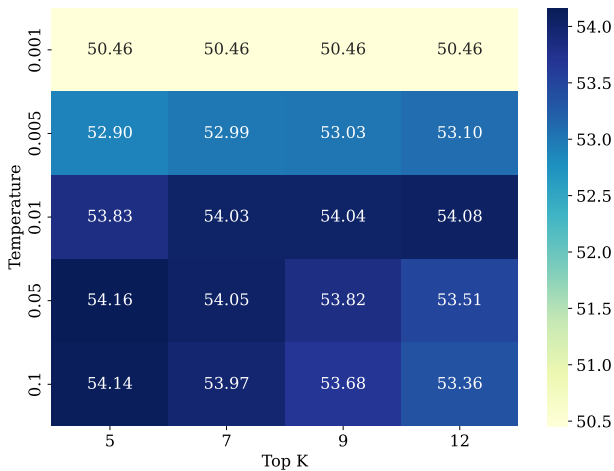


Figure 10. **Hyper-parameters Study.** Impact of  $K$  (number of adapters) and  $\tau$  (temperature) on overall performance (mIoU).

ance estimation becomes numerically unstable for domains with limited samples (fewer than 500 samples), necessitating the exclusion of some adapters and thus degrading performance. Overall, Euclidean distance emerges as the simplest, most robust choice for our method.

### C.3. Full Fine-Tuning (FFT)

We explore whether our library could be constructed using fully fine-tuned models instead of LoRA adapters. Table 6 reports the

results achieved either by deploying and fusing fully fine-tuned models or LoRA adapters in our library. While aggregating fully fine-tuned models is a known practice to merge different knowledge – as explored in [35] – the results indicate no benefits over our LoRA-based approach. Moreover, storing and merging full models is significantly more computationally expensive than operating with adapters, introducing a sizable overhead at inference time. Full fine-tuning is more prone to overfitting, especially on smaller datasets, whereas LoRA adapters are lightweight and can be trained effectively with limited data. This reinforces our choice of using LoRA adapters, which are modular, efficient, and easily combinable.

### C.4. Domain Navigators: DINO vs. CLIP

SemLA uses CLIP [51] to navigate into the LoRA Library and pick the most relevant adapters to combine. However, different visual encoders could serve the same purpose. In Table 7, we test the use of an alternative domain navigator – DINO v2 [48] – and compare the performance achieved by SemLA variants using this latter or CLIP.

On average, the two perform comparably, with CLIP embeddings slightly outperforming DINO ones in guiding adapter selection on average, likely due to their joint text-image embedding space capturing semantic information more effectively. Nonetheless, this experiment proves that SemLA is not bound to use CLIP as the domain navigator, although this latter provides nice properties in terms of explainability – as showcased in Section B.2.



Method	ACDC				MUSES								CS	BDD	MV	A150	IDD	PC59	NYU	COCONutL*	h-mean
	rain	snow	fog	night	clear (d)	clear (n)	rain (d)	rain (n)	fog (d)	fog (n)	snow (d)	snow (n)									
Zero-shot [10]	46.53	48.04	47.09	37.93	44.43	39.29	38.95	27.78	53.73	25.35	43.56	33.29	47.11	47.95	25.69	37.83	35.39	63.33	49.38	(68.26)	39.39
Oracles (LoRA)	70.94	69.22	69.98	51.55	69.36	57.09	54.28	52.11	75.85	61.26	66.25	54.35	67.47	60.06	49.57	53.99	64.34	68.68	61.90	(70.44)	61.05
Oracles (FFT)	70.97	72.02	73.78	53.09	70.49	58.20	55.57	53.18	74.90	62.64	65.83	58.67	70.35	61.03	50.56	52.84	66.38	68.43	64.36	(68.36)	62.38
Uniform [35] (FFT)	69.01	67.91	73.28	51.71	61.44	57.28	54.99	43.13	74.14	36.91	57.81	52.99	62.29	58.05	30.62	36.34	39.73	62.60	48.09	(65.29)	51.43
SemLA (FFT)	69.54	72.07	73.20	52.87	62.78	59.49	57.44	45.59	74.24	53.22	64.54	56.75	65.52	58.28	28.44	31.26	41.33	62.01	46.02	(63.64)	52.79
SemLA (LoRA)	67.71	68.95	71.92	51.73	61.09	60.06	57.60	47.35	72.97	52.38	67.28	55.92	63.91	57.30	31.12	38.18	40.16	64.75	51.35	(67.26)	54.16

Table 6. **Ablation study – Full Fine-Tuning vs LoRA Adaptation (CAT-Seg [10]).** We use full fine-tuned models instead of LoRA adapters and measure the impact on performance over our 20-domain benchmark in leave-one-out setting. On MUSES, (d) and (n) stand for day and night. ( ) means excluded from h-mean. SemLA (LoRA) with  $\tau = 0.05$ , and top- $K = 5$ ; SemLA (FFT) with  $\tau = 0.01$ , and top- $K = 9$ .

