# Diffusion Model is Secretly a Training-free Open Vocabulary Semantic Segmenter

**Jinglong Wang**[1*], **Xiawei Li**[1*], **Jing Zhang**[1†], **Qingyuan Xu**[1],
**Qin Zhou**[1], **Qian Yu**[1], **Lu Sheng**[1], **Dong Xu**[2]

[1]Beihang University
[2]The University of Hong Kong
wjlzy@buaa.edu.cn, zy2121108@buaa.edu.cn, zhang_jing@buaa.edu.cn, ZY2121121@buaa.edu.cn
zhouqin2023@buaa.edu.cn, qianyu@buaa.edu.cn, lsheng@buaa.edu.cn ,dongxu@hku.hk

## Abstract

Recent research has explored the utilization of pre-trained text-image discriminative models, such as CLIP, to tackle the challenges associated with open-vocabulary semantic segmentation. However, it is worth noting that the alignment process based on contrastive learning employed by these models may unintentionally result in the loss of crucial localization information and object completeness, which are essential for achieving accurate semantic segmentation. More recently, there has been an emerging interest in extending the application of diffusion models beyond text-to-image generation tasks, particularly in the domain of semantic segmentation. These approaches utilize diffusion models either for generating annotated data or for extracting features to facilitate semantic segmentation. This typically involves training segmentation models by generating a considerable amount of synthetic data or incorporating additional mask annotations. To this end, we uncover the potential of generative text-to-image conditional diffusion models as highly efficient open-vocabulary semantic segmenters, and introduce a novel training-free approach named DiffSegmenter. Specifically, by feeding an input image and candidate classes into an off-the-shelf pre-trained conditional latent diffusion model, the cross-attention maps produced by the denoising U-Net are directly used as segmentation score maps, which are further refined and completed by the followed self-attention maps. Additionally, we carefully design effective textual prompts and a category filtering mechanism to further enhance the segmentation results. Extensive experiments on three benchmark datasets show that the proposed DiffSegmenter achieves impressive results for open-vocabulary semantic segmentation.

## Introduction

The process of gathering and annotating images with pixel-level labels is labor-intensive and time-consuming. Models trained solely on fully annotated data are restricted to specific categories, which limits their scalability. As a result, there is an increasing focus on developing open vocabulary semantic segmentation methods. These approaches are considered more practical alternatives to fully supervised approaches, as they reduce the reliance on extensive annotation and enable segmentation across a wider range of categories.

---

*These authors contributed equally.
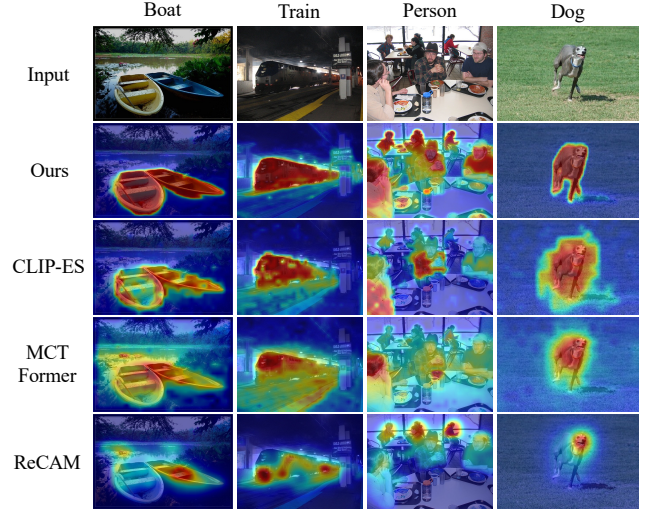†Corresponding author



Figure 1: Segmentation score maps generated by our proposed DiffSegmenter and previous discriminative methods.

Large-scale image-text pre-trained models like CLIP (Radford et al. 2021) have showcased remarkable zero-shot transferability across diverse downstream tasks. As a result, researchers have started to investigate the potential of leveraging these models for open-vocabulary semantic segmentation in zero-shot (Xu et al. 2022c; Ding et al. 2022; Liang et al. 2023) setting or weakly supervised (Xie et al. 2022; Lin et al. 2023a) setting. However, it is worth noting that the alignment process based on contrastive learning employed by the discriminative image-text pre-trained models may unintentionally result in the loss of crucial localization information and object completeness, which are essential for achieving accurate semantic segmentation.

Recently, there has been a growing interest in exploring the application of generative models, such as diffusion models, beyond image generation tasks, particularly in the field of semantic segmentation (Amit et al. 2021; Wu et al. 2022, 2023; Karazija et al. 2023a; Baranchuk et al. 2021; Xu et al. 2023b). While these methods show promise, they still require specific training or complex data synthesis procedures. For instance, SegDiff and MedSegDiff rely on annotated data for supervised segmentation. For open-vocabulary

segmentation tasks, ODISE and SegDiff require additional mask annotations from a large dataset for training. Diffu-Mask and OVDiff involve complex data synthesis process.

In this paper, we delve deeper into the capabilities of generative text-to-image diffusion models for semantic segmentation. In contrast to discriminative image-text pre-trained models that model the probability distribution $p(\boldsymbol{c}|\boldsymbol{x})$ of a whole image, diffusion models model the conditional distribution $p(\boldsymbol{x}|\boldsymbol{c})$ for generating complete objects while preserving the underlying semantic information. Specifically, Stable Diffusion(Rombach et al. 2022) introduces the cross-attention mechanism, effectively transforming diffusion models into conditional generative models. The resulting cross-attention map serves as a vital component in the conditional generation process, providing valuable insights into the conditional probability distribution $p(\boldsymbol{x}|\boldsymbol{c})$. By using Bayes' theorem with a proper prior $p(\boldsymbol{c})$, it becomes straightforward to convert this into pixel-level $p(\boldsymbol{c}|\boldsymbol{x})$, which is the key to semantic segmentation. Unfortunately, this idea has been under-exploited by previous work.

