

# DeCLIP: Decoupled Learning for Open-Vocabulary Dense Perception

## Supplementary Material

### Overview

This material provides supplementary details to the main paper, including the following sections:

- (A) **Details of Proxy Token Phenomenon**
- (B) **Additional Experiments**
  - (B.1) Ablation Studies
  - (B.2) Sanity Checks
  - (B.3) Further Details on Benchmark Results
- (C) **Additional Qualitative Analysis**
  - (C.1) Analyses of Feature Correlations
  - (C.2) Comparison of Semantic Segmentation Results
  - (C.3) Comparison of Attention Maps
- (D) **Details of Experimental Settings**
  - (D.1) Datasets and Evaluation Protocols
  - (D.2) Implementation Details
- (E) **Related Work**
  - (E.1) Open-Vocabulary Dense Prediction
  - (E.2) Transferring VLMs to Dense Prediction Tasks
  - (E.3) Vision Foundation Models

### A. Details of Proxy Token Phenomenon

This section primarily supplements the details of the proxy token phenomenon observed in CLIP, offering deeper insights into the rationale behind our proposed DeCLIP.

**Observation.** As stated in the main paper, ViT-based [14] CLIP utilizes the [CLS] token to represent the overall features of an image and performs image-text contrastive learning accordingly. Therefore, it is commonly believed that the [CLS] token comprehensively attends to all image tokens during the forward pass to obtain a “global view”, thereby enhancing the image classification process.

Unexpectedly, the [CLS] token ceased to focus on the primary object in the image starting from the 7th layer and instead redirected its attention to several image tokens in the background as shown in the first row of Figure 1. These specific image tokens continued to receive significant attention from the [CLS] token in the following encoding layers.

A similar pattern was observed in the attention maps of CLIP’s image tokens. As shown in the second row of Figure 1, we first randomly selected an image token located on the primary object in the image as the anchor image token, and then visualized its attention maps across different encoder layers. The experimental results show that the at-

tention of the anchor image token in layers 1-6 is primarily distributed over the object it belongs to. However, after the 7th layer, which is when the [CLS] token shifted its attention to several specific image tokens in the background, the anchor image token also began to focus on these specific image tokens.

Moreover, as illustrated in the third row of Figure 1, when the position of the anchor image token is shifted, the new anchor image token continues to exhibit high attention towards these specific tokens. This demonstrates that this phenomenon is not limited to a particular image token but is instead widespread across the image tokens in CLIP.

**Analysis.** One possible explanation for this phenomenon could be the redundancy present in image data. Images inherently carry a higher information load than text, encompassing substantial background details that are unrelated to image classification tasks. These specific background tokens may serve as “proxies” for the [CLS] token. This suggests that these tokens aggregate essential information from other image tokens, enabling the [CLS] token to form an approximate “global view” by summarizing content from them, thereby facilitating image classification. This perspective is also supported by recent studies [11, 48].

In over a decade of CNN [22, 36] development, no studies have reported similar phenomena. Therefore, we speculate that the second reason for this phenomenon may stem from the ViT architecture [14]. The classic ResNet [22] architecture consists of four stages, in which the feature resolution is halved and the number of channels is doubled at each stage. This is a process of learning sparse features, where redundant image details are progressively discarded, and feature semantics are continually enhanced. However, CLIP with a ViT architecture lacks this process. After patch embedding, the size and the number of channels in the feature map remain unchanged. As a result, the model spontaneously generates “proxy” tokens to mimic the process of learning sparse features, akin to CNN.

**Effects.** As discussed above, the proxy token phenomenon allows ViT CLIP to learn sparse features, which facilitate the extraction of key information from images, enhance image-text contrastive learning and reduce the optimization burden.

However, this phenomenon causes the image tokens in CLIP to indiscriminately focus on the proxy tokens in the background, rather than on the regions that are spatially or semantically related to them. Consequently, this leads to CLIP’s dense features to lack local discriminability and spatial consistency, affecting its performance in open-

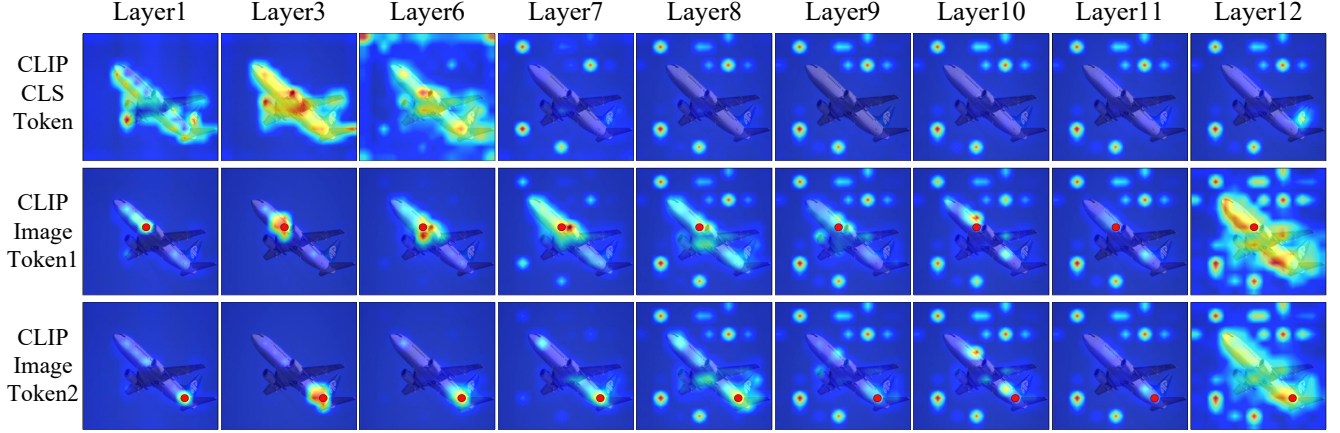


Figure 1. **Visualization of the “proxy” token phenomenon in the attention maps of the CLIP visual encoder.** Specifically, the input image resolution is 224\*224. We extract the attention weights from each attention block of CLIP and average them across the multi-head dimension (after Softmax), yielding attention maps  $\mathbf{M} \in \mathbb{R}^{197 \times 197}$ .  $\mathbf{M}[0, 1:] \in \mathbb{R}^{1 \times 196}$  represents the attention map from the [CLS] token to other image tokens (first row).  $\mathbf{M}[1:197, 1:197] \in \mathbb{R}^{196 \times 196}$  represents the attention map between each image token and all image tokens. We randomly select specific image tokens’ attention map (the second and third rows, indicated by the red dots) for visualization, each with dimensions of 1\*196. We reshape them to 1\*14\*14 and apply bilinear upsampling to 1\*224\*224 for better visualization.

Table 1. Ablation study on types of  $\mathbf{X}_{\text{Context}}$ .

$\mathbf{X}_{\text{Context}}$	Region Classification (mAcc)		Semantic Segmentation (mIoU)	
	COCO (Thing)	COCO (Stuff)	PASCAL Context59	ADE
<b>Q</b>	77.2	52.5	38.7	21.8
<b>K</b>	76.5	51.0	<b>39.4</b>	21.6
<b>Q + K</b>	<b>77.3</b>	<b>53.8</b>	39.2	<b>21.9</b>

Table 2. Ablation study on number of fine-tuning layers.