Method	ACDC				MUSES								CS	BDD	MV	A150	IDD	PC59	NYU	COCONutL*	h-mean
	rain	snow	fog	night	clear (d)	clear (n)	rain (d)	rain (n)	fog (d)	fog (n)	snow (d)	snow (n)									
Zero-shot [10]	46.53	48.04	47.09	37.93	44.43	39.29	38.95	27.78	53.73	25.35	43.56	33.29	47.11	47.95	25.69	37.83	35.39	63.33	49.38	(68.26)	39.39
Oracles	70.94	69.22	69.98	51.55	69.36	57.09	54.28	52.11	75.85	61.26	66.25	54.35	67.47	60.06	49.57	53.99	64.34	68.68	61.90	(70.44)	61.05
Uniform [35]	67.40	66.35	69.71	49.98	58.28	55.78	54.70	45.09	73.75	45.16	61.02	49.08	62.18	58.19	31.51	37.25	38.83	63.06	48.93	(67.62)	51.89
SemLA (DINOv2)	68.40	68.26	73.57	51.18	61.94	59.58	56.06	48.43	73.81	52.60	67.42	56.33	64.04	58.23	31.02	37.45	40.12	64.40	50.52	(67.63)	54.14
SemLA (CLIP)	67.71	68.95	71.92	51.73	61.09	60.06	57.60	47.35	72.97	52.38	67.28	55.92	63.91	57.30	31.12	38.18	40.16	64.75	51.35	(67.26)	54.16

Table 7. **Ablation study – CLIP [51] vs DINOv2 [48] for domain navigation (CAT-Seg [10]).** – We generate the weights for merging the LoRAs based on features extracted from DINOv2 or CLIP, and evaluate the impact on performance over our 20-domain benchmark in a leave-one-out setting. On MUSES, (d) and (n) stand for day and night. ( ) means excluded from h-mean.

## D. Additional Qualitative Results

Figure 11 presents additional qualitative segmentation results comparing the zero-shot baseline, uniform merging, and SemLA across different domains. These examples further confirm the effectiveness of our method in adapting to diverse and challenging domains without any additional training being conducted.

## E. Discussion: Real-World Deployment

While SemLA demonstrates strong performance in controlled experimental settings, deploying it in real-world applications introduces additional challenges and considerations. In this section, we discuss practical aspects related to the use of CLIP as a domain navigator and propose strategies to address potential limitations.

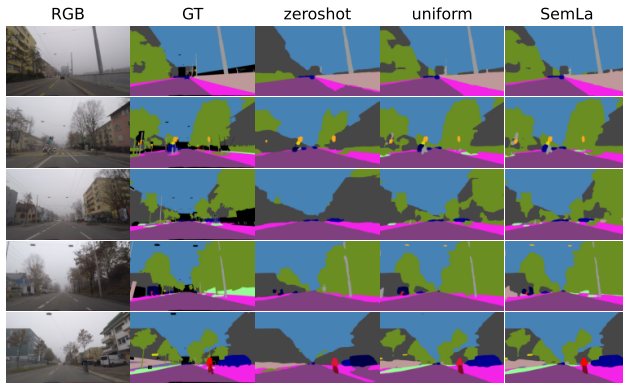
**CLIP as a Domain Navigator for Specific Domains.** Although CLIP has shown remarkable generalization capabilities across diverse domains – as evidenced by our extensive 20-domain benchmark – it may struggle in niche or highly specialized domains [67]. When bringing SemLA into production for such specific use cases, it is important to account for this potential limitation. If the target domain is well-scoped, using a fine-tuned domain navigator, with better semantic understanding, might provide better performance.

Alternatively, a hierarchical approach could be explored: a general CLIP model can provide a coarse understanding of the domain, identifying then a domain-specific CLIP expert. The expert is then tasked with computing the LoRA distances more precisely.