To this end, we uncover the potential of generative text-to-image diffusion models as highly efficient training-free semantic segmenters, and introduce a novel training-free approach named DiffSegmenter. The proposed DiffSegmenter enables open-vocabulary semantic segmentation using off-the-shelf diffusion models without the need for any additional learnable modules or any parameter tuning. Specifically, we feed an input image and candidate classes into an off-the-shelf pre-trained conditional latent diffusion model. By extracting the cross-attention maps produced by the denoising U-Net, we obtain the initial segmentation score maps. We further propose to refine and complete the semantic score maps with the self-attention maps of the U-Net, which effectively capture pairwise pixel affinities to link pixels with similar features. Additionally, we carefully design effective textual prompts and a category filtering mechanism by harnessing the power of the pre-trained BLIP(Li et al. 2022) model to further enhance the segmentation results.

Figure 1 shows a comparison of the segmentation score maps obtained from our proposed DiffSegmenter and previous discriminative model-based methods. This comparison highlights the improvements achieved by the proposed approach in terms of segmentation quality and accuracy. Extensive experiments on PASCAL VOC 2012(Everingham et al. 2010), MS COCO 2014(Lin et al. 2014) datasets, and Pascal Context(Mottaghi et al. 2014) verify that our generative approach to semantic segmentation, leveraging the inherent capabilities of diffusion models, achieves impressive results for open-vocabulary semantic segmentation in both zero-shot setting and image-level weakly supervised setting.

In summary, our contributions are three-fold:

- We uncover the potential of generative text-to-image conditional diffusion models as highly efficient training-free semantic segmenters.

- We propose a novel training-free open-vocabulary semantic segmentation method, DiffSegmenter, which fully exploits the potential of attention layers in the denoising U-Net of diffusion models and carefully designs

the textual prompts for semantic enhancement.

- Extensive experiments on three benchmark datasets show that the proposed DiffSegmenter achieves impressive results for open-vocabulary semantic segmentation.

## Related Work
### CLIP-based Open-vocabulary Segmentation

Large-scale image-text pre-trained models like CLIP (Radford et al. 2021) have shown promising zero-shot classification capabilities. It has recently been extended to dense prediction tasks such as semantic segmentation. CLIP models are discriminative models trained by aligning images with global category label, denoted as $p(\boldsymbol{c}|\boldsymbol{x})$, for an input image $\boldsymbol{x}$. However, semantic segmentation tasks necessitate per-pixel predictions denoted as $p(\boldsymbol{c}|\boldsymbol{x}_{h,w})$, where $\boldsymbol{c}$ is the class label for each individual pixel in the image $\boldsymbol{x}_{h,w}$.

Hence, existing methods tackle the adaptation of CLIP models for semantic segmentation through three primary approaches. 1) A large scale pixel-wise annotated data or pseudo labels are used to finetune the CLIP image encoder (Rao et al. 2022; Zhou, Loy, and Dai 2022; Zabari and Hoshen 2021). 2) Additional pre-trained mask generators (Liang et al. 2023; Ghiasi et al. 2022) super-pixel grouping methods (Xu et al. 2023a; Ding et al. 2022; Lüddecke and Ecker 2022) are required to obtain the object mask, while CLIP is employed solely as a proposal/group classifier. 3) Since discriminative models only activate the most discriminative parts of class objects, complex multi-round refinement mechanisms are employed to enhance the class activation maps obtained from CLIP models as segmentation results (Xie et al. 2022; Lin et al. 2023a).

However, it is worth noting that the alignment process based on contrastive learning employed by the discriminative image-text pre-trained models may unintentionally result in the loss of crucial localization information and object completeness, which are essential for achieving accurate semantic segmentation. In summary, adapting CLIP models as a semantic segmentor is complex.

### Diffusion Models for Perception Tasks

Inspired by non-equilibrium thermodynamics, diffusion models(Sohl-Dickstein et al. 2015) incorporate both diffusion and reverse processes to facilitate data generation. Over time, diffusion models have undergone significant developments and emerged as prominent generative models in contemporary research(Ho, Jain, and Abbeel 2020; Dhariwal and Nichol 2021; Rombach et al. 2022). The cross-attention mechanism that allows for conditional synthesis in diffusion model is also exploited for different tasks by previous work, such as image editing (Hertz et al. 2022), sketch synthesis (Xing et al. 2023), and interpreting the Stable Diffusion model (Tang et al. 2022).

Nevertheless, diffusion models possess not only powerful generative capabilities but also can be cleverly utilized in various perception tasks. The Diffusion Classifier demonstrated that density estimation derived from diffusion models, which generate images from text, can be employed for zero-shot classification tasks without requiring additional
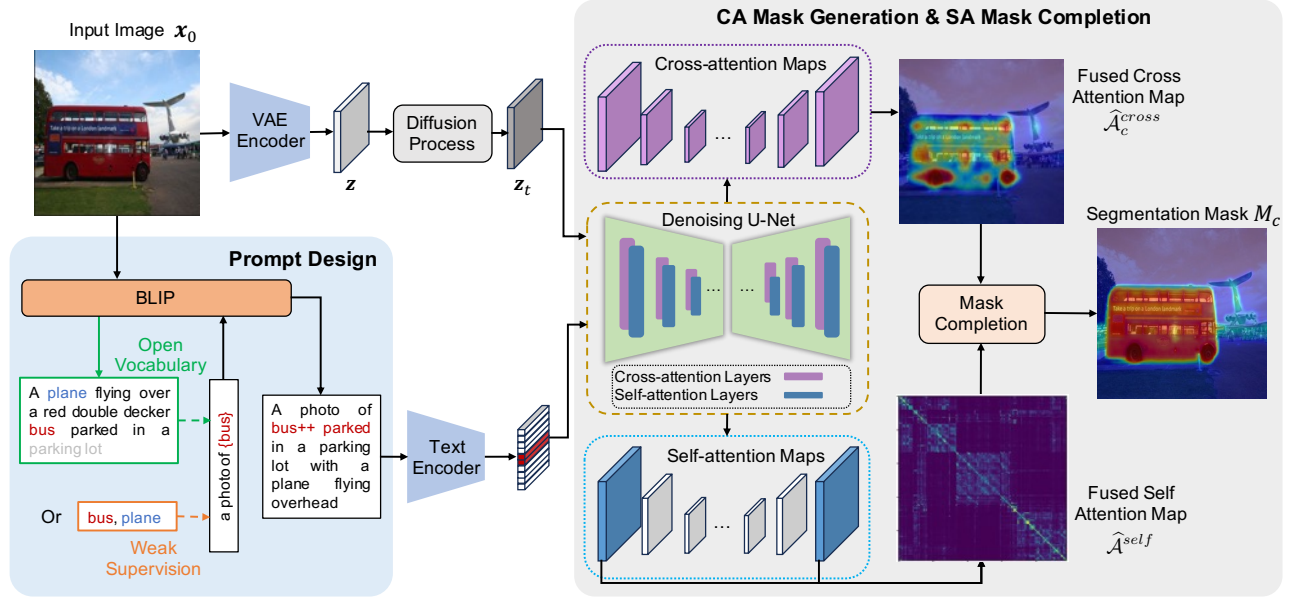
Figure 2: **Overview of the proposed DiffSegmenter.** An input image and enhanced candidate class tokens by the BLIP-based prompt design module are fed into an off-the-shelf pre-trained conditional latent diffusion model. The fused cross-attention maps produced by the denoising U-Net are treated as the initial segmentation score maps, which is further refined and completed by the fused self-attention maps of the U-Net. Note that the parameters of all the involved models are frozen without any tuning.