Fine-tuning Layers	Region Classification (mAcc)		Semantic Segmentation (mIoU)	
	COCO (Thing)	COCO (Stuff)	PASCAL Context59	ADE
3	62.7	47.0	38.0	21.8
6	67.1	47.8	<b>39.0</b>	<b>22.3</b>
9	70.7	50.5	<b>39.0</b>	22.1
12	<b>72.2</b>	<b>51.3</b>	38.7	21.8

vocabulary dense prediction tasks.

## B. Additional Experiments

### B.1. Ablation Studies

In this section, we conduct a thorough ablation study on DeCLIP, encompassing the examination of various  $\mathbf{X}_{\text{context}}$  implementations, the variation in the number of fine-tuning layers, the impact of the hyperparameter  $\lambda$  in the loss function, and the influence of the distillation baseline.

Except for the region classification experiment in Table 1, which was conducted at a resolution of 1024\*1024, the region classification performance in all other experiments was assessed at a resolution of 560\*560. Additionally, the semantic segmentation performance of all ablation experiments was assessed at a resolution of 336\*336.

**Types of Context.** Since there are various implementations of  $\mathbf{X}_{\text{context}}$ , including **Q**, **K**, and **Q + K**, we performed an ablation study on their performance in dense prediction tasks, including region classification (mAcc) and semantic segmentation (mIoU), as shown in Table 1. Specifically, implementing  $\mathbf{X}_{\text{context}}$  based on **K** means that the last attention block of CLIP leverages **K** to compute the attention weight. Additionally, implementing  $\mathbf{X}_{\text{context}}$  based on **Q + K** in-

volves first computing the attention weights of **Q** and **K** separately, and then summing them. The experimental results indicate that the performance differences among the three implementations are minimal, while the **Q** and **K** exhibits slightly better performance in dense prediction tasks.

**Number of fine-tuning layers.** We performed an ablation study to examine the relationship between the number of fine-tuning attention blocks and dense prediction performance. The experiment was conducted on the ViT-B version of CLIP, which comprises a total of 12 attention blocks. we experiment with updating the last 3, 6, 9, and 12 attention blocks. As shown in Table 2, we observed that as the number of fine-tuning layers increased, the performance of region classification continuously improved, reaching its peak at 12 layers. However, the performance of semantic segmentation peaked at 6 layers, and as the number of layers increased further, the performance slightly declined. In practice, to balance the performance of both tasks, we chose to fine-tune all attention blocks in the implementation of DeCLIP.

**Sensitivity Analysis of  $\lambda$ .** In DeCLIP, we employ a hyperparameter  $\lambda$  to balance the weight between  $\mathcal{L}_{\text{content}}$  and  $\mathcal{L}_{\text{context}}$ . We performed an ablation study to examine the

Table 3. Ablation Study on EVA-CLIP for open-vocabulary semantic segmentation

Method	Backbone	Training Set	ADE847	Context459	ADE150	Context59	VOC20	VOC21
CAT-Seg+CLIP [42]	ViT-B/16	COCO-Stuff	12.0	19.0	31.8	57.5	94.6	77.3
CAT-Seg+CLIP [42]	ViT-L/14	COCO-Stuff	16.0	23.8	37.9	63.3	97.0	82.5
CAT-Seg+EVA-CLIP [49]	ViT-B/16	COCO-Stuff	11.9	17.6	30.4	52.3	94.2	74.2
CAT-Seg+EVA-CLIP [49]	ViT-L/14	COCO-Stuff	14.2	21.3	34.8	56.2	95.8	80.1
CAT-Seg+DeCLIP	ViT-B/16	COCO-Stuff	15.3	21.4	36.3	60.6	96.6	81.3
CAT-Seg+DeCLIP	ViT-L/14	COCO-Stuff	<b>17.6</b>	<b>25.9</b>	<b>40.7</b>	<b>63.9</b>	<b>97.7</b>	<b>83.9</b>

Table 4. Ablation Study on EVA-CLIP for open-vocabulary semantic segmentation based on VLM features.

Method	With a background category			Without background category					Avg
	VOC21	Context60	COCO-Obj	VOC20	CityScape	Context59	ADE	COCO-Stf	
CLIP [42]	18.8	9.9	8.1	49.4	6.5	11.1	3.1	5.7	14.1
EVA-CLIP [49]	23.4	12.8	15.3	55.9	12.8	13.9	7.7	9.7	18.9
ClearCLIP [32]	51.8	32.6	33.0	80.9	30.0	35.9	16.7	23.9	38.1
EVA-ClearCLIP	47.0	29.7	30.2	78.3	26.3	29.4	16.7	20.4	34.7
DeCLIP	<b>59.7</b>	<b>35.3</b>	<b>36.4</b>	<b>85.0</b>	<b>32.8</b>	<b>39.2</b>	<b>21.9</b>	<b>25.3</b>	<b>41.9</b>

Table 5. Sentitivity Analysis of hyperparameter  $\lambda$ .

$\lambda$	Region Classification (mAcc)		Semantic Segmentation (mIoU)	
	COCO (Thing)	COCO (Stuff)	PASCAL Context59	ADE
0.1	72.4	50.6	37.9	21.3
0.2	72.4	51.0	38.4	21.7
0.25	72.2	51.3	38.7	21.8
0.3	71.9	51.4	38.7	21.7

relationship between the hyperparameter  $\lambda$  and dense prediction performance. The experimental results demonstrate that our method exhibits strong robustness, and the dense prediction performance of DeCLIP does not fluctuate drastically with changes in  $\lambda$ . Furthermore, the results indicate that  $\lambda = 0.25$  strikes a good balance between region classification capability and image segmentation performance.

**Distillation Baseline.** In our experiments, we used EVA-CLIP [49] as the baseline for DeCLIP, as we found that it demonstrated improved performance after distillation, as shown in Table 6. This can be attributed to two main factors: (1) EVA-CLIP uses the EVA02 [17] model for initializing the visual encoder. EVA02 was trained using Masked Image Modeling (MIM), thereby enhancing its compatibility with Vision Foundation Models (VFM). (2) EVA-CLIP’s [CLS] token exhibits superior zero-shot classification capability compared to OpenAI’s model [55]. In Sec. B.2, we conducted comprehensive sanity checks to verify whether the performance improvement of DeCLIP in dense prediction tasks is due to the use of EVA-CLIP.

## B.2. Sanity Checks

To eliminate potential biases that EVA-CLIP [49] might introduce, we conducted additional sanity check experiments.

Table 6. Comparison of different distillation baselines.

Source	Region Classification (mAcc)		Semantic Segmentation (mIoU)	
	COCO (Thing)	COCO (Stuff)	PASCAL Context59	ADE
OpenAI	65.0	38.8	36.2	18.6
EVA-CLIP	72.2	51.3	38.7	21.8

Specifically, we first apply vanilla EVA-CLIP as the backbone network in the CAT-Seg [9] model and compare its performance with DeCLIP in the Open-Vocabulary Semantic segmentation (OVSS) task, as shown in Table 3. Furthermore, we re-implemented ClearCLIP [32] based on EVA-CLIP and named it EVA-ClearCLIP. Then, we compared the performance between EVA-CLIP, EVA-ClearCLIP, and DeCLIP in the OVSS based on VLM features task, as shown in Table 4. We did not conduct further open-vocabulary detection experiments because the baseline detectors, OV-DQUO [52] and F-ViT [55], have already used EVA-CLIP as the backbone network in their respective studies.