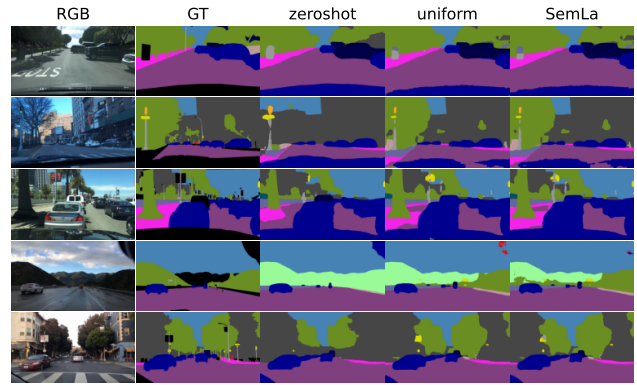
**Efficiency in Production Environments.** In production-intensive applications, dynamically loading and unloading dedicated LoRA adapters for each individual input image may be impractical due to computational overhead. While this overhead is significantly lower than the one introduced by retraining the model at test time – as required by most traditional test-time adaptation

methods – it is still non-negligible. For applications that do not require real-time processing, such as batch processing of large volumes of images (e.g., processing data accumulated over 24 hours), a practical approach involves pre-computing the CLIP embeddings for all images. The images can then be clustered based on their embeddings, and a batch centroid can guide the fusion of relevant LoRA adapters for the entire batch. This reduces the frequency of adapter loading and improves efficiency by applying the same fused model to similar images. In contrast, real-time applications in the field of robotics and autonomous driving cannot rely on batch processing due to their immediate response requirements. In these cases, we propose implementing a debouncing mechanism that triggers adapter swapping only when there is a significant change in the domain. Specifically, the system can monitor the CLIP embeddings of incoming images—or use an exponential moving average (EMA) of these embeddings—and compare them to the embeddings associated with the currently active adapters. If the embedding distance exceeds a predetermined threshold, indicating that a new domain has been encountered, the system triggers the retrieval and fusion of new adapters. This approach ensures that the model adapts only when necessary, minimizing computational overhead while maintaining adaptability. This strategy is analogous to concepts proposed in domain-adaptive systems like HAMLET [3], where adaptation occurs only upon detecting domain shifts. Furthermore, in a real-world deployment, the prediction process can be presented with an average LoRA distance metric. As shown in our analysis, this metric provides an additional source of confidence estimation by indicating how well the selected adapters align with the target domain. Such a heuristic contributes to the study of model calibration and can be valuable for downstream tasks—effectively informing whether to trust the model’s predictions in critical applications.

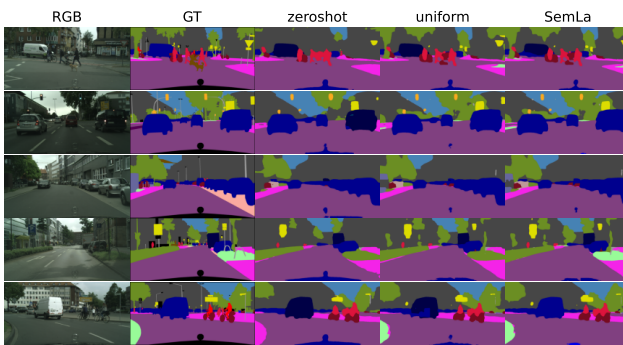
**Scalability and Model Calibration.** Scaling SemLA to handle a vast number of adapters introduces challenges in identifying and addressing library weaknesses. Automated strategies for recognizing gaps in the library – such as monitoring frequent



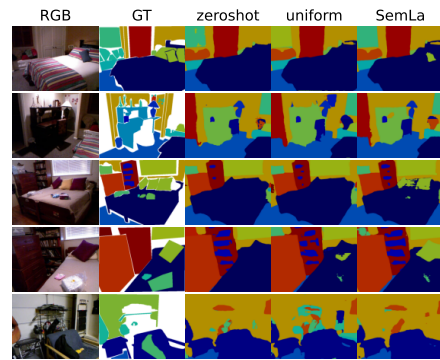
(a) ACDC Fog



(b) BDD



(c) Cityscapes (CS)



(d) NYU

Figure 11. **Additional qualitative results.** The datasets displayed are ACDC Fog, BDD, Cityscapes (CS), and NYU. For each dataset, images are shown in order: Input Image, Zero-Shot, Uniform Merging, SemLA (Ours), Ground Truth. Our method produces more accurate and detailed segmentations across various domains.

occurrences of high embedding distances – can prompt the training of new adapters to fill these gaps. Integrating a LoRA support score into the system allows for continuous monitoring of the model’s performance relative to the domain coverage of the adapter library. This not only enhances scalability but also improves the system’s robustness and reliability in dynamic environments.