training(Li et al. 2023). Leveraging the energy-based modeling nature of diffusion models, they can be applied to unsupervised compositional concept discovery, facilitating effective representation of semantic information in images(Liu et al. 2023). TF-ICON can leverage off-the-shelf diffusion models to perform cross-domain image-guided composition without requiring additional training or finetuning(Lu, Liu, and Kong 2023).

Diffusion models have also exhibited impressive performance in semantic segmentation tasks through various approaches. Directly inputting the latent variables from diffusion models into segmentation networks enables semantic segmentation even with limited training samples(Baranchuk et al. 2021). DiffuMask leverages diffusion models to generate images and pixel-level image annotations, thereby training a high-performing semantic segmentation model(Wu et al. 2023). However, both of these approaches are not open-vocabulary semantic segmentation solutions. ODISE leverages large-scale text-to-image diffusion models and discriminative models to construct a panoramic segmentation scheme(Xu et al. 2023b). OVDiff initially utilizes diffusion models to generate prototypes for multiple classes and then matches pixel features with these prototypes during actual segmentation, thereby determining the class of each pixel(Karazija et al. 2023b). Although these two approaches yield excellent results, they still require model training and the use of a pretrained segmentation network for actual segmentation. In contrast to previous works, our approach does not need addtional training and does not rely on a pretrained segmentation network.

## Methodology

### Diffusion Model for Segmentation: A Baseline.

Unlike discriminative models that focus on modeling the posterior probability $p(\boldsymbol{c}|\boldsymbol{x})$ to predict per-pixel classification results, generative models such as conditional latent diffusion models models the conditional probabilities $p(\boldsymbol{x}|\boldsymbol{c})$. This reversal of the conditional probabilities allows generative models to generate complete objects given a specific category, which is valuable for image synthesis. A generative model that models $p(\boldsymbol{x}|\boldsymbol{c})$ can be easily adapted for semantic segmentation via Bayes' theorem.

The conditional diffusion models such as Stable Diffusion(Rombach et al. 2022) define the conditional probability of $\boldsymbol{x}_0$ as,

$$p_\theta(\boldsymbol{x}_0|\boldsymbol{c}) = \int_{\boldsymbol{x}_{1:T}} p(\boldsymbol{x}_T) \prod_{t=1}^{T} p_\theta(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t, \boldsymbol{c}) d\boldsymbol{x}_{1:T}, \quad (1)$$

where $p(\boldsymbol{x}_T)$ is fixed to $\mathcal{N}(0, I)$. Diffusion models are trained to minimize the variational lower bound (ELBO)(Blei, Kucukelbir, and McAuliffe 2017) of the log-likelihood,

$$\log p_\theta(\boldsymbol{x}_0|\boldsymbol{c}) \le \mathbb{E}_q \left[ \log \frac{p_\theta(\boldsymbol{x}_{0:T}, \boldsymbol{c})}{q(\boldsymbol{x}_{1:T}|\boldsymbol{x}_0)} \right]. \quad (2)$$

Inspired by diffusion classifier (Li et al. 2023), the ELBO can be written as,

$$-\mathbb{E}_{t,\epsilon} \left[ \|\epsilon - \epsilon_\theta(\boldsymbol{x}_t, \boldsymbol{c})\|^2 \right]. \quad (3)$$

As the loss for the ELBO defined in Eq. 3 is calculated on a per-pixel basis, we can leverage Bayes' theorem to obtain pixel-level classification results. The baseline method

demonstrates reasonable results by carefully set the diffusion time steps and utilizing a large number of samples for Monte Carlo estimation. However, it is worth noting that this approach requires a substantial computational cost due to the high number of samples required.

Furthermore, it is important to note that the diffusion process has the potential to compromise the structural information of the original image. And thus the fine-grained details and local pixel-level semantics may become less distinguishable or even blurred during the diffusion progresses, which is detrimental to the semantic segmentation results.

By delving deeper into the capabilities of conditional latent diffusion models and their ability to model conditional distributions $p(\boldsymbol{x}|\boldsymbol{c})$ through the utilization of cross-attention modules, this paper introduces a new attention-based method, named DiffSegmenter.

## Proposed DiffSegmenter.

**Motivation.** By introducing the conditional denoising autoencoder $\epsilon_\theta(\boldsymbol{x}_t, t, \boldsymbol{c})$, Stable Diffusion utilizes cross-attention modules to establish the relationship between text features and latent space features. This allows us to control the denoising process based on the text condition, thereby influencing the image synthesis. In other words, Stable Diffusion can model conditional distributions of the form $p(\boldsymbol{x}|\boldsymbol{c})$ thanks to the cross-attention modules.

In the case of segmentation tasks, if we assume a uniform prior over $\{\boldsymbol{c}_i\}$, the formulation of Bayes' theorem simplifies to $p(\boldsymbol{c}_i|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|\boldsymbol{c}_i)}{\sum_j p(\boldsymbol{x}|\boldsymbol{c}_j)}$. It can be concluded that $p(\boldsymbol{c}_i|\boldsymbol{x}) \propto p(\boldsymbol{x}|\boldsymbol{c}_i)$. Since Stable Diffusion can approximate the distribution $p(\boldsymbol{x}|\boldsymbol{c}_i)$ with ELBO, it can also be used for discriminative tasks. Building upon this insight, we utilize cross-attention from the diffusion model to establish the relationship between the distributions $p(\boldsymbol{c}_i|\boldsymbol{x})$ and $p(\boldsymbol{x}|\boldsymbol{c}_i)$.