**OVSS.** As shown in Table 3, experimental results demonstrate that directly applying EVA-CLIP to CAT-Seg performs worse than OpenAI’s model. In contrast, DeCLIP significantly improves CAT-Seg’s performance across all semantic segmentation benchmarks.

**OVSS based on VLM feautures.** As shown in Table 4, experimental results indicate that EVA-CLIP performs slightly better than CLIP in this task, while EVA-ClearCLIP underperforms in comparison to ClearCLIP. However, both EVA-CLIP and EVA-ClearCLIP fall significantly short of DeCLIP’s average performance of 41.9 across the eight benchmarks.



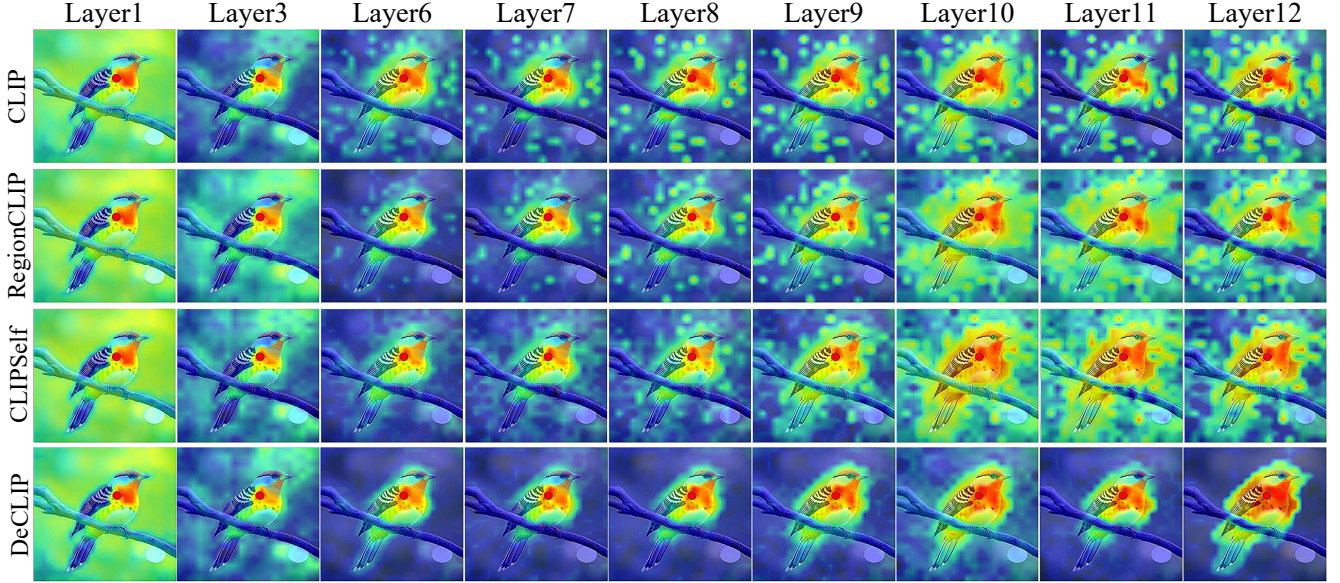


Figure 2. **Qualitative comparison of feature correlations between DeCLIP and existing pre-fine-tuning approaches** [55, 68]. Specifically, the input image resolution is 336\*336. We extract the output features from each attention block of CLIP, where each feature  $\mathbf{F} \in \mathbb{R}^{441 \times D}$ . Then, we compute the feature correlations  $\mathbf{FC} \in \mathbb{R}^{441 \times 441}$  between the image tokens within  $\mathbf{F}$  using cosine similarity. We randomly select a specific image token’s feature correlation (indicated by the red dots) and upsample it to a resolution of 336\*336 for visualization.

Based on the results of the aforementioned experiments, we conclude that the performance improvement of DeCLIP is not attributable to the introduction of EVA-CLIP, but is instead due to the superiority of the decoupled feature enhancement strategy.

### B.3. Further Details on Benchmark Results

We present detailed results for the OV-COCO, OV-LVIS, and cross-dataset benchmarks to provide a comprehensive comparison of the open-vocabulary object detection task, as shown in Tables 7 and 8.

## C. Additional Qualitative Analysis

This section further presents a qualitative experimental analysis of our proposed DeCLIP method in comparison to existing methods, including feature correlation analysis, semantic segmentation results, and attention map comparisons, thereby providing a more comprehensive demonstration of the superiority of DeCLIP’s decoupled feature enhancement strategy.

### C.1. Analyses of Feature Correlations

We have analyzed CLIP and found that its limitation in open-vocabulary dense prediction arises from image tokens failing to aggregate information from spatially or semantically related regions. Figure 2 presents a comparison of feature correlations among CLIP [42], DeCLIP, and exist-

ing pre-finetuning methods [55, 68] at each vision encoder layer.

This experiment provide insight into how the output features of each layer in CLIP’s visual encoder changed after fine-tuning. In this experiment, we randomly select an image token from the primary object within the image (i.e., the bird) as the anchor and visualize the cosine similarity between the anchor and the other image tokens. The experimental results indicate that the impact of various fine-tuning methods on the correlation of CLIP’s output features becomes noticeable starting from the 6th encoder layer.

**CLIP vs. existing pre-fine-tuning methods.** Rows 1, 2, and 3 of Figure 2 exhibit the changes in feature correlations of CLIP after region-level fine-tuning [55, 68]. The experimental results indicate that region-level fine-tuning enhances the feature correlations of the anchor image token to start converging towards the object it belongs to (rows 2 and 3), rather than being randomly scattered across the image (row 1).

This change is highly effective for open-vocabulary object detection tasks. As relevant features become more focused, region features exhibit enhanced discriminative power in the visual-language space when extracting the object’s region features from the image for recognition. However, these methods remain constrained in image segmentation tasks that demand pixel-level precision. As shown in the feature correlation results in rows 2 and 3 of Figure 2, most of the pixels surrounding the bird will be misclassified



Table 7. Detailed comparison on OV-COCO and OV-LVIS benchmarks. Caption supervision indicates that the method learns from extra image-text pairs, while CLIP supervision refers to transferring knowledge from CLIP. <sup>†</sup>: Detection Transformer based detectors.