The cross-attention map calculates the similarity between the query and key matrices (Q and K) and applies the weighted sum to the value matrix (V). This means that the spatial features are augmented with the most relevant textual features for each pixel. A higher similarity leads to larger activation values in $\mathcal{A}$, indicating a closer relationship between the current pixel and the corresponding text. Therefore, we simply use the cross-attention maps calculated by the pre-trained conditional latent diffusion models as the mask basis for semantic segmentation, without the need for any additional learnable modules or any parameter tuning.

**Overview.** An overview of our method is shown in Figure 2. We feed an input image and candidate classes into an off-the-shelf pre-trained conditional latent diffusion model. By extracting the cross-attention maps produced by the denoising U-Net, we obtain the initial segmentation score maps. We further propose to refine and complete the semantic score maps with the self-attention maps of the U-Net, which effectively capture pairwise pixel affinities to link pixels with similar features. Moreover, we carefully design effective textual prompts and category filtering mechanism based on BLIP model to further enhance the results.

## Cross-attention-based Score Map Generation

We feed an input image to be segmented to the VAE encoder $\mathcal{E}$ to obtain a latent variable $\boldsymbol{z} = \mathcal{E}(\boldsymbol{x}_0)$. In each time-step $t$, we add Gaussian noise to $\boldsymbol{z}$, resulting in $\boldsymbol{z}_t = \sqrt{\bar{\alpha}_t}\boldsymbol{z} + \sqrt{1 - \bar{\alpha}_t}\epsilon$, where $\epsilon \sim N(0, I)$. The deep spatial features are represented as $\varphi(\boldsymbol{z}_t) \in \mathbb{R}^{H \times W \times C}$, where $H$ is the height of the feature map, $W$ is the width, and $C$ is the number of channels. These features are linearly mapped to the query matrix $Q = \ell_Q(\varphi(\boldsymbol{z}_t))$. The text prompt $\mathcal{P}$ is encoded by the text encoder and mapped to $\tau_\theta(\mathcal{P}) \in \mathbb{R}^{N \times D}$, where $N$ is the length of the text tokens and $D$ is the latent projection dimension. The text is further linearly mapped to the key matrix $K = \ell_K(\tau_\theta(\mathcal{P}))$ and the value matrix $V = \ell_V(\tau_\theta(\mathcal{P}))$.

The cross-attention maps are computed as:

$$\mathcal{A}^{cross} = \text{Softmax}(\frac{QK^T}{\sqrt{d}}), \qquad (4)$$

where $\mathcal{A}^{cross} \in \mathbb{R}^{H \times W \times N(\text{reshaped})}$, and $\mathcal{A}_c^{cross} \in \mathbb{R}^{H \times W}$ represents the cross-attention for a specific class token $\boldsymbol{c}$. The output of the cross-attention $\widehat{\varphi}(\boldsymbol{z}_t) = \mathcal{A}^{cross}V$ is used to update the spatial features $\varphi(\boldsymbol{z}_t)$.

We observed that different cross-attention layers have complementary abilities to capture semantic information as shown in Figure 3. The attention maps from the smaller feature maps can precisely localize the objects-of-interest while the attention maps from the larger feature maps capture more fine-grained object details. Therefore, we fuse the cross-attention maps of class $\boldsymbol{c}$ from multiple layers with different importance weights $w_l$, which obtains:

$$\widehat{\mathcal{A}}_c^{cross} = \sum_{l \in L} w_l \cdot \mathcal{A}_{c,l}^{cross} \in \mathbb{R}^{H \times W}, \qquad (5)$$

where $\mathcal{A}_{c,l}^{cross}$ is the cross-attention map of class token $\boldsymbol{c}$ of the $l$-th layer in the UNet, and $\sum_{l \in L} w_l = 1$ is the importance weights. Attention maps from four different layers (e.g., $8 \times 8$, $16 \times 16$, $32 \times 32$, $64 \times 64$) of the UNet are used and interpolated to the same size for fusion in our method. The fused cross-attention map serves as the initial segmentation score maps in our method.
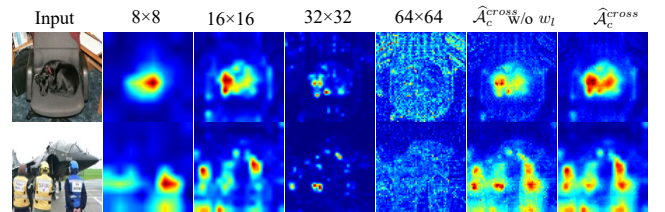


Figure 3: Cross-attention maps of different layers.

## Self-attention-based Score Map Completion

Although the score maps generated by cross-attention maps could successfully localize the objects of interest, it has been observed that the produced score maps often lack clear object boundaries and may exhibit internal holes. It is worth noting that there exists attention between latent representations in UNet, namely the self-attention map $\mathcal{A}_l^{self} \in$

$\mathbb{R}^{HW \times HW}$. The Q, K, and V (i.e., Query, Key, and Value) in self-attention are all based on latent visual features, enabling the establishment of correlations between different pixels. This endows self-attention with the ability to perform region completion, which can compensate for the incomplete activation regions in cross-attention. By considering different time steps and layers, we obtain the fused self-attention map,

$$\widehat{\mathcal{A}}^{self} = \frac{1}{L} \sum_{l \in L} \mathcal{A}_l^{self} \in \mathbb{R}^{HW \times HW}, \qquad (6)$$

In our method, instead of fusing self-attention maps from all the layers of the UNet, we specifically select and fuse the attention maps of the encoder and decoder with a latent feature size of $64 \times 64$, where more fine-grained pixel-level affinity weights are captured.

To address the issues of incomplete segmentation score maps generated by cross-attention maps, we introduce a mask refinement mechanism by multiplying the cross-attention maps with pixel affinity weights obtained from self-attention maps. The refined segmentation score maps for class $c$ is given by,

$$M_c = \text{norm}(\widehat{\mathcal{A}}^{self} \cdot vec(\widehat{\mathcal{A}}_c^{cross})), \qquad (7)$$

where $\text{norm}(\cdot)$ is min-max normalization to ensure the segmentation score maps are appropriately scaled, $vec(\widehat{\mathcal{A}}_c^{cross}) \in \mathbb{R}^{HW \times 1}$, and $vec(\cdot)$ is a vectorization operation of a matrix. Therefore, the final segmentation score maps $M_c$ are simply produced during the denoising inference process of the pre-trained text-conditional diffusion models, without any additional additional learning modules or any parameter tuning.

### Prompt Design for Semantic Enhancement

**Augmented Prompts.** The quality of the attention-based segmentation score maps is also highly influenced by the textual inputs. Intuitively, providing more comprehensive and detailed descriptions of the objects-of-interest can lead to better localization. Building on this intuition, we extend the use of cross-attention maps beyond just the class token. We also incorporate corresponding adverbs and adjectives that provide more detailed attributes of the objects. The cross-attention maps of the class names and the adverbs or adjectives are fused to obtain the segmentation score maps.

**Class Token Re-weighting** To place greater emphasis on the class token corresponding to the object-of-interest, we introduce a re-weighting mechanism for the class token embedding. By assigning a higher weight to the class token, we prioritize and highlight the target objects while suppress the background for better segmentation results. The re-weighted class token is denoted as "[class]++".

### Open-vocabulary Semantic Segmentation

In zero-shot open-vocabulary semantic segmentation, we generate $C$ candidate score maps by individually feeding each candidate class into the text encoder and compute the $M_c$ for the input image. The final segmentation results are obtained by taking the maximum value of each pixel by comparing $C$ score maps. To reduce the computational cost, we pre-filter the candidate categories using two strategies. Firstly, we feed the input image to the BLIP model, which generates captions describing the image content. From these captions, we extract the nouns, which represent the main objects or entities present in the image. These extracted nouns then serve as the filtered categories, narrowing down the candidate categories for semantic segmentation. Secondly, we feed the input image and all the candidate classes into CLIP model and select the classes with cosine similarity larger than 0.97. The union of the selected classes by BLIP and CLIP models is used as our final filtered categories for a given image.

In weakly-supervised open-vocabulary semantic segmentation, the final segmentation results of the training data are produced from the prompted ground-truth object labels of the input image. And any off-the-shelf segmentation models can be used for training the segmentation models.

## Experiments

### Datesets and Implementation Details

To validate our methodology, we conducted experiments on tasks involving open-vocabulary semantic segmentation in both weakly-supervised setting with image-level class supervisions and completely zero-shot setting.

For the task of weakly supervised semantic segmentation, we evaluated our approach on the PASCAL VOC 2012(Everingham et al. 2010). PASCAL VOC 2012 consists of 21 categories (including one background class). The augmented set comprises 10,582 images for training and 1449 images for testing. For zero-shot open-vocabulary semantic segmentation, we evaluate our approach on PASCAL VOC 2012(VOC), Pascal Context(Context), and COCO-Object(Object) datasets. Pascal Context(Mottaghi et al. 2014) have 60 classes(including backgroud class). Since our method is entirely train-free, we solely validate our approach on the validation set. The VOC, Context, and Object datasets comprise 1449, 5105, and 5000 images respectively.

We employed the BLIP(Li et al. 2022) model to generate textual descriptions for input images. For zero-shot open-vocabulary segmentation, we extracted nouns from the generated captions and compared them to all the candidate categories to compute similarity. We utilize the outcomes of (Lin et al. 2023b)CLIP classification to complement the extraction of nouns. This process yielded labels for the images, which were subsequently incorporated as prompts into the BLIP model to obtain comprehensive descriptions of the image. We validated our approach using Stable Diffusion v1.5 with frozen pre-trained parameters.

For the weakly supervised semantic segmentation task, we followed the approach (Lin et al. 2023a; Xie et al. 2022) to assess the quality of the initial Class Activation Maps (CAMs). After applying DenseCRF(Krähenbühl and Koltun 2011) to obtain pseudo-labels, we proceeded to train a fully supervised segmentation network, such as DeepLabV2. We employ the mean Intersection over Union (mIoU) as the evaluation metric for all experiments.

**Results on Open-vocabulary Segmentation**

| Method | VOC | Context | Object |
|---|---|---|---|
| *Training-involved* | | | |
| ReCo(Shin, Xie, and Albanie 2022) | 25.1 | 19.9 | 15.7 |
| ViL-Seg(Liu et al. 2022) | 37.3 | 18.9 | - |
| MaskCLIP(Zhou, Loy, and Dai 2022) | 38.8 | 23.6 | 20.6 |
| TCL(Cha, Mun, and Roh 2023) | 51.2 | 24.3 | 30.4 |
| CLIPpy(Ranasinghe et al. 2022) | 52.2 | - | 32.0 |
| GroupViT(Xu et al. 2022a) | 52.3 | 22.4 | - |
| ViewCo(Ren et al. 2023) | 52.4 | 23.0 | 23.5 |
| SegCLIP(Luo et al. 2023) | 52.6 | 24.7 | 26.5 |
| OVSegmentor(Xu et al. 2023a) | 53.8 | 20.4 | 25.1 |
| *Training-free* | | | |
| OVDiff(Karazija et al. 2023a) | **67.1** | **30.1** | 34.8 |
| DiffSegmenter (Ours) | 60.1 | 27.5 | **37.9** |

Table 1: Results of zero-shot open-vocabulary semantic segmentation on three benchmark datasets.

Table 1 shows a comparison of our method to previous work on zero-shot open-vocabulary semantic segmentation. We evaluate mIoU on three datasets: VOC, Context, and Object. Compared to previous training-based approaches, our method exhibits a significant improvement in mIoU across different datasets, showing the strong zero-shot generalization capability of our method. Our method only performs poorer than a concurrent work OVDiff on VOC and Context datasets. We argue that OVDiff necessitates a complex image synthesis process and involve additional pre-trained segmenters and feature extractors for prototype generation. By contrast, our method is simple and efficient.