Method	Supervision	Backbone	AP <sub>50</sub> <sup>Novel</sup>	AP <sub>50</sub> <sup>Base</sup>	AP <sub>50</sub>
ViLD [19]	CLIP	RN50	27.6	59.5	51.2
Detic [72]	Caption	RN50	27.8	51.1	45.0
OV-DETR <sup>†</sup> [64]	CLIP	RN50	29.4	61.0	52.7
ProxyDet [23]	Caption	RN50	30.4	52.6	46.8
RegionCLIP [68]	Caption	RN50	31.4	57.1	50.4
RTGen [6]	Caption	RN50	33.6	51.7	46.9
BARON-KD [54]	CLIP	RN50	34.0	60.4	53.5
CLIM [56]	CLIP	RN50	36.9	-	-
SAS-Det [67]	CLIP	RN50	37.4	58.5	53.0
RegionCLIP [68]	Captions	RN50x4	39.3	61.6	55.7
CORA <sup>†</sup> [57]	CLIP	RN50x4	41.7	44.5	43.8
OV-DQUO <sup>†</sup> [52]	CLIP	RN50x4	45.6	-	-
RO-ViT [28]	CLIP	ViT-L/16	33.0	-	47.7
CFM-ViT [27]	CLIP	ViT-L/16	34.1	-	46.0
F-ViT [55]	CLIP	ViT-B/16	37.6	54.9	50.4
BIND [66]	CLIP	ViT-L/16	41.5	58.3	54.8
F-ViT [55]	CLIP	ViT-L/14	44.3	64.1	59.0
F-ViT+DeCLIP	CLIP	ViT-B/16	41.1	57.8	53.5
F-ViT+DeCLIP	CLIP	ViT-L/14	46.2	65.2	60.3
OV-DQUO+DeCLIP <sup>†</sup>	CLIP	ViT-B/16	46.1	56.3	53.6
OV-DQUO+DeCLIP <sup>†</sup>	CLIP	ViT-L/14	<b>48.3</b>	60.0	56.9

Method	Supervision	Backbone	mAP <sub>r</sub>	mAP <sub>c</sub>	mAP <sub>f</sub>	mAP
ViLD [19]	CLIP	RN50	16.6	24.6	30.3	25.5
OV-DETR <sup>†</sup> [64]	CLIP	RN50	17.4	25.0	32.5	26.6
BARON-KD [54]	CLIP	RN50	22.6	27.6	29.8	27.6
RegionCLIP [68]	Caption	RN50x4	22.0	32.1	36.9	32.3
CORA <sup>†</sup> [57]	Caption	RN50x4	28.1	-	-	-
SAS-Det [67]	CLIP	RN50x4	29.1	32.4	36.8	33.5
CLIM [56]	CLIP	RN50x64	32.3	-	-	-
F-VLM [31]	CLIP	RN50x64	32.8	-	-	34.9
F-ViT [55]	CLIP	ViT-B/16	25.3	21.8	29.1	25.2
RTGen [6]	Caption	Swin-B	30.2	39.9	41.3	38.8
BIND [66]	CLIP	ViT-L/16	32.5	33.4	35.3	33.2
Detic [72]	Caption	Swin-B	33.8	-	-	47.0
CFM-ViT [27]	CLIP	ViT-L/14	33.9	-	-	36.6
RO-ViT [28]	CLIP	ViT-H/16	34.1	-	-	35.1
F-ViT [55]	CLIP	ViT-L/14	34.9	34.6	35.6	35.1
ProxyDet [23]	Caption	Swin-B	36.7	-	-	41.5
CoDet [38]	Caption	ViT-L/14	37.0	46.3	46.3	44.7
OV-DQUO <sup>†</sup> [52]	CLIP	ViT-L/14	39.3	-	-	-
F-ViT+DeCLIP	CLIP	ViT-B/16	26.8	22.4	29.8	26.0
F-ViT+DeCLIP	CLIP	ViT-L/14	37.2	35.2	36.5	36.0
OV-DQUO+DeCLIP <sup>†</sup>	CLIP	ViT-B/16	31.0	-	-	27.7
OV-DQUO+DeCLIP <sup>†</sup>	CLIP	ViT-L/14	<b>41.5</b>	-	-	34.6

Table 8. Detailed comparison of transferring LVIS-trained detectors to the COCO and Objects365 datasets.

Method	COCO [35]			Objects365 [47]						
	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>s</sub>	AP <sub>m</sub>	AP <sub>l</sub>	
Supervised Baseline [19]	46.5	67.6	50.9	25.6	38.6	28.0	-	-	-	-
ViLD [19]	36.6	55.6	39.6	11.8	18.0	12.6	-	-	-	-
DetPro [15]	34.9	53.8	37.4	12.1	18.8	12.9	4.5	11.5	18.6	-
BARON [54]	36.2	55.7	39.1	13.6	21.0	14.5	5.0	13.1	20.7	-
F-VLM [31]	37.9	59.6	41.2	16.2	25.3	17.5	-	-	-	-
CoDet [38]	39.1	57.0	42.3	14.2	20.5	15.3	-	-	-	-
RO-ViT [29]	-	-	-	17.7	27.4	19.1	-	-	-	-
CLIPSelf [55]	40.5	63.8	44.3	19.5	31.3	20.7	9.7	23.2	35.5	-
<b>DeCLIP</b>	<b>41.0</b>	<b>64.6</b>	<b>44.8</b>	<b>20.0</b>	<b>32.2</b>	<b>21.2</b>	<b>10.0</b>	<b>24.4</b>	<b>36.7</b>	-

as “bird” rather than to be “background”.

**CLIP vs. DeCLIP.** Rows 1 and 4 of Figure 2 exhibit the changes in feature correlations of CLIP after decoupled feature enhancement strategy. The experimental results indicate that DeCLIP enhances the feature correlations of the anchor image token to closely align with the object it represents, in clear contrast with other existing pre-fine-tuning approaches (row 2 and 3). This experiment reveals why DeCLIP is better suited for image segmentation tasks than existing methods. Additionally, the experiment demonstrates DeCLIP’s also superiority over current pre-finetuning approaches in region classification tasks. As shown in the feature correlation map of DeCLIP’s 12th layer, the image regions corresponding to the same object as the anchor image token display a strong red color, indicating a very high feature correlation strength in these regions, thereby enhancing the discriminative power of region features within the visual-language space.

## C.2. Comparison of Semantic Segmentation Results

Figure 3 shows a qualitative comparison of MaskCLIP [70], SCLIP [50], ClearCLIP [32], and our proposed DeCLIP across the Context59 [39], COCO-Stuff [3], Cityscapes [10], and ADE20K [69] datasets. We observe that, compared to other methods, DeCLIP consistently produces higher-quality and more precise segmentation maps.

Specifically, benefiting from content feature distillation, which improves the discriminability of local features, DeCLIP successfully recognizes trees, people, and curbs in the images, as shown in columns 1, 5, and 6 of Figure 3, whereas other models fail. Furthermore, our observation indicates that the distillation of context features improves the spatial consistency of DeCLIP’s local features, leading to smoother and less noisy segmentation results compared to other models, as demonstrated in columns 2, 3, 4, and 7 of Figure 3. This demonstrates the superiority of our decoupled feature enhancement strategy.

## C.3. Comparison of Attention Maps

Figure 4 offers a detailed comparison of attention maps between CLIP and our proposed DeCLIP approach. As DeCLIP involves unsupervised fine-tuning, we conducted tests using diverse cross-domain image styles to thoroughly assess its generalization capability. Specifically, we utilized generative models [46] to generate test images in various styles such as ink painting, watercolor, sketch, animation, and oil painting, which are depicted on the left side of Figure 4. These cross-domain test images were not part of the fine-tuning dataset for DeCLIP (i.e., COCO2017 [35]).