Figure 4 illustrates the qualitative segmentation outcomes through a comparison of DiffSegmenter and the SegCLIP baseline. The results indicate that our approach achieves a more complete segmentation results of the objects of interest, with clearer mask boundaries. More results can be found in the supplementary material.

| Method | VOC train |
|---|---|
| *Image-level Supervsion* | |
| IRN(Ahn, Cho, and Kwak 2019) | 48.8 |
| SC-CAM(Chang et al. 2020) | 50.9 |
| SEAM(Wang et al. 2020) | 55.4 |
| AdvCAM(Lee, Kim, and Yoon 2021) | 55.6 |
| RIB(Lee et al. 2021) | 56.5 |
| OoD(Lee et al. 2022) | 59.1 |
| MCTfomer(Xu et al. 2022b) | 61.7 |
| DiffSegmenter (Ours) | **70.5** |
| *Image-level Supervision+Language Supervision* | |
| CLIMS(Xie et al. 2022) | 56.6 |
| CLIP-ES(Lin et al. 2023a) | 70.8 |

Table 2: Segmentation results of on PASCAL VOC 2012 train sets with image-level object labels.

**Results on Weakly-supervised Segmentation**

Table 2 shows the segmentation results of the training set of VOC with image-level class labels. Our method results in an mIoU of 70.5%, which significantly outperforms most



Figure 4: Qualitative results of DiffSegmenter and SegCLIP baseline for zero-shot open-vocabulary segmentation.

of the discriminative model-based baseline methods with image-level supervision. Our method is also comparable to the state-of-the-art CLIP-ES model, which is a very strong baseline. However, CLIP-ES rely on both image-level supervision and language supervision while our method only rely on image-level supervision. Moreover, CLIP-ES uses a complex synonym fusion strategy to enrich the category names, which can also be used together with our method to further improve the results (we will leave this in our future work). In Figure 1, we showcase segmentation score maps for some example images, illustrating that our method generates more complete semantic score maps directly from the initial CAMs than baseline methods.

To further evaluate the quality of the generated segmentation masks, we utilized the segmentation masks in the training set as pseudo-labels for training a segmentation network. By following the state-of-the art methods, such as CLIMS and CLIP-ES, we employed the DeepLabV2 architecture pretrained on ImageNet-1k, with a ResNet-101 backbone, as our segmentation network. The results are presented in Table 3, where the segmentation model trained with the masks generated by our approach achieved 69.1% and 68.6% on the validation and test sets, which is comparable to the state-of-the-art methods.

**Ablation Studies and Analyses**

We present the results of ablation study to our method in Table 4 to evaluate the effectiveness of each component.

**Effect of different mask generation strategies.** Table 4 presents the mIoU evaluation results obtained on the PASCAL VOC 2012 training set when using different mask generation strategies proposed in our method. The visualized comparison results are shown in Figure 5. When using only

| Method | Backbone | Val | Test |
|---|---|---|---|
| *Image-level Supervsion* | | | |
| AdvCAM(Lee, Kim, and Yoon 2021) | R101 | 68.1 | 68.0 |
| RIB(Lee et al. 2021) | R101 | 68.3 | **69.1** |
| ReCAM(Chen et al. 2022) | R101 | 68.5 | 68.4 |
| DiffSegmenter (Ours) | R101 | **69.1** | **68.6** |
| *Image-level Supervision+Language Supervision* | | | |
| CLIMS(Xie et al. 2022) | R101 | 69.3 | 68.7 |
| CLIP-ES(Lin et al. 2023a) | R101 | 71.1 | 71.4 |

Table 3: Segmentation results on PASCAL VOC 2012 validation and test sets.

| Method | | | | | VOC train |
|---|---|---|---|---|---|
| $\widehat{\mathcal{A}}_c^{cross}$ | $\widehat{\mathcal{A}}_{all}^{self}$ | $\widehat{\mathcal{A}}^{self}$ | BLIP | "++" | mIoU |
| w/o $w_l$ | | | ✓ | ✓ | 61.25 |
| w/o $w_l$ | ✓ | | ✓ | ✓ | 65.01 |
| w/o $w_l$ | | ✓ | ✓ | ✓ | 67.89 |
| ✓ | | ✓ | | | 65.32 |
| ✓ | | ✓ | ✓ | | 67.99 |
| ✓ | | ✓ | | ✓ | 69.46 |
| ✓ | | ✓ | ✓ | ✓ | **70.49** |

Table 4: Ablations study.

cross-attention and averaging the values of the four cross-attention layers (i.e., $\widehat{\mathcal{A}}_c^{cross}$ w/o $w_l$), we can obtain an mIoU of 61.25%, which is already comparable to the state-of-the-art methods with image-level supervision. The incorporation of self-attention in all layers of the model (i.e., $\widehat{\mathcal{A}}_{all}^{self}$), the obtained $M_c$ with $\widehat{\mathcal{A}}_{all}^{self}$ leads to a significant improvement in performance. This verifies that self-attention provides rich image information that complements the limitations of cross-attention and enhances the activation maps. Furthermore, focusing solely on the highest resolution self-attention $\widehat{\mathcal{A}}^{self}$, the obtained $M_c$ leads to additional advancements, resulting in more complete image boundaries.
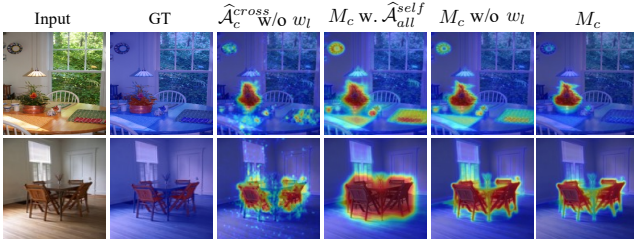


Figure 5: Comparison of different mask generation variants.

As shown in Figure 3, we observe that different cross-attention layers contain varying degrees of semantic information. For example, the strongest semantic information is captured by the cross-attention at resolutions of $8 \times 8$ and $16 \times 16$. Instead of the previous averaging operation, we applied weighted summation to the cross-attention from different layers. We empirically used weights of $[0.3, 0.5, 0.1, 0.1]$ for the cross-attention at resolutions of $8\times8, 16\times16, 32\times32$, and $64 \times 64$, respectively. By reinforcing the semantic infor-

mation, the final results are further improved compared to the averaging operation.