In addition, we performed a detailed comparison of at-

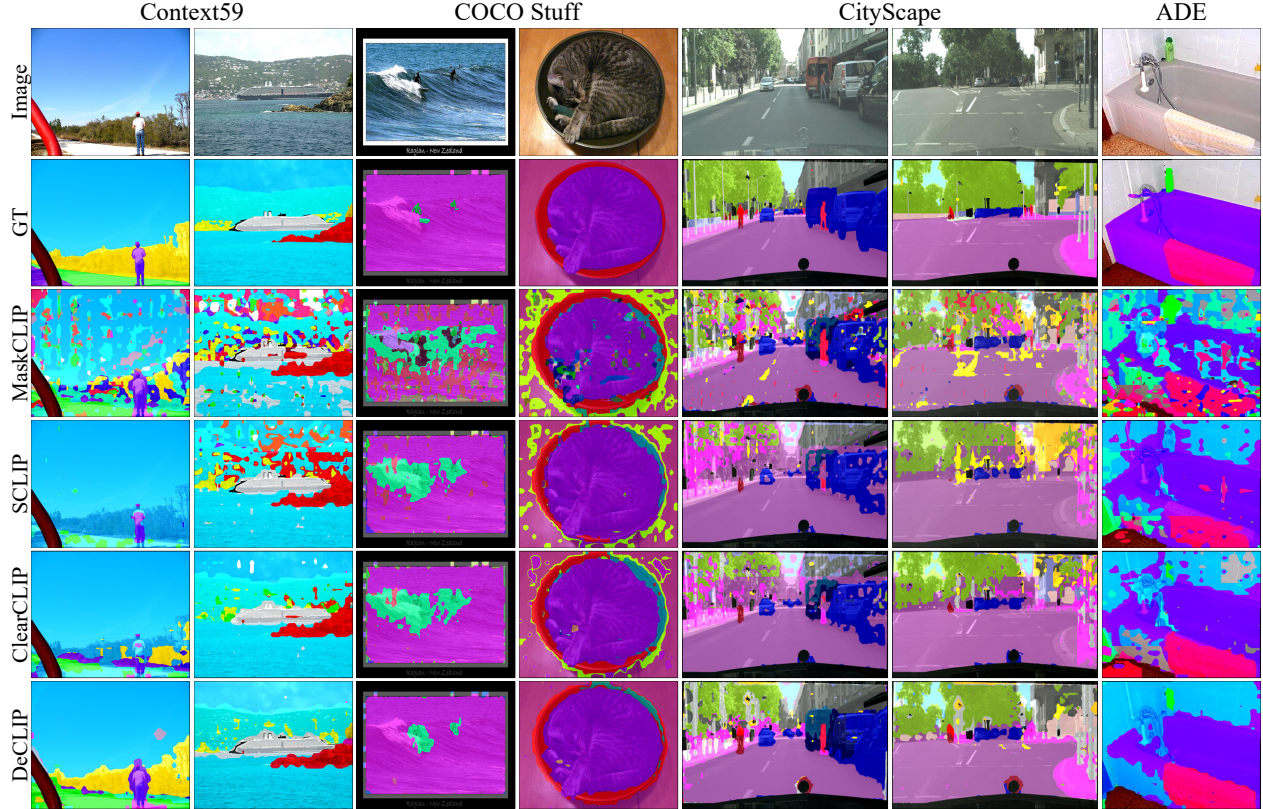


Figure 3. Qualitative comparison of the open-vocabulary semantic segmentation results between DeCLIP and existing approaches [32, 50, 70].

tention maps between CLIP and DeCLIP on in-domain images. Specifically, we selected a subset of images from the Object365 [47] validation set for testing, with the results shown on the right-hand side of Figure 4. During the testing phase, we first resized the images to 336×336 pixels and then fed them into the model to extract features. Subsequently, we randomly selected an anchor image token and visualized its attention map in the 12th attention block, as indicated by the red dots on the test images in Figure 4. For details on the calculation process of the attention map, please refer to Figure 1.

As depicted in Figure 4, due to the proxy token phenomenon, the heatmap generated by the anchor image token in vanilla CLIP frequently lacks semantic consistency with its corresponding object. In contrast, despite being fine-tuned only on the natural scene dataset COCO, DeCLIP demonstrates significant semantic relevance for both in-domain and cross-domain test images. Moreover, benefiting from context feature distillation, DeCLIP’s semantic correlations demonstrate remarkably fine granularity, effectively outlining the boundaries of each object semantically associated with the anchor image token.

## D. Details of Experimental Settings

In this section, we present further details and configurations utilized in our experiments.

### D.1. Datasets and Evaluation Protocols

**Open-Vocabulary Detection.** Following established settings [55, 57, 65], we evaluated our model on the OV-COCO [35], OV-LVIS [20], COCO, and Object365 [47] datasets. The OV-COCO dataset includes 48 base categories and 17 novel categories. The training set contains only base categories, totaling 107,761 images, while the validation set comprises 4,836 images featuring both base and novel categories. We report the mean Average Precision (mAP) at an Intersection over Union (IoU) threshold of 0.5 for novel categories. The OV-LVIS dataset consists of 1,203 categories. Its training set includes only 461 common and 405 frequent categories, totaling 100,170 images. The validation set contains 19,809 images with common, frequent, and rare categories. We report the mAP for rare categories at IoU thresholds ranging from 0.5 to 0.95. Additionally, we provide cross-dataset evaluation results on the COCO and Object365 validation sets for models trained on OV-LVIS to assess generalization across domains.



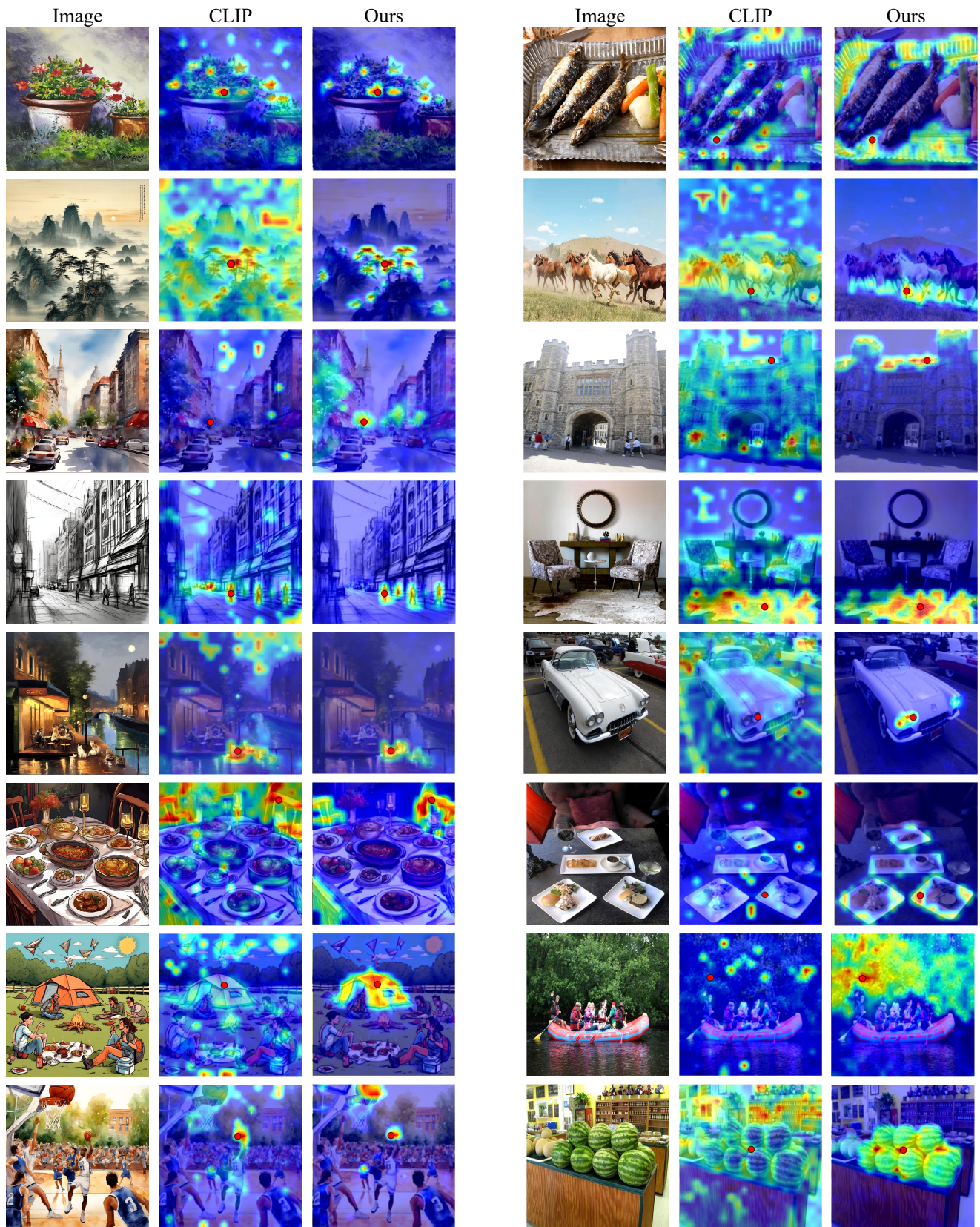


Figure 4. Comprehensive comparison of attention maps between CLIP and DeCLIP. The left side presents images of various styles generated by generative models [46]. The images presented on the right-hand side comes from a subset of images in the Object365 [47] validation set. Anchor image token marked in red.



**Open-Vocabulary Semantic Segmentation.** In line with prior studies [9], we trained our model on the COCO-Stuff dataset [3], which comprises 118,000 images with dense annotations across 171 categories. We then evaluated the model on the ADE20K [69], PASCAL VOC [16], and PASCAL-Context [39] datasets. ADE20K [69] includes 20,000 training images and 2,000 validation images, with two category sets: A-150 (150 common categories) and A-847 (847 categories) [13]. PASCAL-Context consists of 5,000 training and validation images, with category sets PC-59 (59 categories) and PC-459 (459 categories). The PASCAL VOC dataset includes 1,500 images for training and validation, featuring category sets PAS-20 (20 categories) and PAS-21 (20 object categories plus one background class). We used mean Intersection over Union (mIoU) as the evaluation metric in all experiments.

**Open-Vocabulary Semantic Segmentation Based on VLM Features.** To further evaluate DeCLIP, we assessed it on six commonly used semantic segmentation benchmarks: PASCAL VOC 2012 [16], PASCAL Context [39], Cityscapes [10], ADE20K [69], COCO-Stuff [35], and COCO-Object [3]. For datasets including a background category, we refer to them as VOC21 and Context60; those without a background category are termed VOC20 and Context59. Consistent with previous experiments, we used mIoU as the evaluation metric across these benchmarks.

## D.2. Implementation Details

**DeCLIP.** DeCLIP was trained on training set images from the COCO2017 [35] dataset using 8 GPUs, each with a batch size of 2, for 6 epochs (about 44 min/epoch on 8×4090 GPUs). The AdamW [37] optimizer with a learning rate of  $1e-5$  and a weight decay of 0.1 was employed during the training process.

During the content feature distillation process, the image is divided into  $k$  blocks, where  $k = m \times n$ , and  $m$  and  $n$  are randomly sampled from the range [1, 6]. After cropping  $k$  image blocks from the original image, the patches are resized to a resolution of 224×224 and subsequently fed into the teacher model to generate the corresponding [CLS] tokens for content feature distillation. Unless stated otherwise, our experiments were conducted using EVA-CLIP [49].

In the process of context feature distillation, given the distinct image preprocessing methods with varying means and standard deviations used by CLIP and VFM during pretraining, we incorporated the corresponding parameters during the distillation process. Additionally, to address the potential variation in patch sizes between CLIP and VFM (e.g., CLIP uses a 16-patch size while DINOv2 uses a 14-patch size), we adjusted the image resolutions to maintain consistency in the number of image tokens. For example, we set the resolution of CLIP to 1024 and that of DINOv2

to 896, ensuring both models possess 4096 image tokens. The weight  $\lambda$  for context feature distillation is established at 0.25. Unless specified otherwise, our default VFM is DINOv2 [41].

**Open-vocabulary detection.** In the open-vocabulary detection experiment, DeCLIP was evaluated in two model baselines: F-ViT [55] and OV-DQUO [52]. These baselines are constructed based on transfer learning principles, utilizing the image encoder of CLIP for feature extraction while maintaining the backbone network frozen during training and only training the task-specific components. The two baseline models utilize distinct detector architectures: F-ViT employs the traditional Faster R-CNN [45] architecture, whereas OV-DQUO utilizes the modern Detection Transformer [4] architecture. This enables a thorough assessment of the efficacy of our proposed approach.

We maintained the default training strategies and hyperparameter configurations from the original studies for both baseline models to uphold experiment fairness. The only modification was to the temperature parameter when integrating DeCLIP for object detection. For F-ViT, the temperature was set to 45 for the OV-COCO benchmark and 90 for the OV-LVIS benchmark. In OV-DQUO, the temperature was set to 50 for both the OV-COCO and OV-LVIS benchmarks.

**Open-Vocabulary Semantic Segmentation.** In the open-vocabulary semantic segmentation experiments, we applied DeCLIP to the CAT-Seg [9] baseline. For all experiments, we adhered to the default training and inference settings of vanilla CAT-Seg, replacing only the image encoder with DeCLIP.

**Open-Vocabulary Semantic Segmentation Based on VLM Features.** During inference, we resized the shorter side of images to 448 pixels and employed a sliding window strategy with a window size of 336×336 and a stride of 112×112. For all datasets, we generate textual descriptions by utilizing the standard ImageNet prompts [42] in conjunction with their respective class names. No post-processing steps were applied.

## E. Related Work

### E.1. Open-Vocabulary Dense Prediction

Open-vocabulary dense prediction aims to detect and segment visual concepts from novel categories using textual descriptions, extending beyond the base categories on which the model was trained. According to recent surveys [73], methods in this field can be broadly classified into four categories: knowledge distillation-based [21, 53, 54, 64], pseudo-labeling [52, 63, 67, 68, 72], region-aware training [18, 27, 29, 57, 59], and transfer learning-based approaches [12, 26, 31, 33, 34, 52, 55].