**Evaluations of the Prompt Design** Table 4 also provides the related ablation study for prompt design while fixing the attention map variant. Introducing BLIP allows us to generate conditional captions, providing additional textual descriptions for target categories. Attention maps corresponding to modified adjectives or adverbs that describe the target category can help in segmenting the target object. By changing the re-weighting the target category words, we can make the attention map focus more on the target category while suppressing the activation of the background or other objects. This strategy greatly enhances the segmentation results. In Figure 6, we can have consistent observations. "BLIP" refers to using BLIP to complete the prompt and generate a conditional caption. The meaning of "++" is to apply the class token reweighting method specifically to target category word.
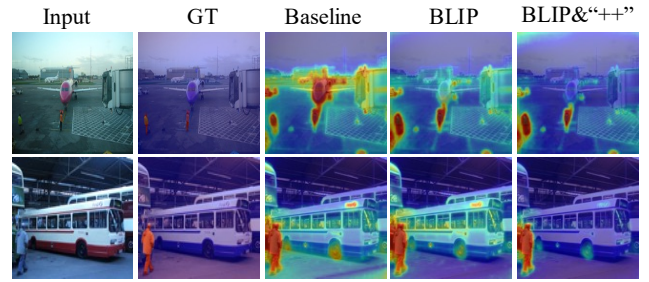


Figure 6: Comparisons of different prompt designs.

**Different Timesteps** In our method, generating attention maps only requires one denoising timestep using the diffusion model. To investigate the impact of the denoising timestep on the results, we conducted more experiments specifically focusing on the timestep. The experimental results are shown in Table 5. It can be observed that the segmentation performance is the best at t=100. If the timestep is too large or too small, it will adversely affect the segmentation results to some extent. By averaging the results from multiple sampling timesteps, the best performance can be achieved. Thus, in weakly supervised setting, we employ average of multiple sampling timesteps, while in zero-shot setting, we simply use t=100 to accelerate the inference speed.

| Method | t=1 | t=50 | t=100 | t=150 | **Avg.** |
|---|---|---|---|---|---|
| **mIoU** | 69.10 | 69.94 | 70.30 | 69.69 | 70.49 |

Table 5: Results of different timesteps. **Avg.** is calculated by averaging the results of t=1,t=50,t=100 and t=150.

## Conclusion

This paper presents an innovative technique called DiffSegmenter, which allows for open-vocabulary semantic segmentation using readily available diffusion models, without requiring additional learnable modules or parameter tunings.

The proposed methodology maximizes the capabilities of attention layers within the denoising U-Net of diffusion models. By harnessing the power of the pre-trained BLIP model, carefully designed textual prompts are incorporated to enhance semantic quality. One limitation of our method is that the Stable Diffusion model relies on the latent features of the input image, which may result in the disappearance of small objects in the small feature maps of the latent space. This issue might be potentially addressed by employing other text-to-image diffusion models that are based on the original images or have larger latent feature map sizes.

# References

Ahn, J.; Cho, S.; and Kwak, S. 2019. Weakly supervised learning of instance segmentation with inter-pixel relations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2209–2218.

Amit, T.; Shaharbany, T.; Nachmani, E.; and Wolf, L. 2021. Segdiff: Image segmentation with diffusion probabilistic models. *arXiv preprint arXiv:2112.00390*.

Baranchuk, D.; Voynov, A.; Rubachev, I.; Khrulkov, V.; and Babenko, A. 2021. Label-Efficient Semantic Segmentation with Diffusion Models. In *International Conference on Learning Representations*.

Blei, D. M.; Kucukelbir, A.; and McAuliffe, J. D. 2017. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518): 859–877.

Cha, J.; Mun, J.; and Roh, B. 2023. Learning to generate text-grounded mask for open-world semantic segmentation from only image-text pairs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11165–11174.

Chang, Y.-T.; Wang, Q.; Hung, W.-C.; Piramuthu, R.; Tsai, Y.-H.; and Yang, M.-H. 2020. Weakly-supervised semantic segmentation via sub-category exploration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8991–9000.

Chen, Z.; Wang, T.; Wu, X.; Hua, X.-S.; Zhang, H.; and Sun, Q. 2022. Class re-activation maps for weakly-supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 969–978.

Dhariwal, P.; and Nichol, A. 2021. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34: 8780–8794.

Ding, J.; Xue, N.; Xia, G.-S.; and Dai, D. 2022. Decoupling zero-shot semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11583–11592.

Everingham, M.; Van Gool, L.; Williams, C. K.; Winn, J.; and Zisserman, A. 2010. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88: 303–338.

Ghiasi, G.; Gu, X.; Cui, Y.; and Lin, T.-Y. 2022. Scaling open-vocabulary image segmentation with image-level labels. In *European Conference on Computer Vision*, 540–557. Springer.

Hertz, A.; Mokady, R.; Tenenbaum, J.; Aberman, K.; Pritch, Y.; and Cohen-Or, D. 2022. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*.

Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33: 6840–6851.

Karazija, L.; Laina, I.; Vedaldi, A.; and Rupprecht, C. 2023a. Diffusion Models for Zero-Shot Open-Vocabulary Segmentation. *arXiv preprint arXiv:2306.09316*.

Karazija, L.; Laina, I.; Vedaldi, A.; and Rupprecht, C. 2023b. Diffusion Models for Zero-Shot Open-Vocabulary Segmentation. *arXiv preprint arXiv:2306.09316*.

Krähenbühl, P.; and Koltun, V. 2011. Efficient inference in fully connected crfs with gaussian edge potentials. *Advances in neural information processing systems*, 24.

Lee, J.; Choi, J.; Mok, J.; and Yoon, S. 2021. Reducing information bottleneck for weakly supervised semantic segmentation. *Advances in Neural Information Processing Systems*, 34: 27408–27421.

Lee, J.; Kim, E.; and Yoon, S. 2021. Anti-adversarially manipulated attributions for weakly and semi-supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4071–4080.

Lee, J.; Oh, S. J.; Yun, S.; Choe, J.; Kim, E.; and Yoon, S. 2022. Weakly supervised semantic segmentation using out-of-distribution data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16897–16906.