Knowledge distillation-based methods, such as ViLD [19], BARON [54], and OADP [53], propose various distillation frameworks to transfer the generalized classification knowledge of VLMs [42, 49] into dense prediction models. Pseudo-labeling methods like RegionCLIP [68] and SAS-Det [67] enhance region-text alignment by generating pseudo-labels for image-text pairs using VLMs or self-training techniques. Region-Aware Training methods, exemplified by CORA [57], improve the object classification accuracy of CLIP by learning region prompts.

Transfer Learning-Based methods [9, 12, 24–26, 34, 52, 55, 60, 61] utilize the image encoder of VLM as a feature extractor and exclusively train lightweight task-specific components. These methods have become mainstream in open-vocabulary dense prediction due to their broad applicability. While leveraging VLMs as feature extractors offers significant advantages due to their comprehensive pre-training, directly applying these image-level models to dense prediction tasks often results in domain shift issues [55, 57], thereby limiting their performance. In this paper, we integrate DeCLIP into transfer learning-based object detection baselines F-ViT and OV-DQUO, as well as the image segmentation baseline CATSeg, to enhance their performance in open-vocabulary dense prediction tasks.

## E.2. Transferring VLMs to Dense Prediction Tasks

As VLMs [42, 49] were initially trained on image-text pairs, the direct application of these image-level models to dense prediction tasks, which require region-level or pixel-level semantic understanding, results in significant performance degradation. Several studies have attempted to address this limitation through fine-tuning strategies. These approaches can be broadly categorized into joint fine-tuning and pre-fine-tuning approaches.

Joint fine-tuning methods fine-tune CLIP while training task-specific components [9, 24, 25, 33, 34, 58, 60]. For instance, CAT-Seg [9] proposes an attention fine-tuning strategy based on ViT CLIP, which generalizes well to unseen categories. MAFT [24] leverages attention bias to fine-tune CLIP for mask classification.

Pre-fine-tuning methods directly fine-tune CLIP using cost-efficient techniques [40, 55–57, 68]. For instance, CLIM [56] employs a mosaic augmentation technique to stitch multiple images into a single image, enabling each sub-image to serve as a pseudo-region for region-text contrastive learning. CLIPSelf [55] enhances CLIP’s region classification accuracy by maximizing cosine similarity between its region representations and the corresponding image crop representations.

Despite the promising results of the two categories of fine-tuned methods, they continue to exhibit certain limitations. In contrast to these studies, we conduct an analysis of CLIP and identify that its limitation in open-vocabulary

dense prediction stems from the inability of image tokens to effectively aggregate information from spatially or semantically related regions. To address this, we propose integrating VFMs into the pre-fine-tuning process and decoupling features for distillation, thereby improving the discriminability and spatial consistency of CLIP’s local features.

## E.3. Vision Foundation Models

Vision foundation models, including the Self-Supervised Representation Learning (SSL) series [1, 2, 5, 7, 41, 71] and the SAM series [30, 44], which are trained on large-scale segmentation data, demonstrate the ability to extract features that exhibit strong spatial consistency.

SSL is a key area in computer vision that focuses on learning meaningful visual features without manual annotations [1, 2, 5, 7, 41, 71]. Vision models trained through SSL can extract image features with excellent spatial understanding. For example, the DINO series [5, 41] can identify similar semantic regions across different images and segment main objects without explicit supervision. Another prominent vision foundation model is SAM [30, 44], which demonstrates similarly outstanding spatial understanding. Trained on the extensive SA-1B segmentation dataset, SAM can accurately capture and segment objects regions in images based on prompts.

Recently, some studies have explored the combination of CLIP with VFM, such as SAM-CLIP [51], OV-SAM [62], and FrozenSeg [8], with the goal of integrating SAM’s powerful image segmentation capabilities and CLIP’s zero-shot semantic perception capabilities. AM-RADIO [43] trains a unified vision model through multi-teacher distillation from multiple foundational vision models such as CLIP, DINOv2, and SAM. However, SAM-CLIP, OV-SAM, and FrozenSeg focus on integrating CLIP into SAM rather than enhancing CLIP itself as DeCLIP does. AM-RADIO does not support OVSS, as confirmed by its authors in Github issues (No. 81, 55, and 42). Another study that solves similar problems to DeCLIP is ViT-Register [11]. However, unlike DeCLIP, ViT-Register [11] does not solve the dense perception deficiency arising from CLIP’s image-text alignment.

## References

- [1] Benedikt Alkin, Lukas Miklautz, Sepp Hochreiter, and Johannes Brandstetter. Mim-refiner: A contrastive learning boost from intermediate pre-trained representations. *arXiv preprint arXiv:2402.10093*, 2024. 9
- [2] Xiang An, Kaicheng Yang, Xiangzi Dai, Ziyong Feng, and Jiankang Deng. Multi-label cluster discrimination for visual representation learning. In *European Conference on Computer Vision*, pages 428–444. Springer, 2025. 9
- [3] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *Proceedings*

- of the IEEE conference on computer vision and pattern recognition, pages 1209–1218, 2018. 5, 8
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 8
  - [5] 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. 9
  - [6] Fangyi Chen, Han Zhang, Zhantao Yang, Hao Chen, Kai Hu, and Marios Savvides. Rtgen: Generating region-text pairs for open-vocabulary object detection. *arXiv preprint arXiv:2405.19854*, 2024. 5
  - [7] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9640–9649, 2021. 9
  - [8] Xi Chen, Haosen Yang, Sheng Jin, Xiatian Zhu, and Hongxun Yao. Frozenseg: Harmonizing frozen foundation models for open-vocabulary segmentation. *arXiv preprint arXiv:2409.03525*, 2024. 9
  - [9] 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. 3, 8, 9
  - [10] 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. 5, 8
  - [11] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. *arXiv preprint arXiv:2309.16588*, 2023. 1, 9
  - [12] 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. 8, 9
  - [13] 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. 8
  - [14] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 1
  - [15] Yu Du, Fangyun Wei, Ziheng Zhang, Miaojing Shi, Yue Gao, and Guoqi Li. Learning to prompt for open-vocabulary object detection with vision-language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14084–14093, 2022. 5
  - [16] 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. 8
  - [17] Yuxin Fang, Quan Sun, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva-02: A visual representation for neon genesis. *Image and Vision Computing*, 149:105171, 2024. 3
  - [18] 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. 8
  - [19] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. *arXiv preprint arXiv:2104.13921*, 2021. 5, 9
  - [20] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5356–5364, 2019. 5, 6
  - [21] Kunyang Han, Yong Liu, Jun Hao Liew, Henghui Ding, Jiajun Liu, Yitong Wang, Yansong Tang, Yujiu Yang, Jiashi Feng, Yao Zhao, et al. Global knowledge calibration for fast open-vocabulary segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 797–807, 2023. 8
  - [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1
  - [23] Joonhyun Jeong, Geondo Park, Jayeon Yoo, Hyungsik Jung, and Heesu Kim. Proxydet: Synthesizing proxy novel classes via classwise mixup for open-vocabulary object detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2462–2470, 2024. 5
  - [24] Siyu Jiao, Yunchao Wei, Yaowei Wang, Yao Zhao, and Humphrey Shi. Learning mask-aware clip representations for zero-shot segmentation. *Advances in Neural Information Processing Systems*, 36:35631–35653, 2023. 9
  - [25] Siyu Jiao, Hongguang Zhu, Jiannan Huang, Yao Zhao, Yunchao Wei, and Humphrey Shi. Collaborative vision-text representation optimizing for open-vocabulary segmentation. In *European Conference on Computer Vision*, pages 399–416. Springer, 2025. 9
  - [26] Laurynas Karazija, Iro Laina, Andrea Vedaldi, and Christian Rupprecht. Diffusion models for zero-shot open-vocabulary segmentation. *arXiv preprint arXiv:2306.09316*, 2023. 8, 9
  - [27] Dahun Kim, Anelia Angelova, and Weicheng Kuo. Contrastive feature masking open-vocabulary vision transformer. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15556–15566, 2023. 5, 8
  - [28] Dahun Kim, Anelia Angelova, and Weicheng Kuo. Region-aware pretraining for open-vocabulary object detection with vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11144–11154, 2023. 5
  - [29] Dahun Kim, Anelia Angelova, and Weicheng Kuo. Region-aware pretraining for open-vocabulary object detection with



- vision transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11144–11154, 2023. 5, 8
- [30] 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. 9
- [31] Weicheng Kuo, Yin Cui, Xiuye Gu, AJ Piergiovanni, and Anelia Angelova. F-vm: Open-vocabulary object detection upon frozen vision and language models. *arXiv preprint arXiv:2209.15639*, 2022. 5, 8
- [32] Mengcheng Lan, Chaofeng Chen, Yiping Ke, Xinjiang Wang, Litong Feng, and Wayne Zhang. Clearclip: Decomposing clip representations for dense vision-language inference. *arXiv preprint arXiv:2407.12442*, 2024. 3, 5, 6
- [33] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. *arXiv preprint arXiv:2201.03546*, 2022. 8, 9
- [34] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yanan 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. 8, 9
- [35] 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, 6, 8
- [36] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1
- [37] I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 8
- [38] Chuofan Ma, Yi Jiang, Xin Wen, Zehuan Yuan, and Xiaojuan Qi. Codet: Co-occurrence guided region-word alignment for open-vocabulary object detection. *Advances in Neural Information Processing Systems*, 36, 2024. 5
- [39] 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, 8
- [40] 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. 9
- [41] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafranec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 8, 9
- [42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 3, 4, 8, 9
- [43] Mike Ranzinger, Greg Heinrich, Jan Kautz, and Pavlo Molchanov. Am-radio: Agglomerative vision foundation model reduce all domains into one. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12490–12500, 2024. 9
- [44] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 9
- [45] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. 8
- [46] 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. 5, 7
- [47] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8430–8439, 2019. 5, 6, 7
- [48] Tong Shao, Zhuotao Tian, Hang Zhao, and Jingyong Su. Explore the potential of clip for training-free open vocabulary semantic segmentation. In *European Conference on Computer Vision*, pages 139–156. Springer, 2025. 1
- [49] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023. 3, 8, 9
- [50] Feng Wang, Jieru Mei, and Alan Yuille. Sclip: Rethinking self-attention for dense vision-language inference. *arXiv preprint arXiv:2312.01597*, 2023. 5, 6
- [51] 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. 9
- [52] Junjie Wang, Bin Chen, Bin Kang, Yulin Li, YiChi Chen, Weizhi Xian, and Huifeng Chang. Ov-dquo: Open-vocabulary detr with denoising text query training and open-world unknown objects supervision. *arXiv preprint arXiv:2405.17913*, 2024. 3, 5, 8, 9
- [53] Luting Wang, Yi Liu, Penghui Du, Zihan Ding, Yue Liao, Qiaosong Qi, Bialong Chen, and Si Liu. Object-aware

- distillation pyramid for open-vocabulary object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11186–11196, 2023. 8, 9
- [54] Size Wu, Wenwei Zhang, Sheng Jin, Wentao Liu, and Chen Change Loy. Aligning bag of regions for open-vocabulary object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15254–15264, 2023. 5, 8, 9
- [55] Size Wu, Wenwei Zhang, Lumin Xu, Sheng Jin, Xiangtai Li, Wentao Liu, and Chen Change Loy. CLIPSelf: Vision transformer distills itself for open-vocabulary dense prediction. In *The Twelfth International Conference on Learning Representations*, 2024. 3, 4, 5, 6, 8, 9
- [56] Size Wu, Wenwei Zhang, Lumin Xu, Sheng Jin, Wentao Liu, and Chen Change Loy. Clim: Contrastive language-image mosaic for region representation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 6117–6125, 2024. 5, 9
- [57] Xiaoshi Wu, Feng Zhu, Rui Zhao, and Hongsheng Li. Cora: Adapting clip for open-vocabulary detection with region prompting and anchor pre-matching. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7031–7040, 2023. 5, 6, 8, 9
- [58] 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*, pages 3426–3436, 2024. 9
- [59] 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. 8
- [60] Xin Xu, Tianyi Xiong, Zheng Ding, and Zhuowen Tu. Masq-clip for open-vocabulary universal image segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 887–898, 2023. 9
- [61] Qihang Yu, Ju He, Xueqing Deng, Xiaohui Shen, and Liang-Chieh Chen. Convolutions die hard: Open-vocabulary segmentation with single frozen convolutional clip. *Advances in Neural Information Processing Systems*, 36, 2024. 9
- [62] Haobo Yuan, Xiangtai Li, Chong Zhou, Yining Li, Kai Chen, and Chen Change Loy. Open-vocabulary sam: Segment and recognize twenty-thousand classes interactively. In *European Conference on Computer Vision*, pages 419–437. Springer, 2025. 9
- [63] Nir Zabari and Yedid Hoshen. Open-vocabulary semantic segmentation using test-time distillation. In *European Conference on Computer Vision*, pages 56–72. Springer, 2022. 8
- [64] Yuhang Zang, Wei Li, Kaiyang Zhou, Chen Huang, and Chen Change Loy. Open-vocabulary detr with conditional matching. In *European Conference on Computer Vision*, pages 106–122. Springer, 2022. 5, 8
- [65] Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14393–14402, 2021. 6
- [66] Heng Zhang, Qiuyu Zhao, Linyu Zheng, Hao Zeng, Zhiwei Ge, Tianhao Li, and Sulong Xu. Exploring region-word alignment in built-in detector for open-vocabulary object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16975–16984, 2024. 5
- [67] Shiyu Zhao, Samuel Schuster, Long Zhao, Zhixing Zhang, Vijay Kumar B G, Yumin Suh, Manmohan Chandraker, and Dimitris N. Metaxas. Taming self-training for open-vocabulary object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13938–13947, 2024. 5, 8, 9
- [68] Yiwu Zhong, Jianwei Yang, Pengchuan Zhang, Chunyuan Li, Noel Codella, Liunian Harold Li, Luowei Zhou, Xiyang Dai, Lu Yuan, Yin Li, et al. Regionclip: Region-based language-image pretraining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16793–16803, 2022. 4, 5, 8, 9
- [69] 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. 5, 8
- [70] 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. 5, 6
- [71] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*, 2021. 9
- [72] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *European Conference on Computer Vision*, pages 350–368. Springer, 2022. 5, 8
- [73] Chaoyang Zhu and Long Chen. A survey on open-vocabulary detection and segmentation: Past, present, and future. *arXiv preprint arXiv:2307.09220*, 2023. 8