Li, A. C.; Prabhudesai, M.; Duggal, S.; Brown, E.; and Pathak, D. 2023. Your diffusion model is secretly a zero-shot classifier. *arXiv preprint arXiv:2303.16203*.

Li, J.; Li, D.; Xiong, C.; and Hoi, S. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, 12888–12900. PMLR.

Liang, F.; Wu, B.; Dai, X.; Li, K.; Zhao, Y.; Zhang, H.; Zhang, P.; Vajda, P.; and Marculescu, D. 2023. Open-vocabulary semantic segmentation with mask-adapted clip. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 7061–7070.

Lin, T.-Y.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Dollár, P.; and Zitnick, C. L. 2014. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, 740–755. Springer.

Lin, Y.; Chen, M.; Wang, W.; Wu, B.; Li, K.; Lin, B.; Liu, H.; and He, X. 2023a. CLIP Is Also an Efficient Segmenter: A Text-Driven Approach for Weakly Supervised Semantic Segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 15305–15314.

Lin, Y.; Chen, M.; Wang, W.; Wu, B.; Li, K.; Lin, B.; Liu, H.; and He, X. 2023b. Clip is also an efficient segmenter:

A text-driven approach for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 15305–15314.

Liu, N.; Du, Y.; Li, S.; Tenenbaum, J. B.; and Torralba, A. 2023. Unsupervised Compositional Concepts Discovery with Text-to-Image Generative Models. *arXiv preprint arXiv:2306.05357*.

Liu, Q.; Wen, Y.; Han, J.; Xu, C.; Xu, H.; and Liang, X. 2022. Open-world semantic segmentation via contrasting and clustering vision-language embedding. In *European Conference on Computer Vision*, 275–292. Springer.

Lu, S.; Liu, Y.; and Kong, A. W.-K. 2023. TF-ICON: Diffusion-Based Training-Free Cross-Domain Image Composition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.

Lüddecke, T.; and Ecker, A. 2022. Image segmentation using text and image prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 7086–7096.

Luo, H.; Bao, J.; Wu, Y.; He, X.; and Li, T. 2023. Segclip: Patch aggregation with learnable centers for open-vocabulary semantic segmentation. In *International Conference on Machine Learning*, 23033–23044. PMLR.

Mottaghi, R.; Chen, X.; Liu, X.; Cho, N.-G.; Lee, S.-W.; Fidler, S.; Urtasun, R.; and Yuille, A. 2014. 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*, 891–898.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.

Ranasinghe, K.; McKinzie, B.; Ravi, S.; Yang, Y.; Toshev, A.; and Shlens, J. 2022. Perceptual grouping in vision-language models. *arXiv preprint arXiv:2210.09996*.

Rao, Y.; Zhao, W.; Chen, G.; Tang, Y.; Zhu, Z.; Huang, G.; Zhou, J.; and Lu, J. 2022. Denseclip: Language-guided dense prediction with context-aware prompting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 18082–18091.

Ren, P.; Li, C.; Xu, H.; Zhu, Y.; Wang, G.; Liu, J.; Chang, X.; and Liang, X. 2023. Viewco: Discovering text-supervised segmentation masks via multi-view semantic consistency. *arXiv preprint arXiv:2302.10307*.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10684–10695.

Shin, G.; Xie, W.; and Albanie, S. 2022. Reco: Retrieve and co-segment for zero-shot transfer. *Advances in Neural Information Processing Systems*, 35: 33754–33767.

Sohl-Dickstein, J.; Weiss, E.; Maheswaranathan, N.; and Ganguli, S. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, 2256–2265. PMLR.

Tang, R.; Pandey, A.; Jiang, Z.; Yang, G.; Kumar, K.; Lin, J.; and Ture, F. 2022. What the daam: Interpreting stable diffusion using cross attention. *arXiv preprint arXiv:2210.04885*.

Wang, Y.; Zhang, J.; Kan, M.; Shan, S.; and Chen, X. 2020. Self-supervised equivariant attention mechanism for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 12275–12284.

Wu, J.; Fang, H.; Zhang, Y.; Yang, Y.; and Xu, Y. 2022. Medsegdiff: Medical image segmentation with diffusion probabilistic model. *arXiv preprint arXiv:2211.00611*.

Wu, W.; Zhao, Y.; Shou, M. Z.; Zhou, H.; and Shen, C. 2023. Diffumask: Synthesizing images with pixel-level annotations for semantic segmentation using diffusion models. *arXiv preprint arXiv:2303.11681*.

Xie, J.; Hou, X.; Ye, K.; and Shen, L. 2022. Clims: Cross language image matching for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4483–4492.

Xing, X.; Wang, C.; Zhou, H.; Zhang, J.; Yu, Q.; and Xu, D. 2023. DiffSketcher: Text Guided Vector Sketch Synthesis through Latent Diffusion Models. *arXiv preprint arXiv:2306.14685*.

Xu, J.; De Mello, S.; Liu, S.; Byeon, W.; Breuel, T.; Kautz, J.; and Wang, X. 2022a. Groupvit: Semantic segmentation emerges from text supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 18134–18144.

Xu, J.; Hou, J.; Zhang, Y.; Feng, R.; Wang, Y.; Qiao, Y.; and Xie, W. 2023a. Learning open-vocabulary semantic segmentation models from natural language supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2935–2944.

Xu, J.; Liu, S.; Vahdat, A.; Byeon, W.; Wang, X.; and De Mello, S. 2023b. Open-vocabulary panoptic segmentation with text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2955–2966.

Xu, L.; Ouyang, W.; Bennamoun, M.; Boussaid, F.; and Xu, D. 2022b. Multi-class token transformer for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4310–4319.

Xu, M.; Zhang, Z.; Wei, F.; Lin, Y.; Cao, Y.; Hu, H.; and Bai, X. 2022c. A simple baseline for open-vocabulary semantic segmentation with pre-trained vision-language model. In *European Conference on Computer Vision*, 736–753. Springer.

Zabari, N.; and Hoshen, Y. 2021. Semantic segmentation in-the-wild without seeing any segmentation examples. *arXiv preprint arXiv:2112.03185*.

Zhou, C.; Loy, C. C.; and Dai, B. 2022. Extract free dense labels from clip. In *European Conference on Computer Vision*, 696–712. Springer.