

# Diffusion Models for Open-Vocabulary Segmentation

Anonymous CVPR submission

Paper ID \*\*\*\*\*

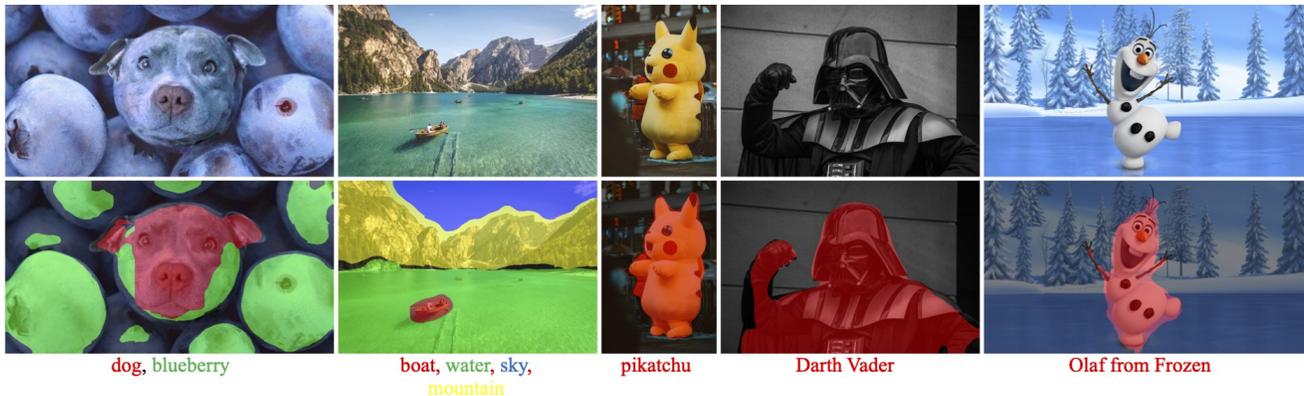


Figure 1. OVDiff is an open-vocabulary segmentation method that, given an image and a free-form set of class names, can segment any user-defined classes. It is fully automatic and does not require any further training.

## Abstract

Open-vocabulary segmentation is the task of segmenting anything that can be named in an image. Recently, large-scale vision-language modelling has led to significant advances in open-vocabulary segmentation, but at the cost of gargantuan and increasing training and annotation efforts. Hence, we ask if it is possible to use existing foundation models to synthesise on-demand efficient segmentation algorithms for specific class sets, making them applicable in an open-vocabulary setting without the need to collect further data, annotations or perform training. To that end, we present OVDiff, a novel method that leverages generative text-to-image diffusion models for unsupervised open-vocabulary segmentation. OVDiff synthesises support image sets for arbitrary textual categories, creating for each a set of prototypes representative of both the category and its surrounding context (background). It relies solely on pre-trained components and outputs the synthesised segmenter directly, without training. Our approach shows strong performance on a range of benchmarks, obtaining a lead of more than 5% over prior work on PASCAL VOC.

## 1. Introduction

Open-vocabulary semantic segmentation is the task of segmenting images into regions matching several free-form textual categories. As the field of Computer Vision moves towards large-scale general-purpose models, open-vocabulary “foundation” models have similarly emerged. Yet, the development of ones suitable for dense localisation tasks such as semantic segmentation incurs both enormous training costs and requires expensive mask annotations. Instead, we show that the open-vocabulary segmentation task can be effectively tackled starting from a set of frozen foundation models, without requiring additional data or even fine-tuning.

In order to do so, we introduce OVDiff, a method that turns existing foundation models into a “factory” of image segmenters, *i.e.*, using foundation models to synthesise on-demand a segmenter for any new concepts specified in natural language. Thus, OVDiff can be used for open-vocabulary segmentation, where it achieves state-of-the-art results in standard benchmarks. Moreover, once synthesised, the segmenters can be efficiently applied to any number of images and easily extended to new categories.

Specifically, segmenting an image using OVDiff can be done in three steps: *generation*, *representation*, and *matching*. Given a textual prompt, OVDiff uses an off-the-shelf text-to-image generator like StableDiffusion [50] to generate

046 a support set of images. In the representation step, we use a  
047 feature extractor (that can be the same network as in the gen-  
048 eration step) to extract feature prototypes that represent the  
049 textual category. Finally, we use simple nearest-neighbour  
050 *matching* scheme to segment the target image using the pro-  
051 totypes computed in the previous step.

052 This approach differs from prior work that largely ap-  
053 proaches the problem in either of two ways. Starting from  
054 multi-modal representations (*e.g.*, CLIP [46]) to bridge vi-  
055 sion and language, the first way relies on labelled data to  
056 fine-tune image-level representations for the segmentation  
057 task. Hence, in line with the zero-shot setting [6], these  
058 methods require costly dense annotations for some known  
059 categories while also extending the segmentation to unseen  
060 categories by incorporating language.

061 The second category of prior work [9, 37, 43, 49, 70, 71]  
062 observes that large-scale vision-language models such as  
063 CLIP have a limited understanding of the positioning of  
064 objects within an image and extend these models with ad-  
065 ditional grouping mechanisms for better localisation using  
066 only image-level captions, but no mask supervision. This,  
067 however, requires expensive additional contrastive training  
068 at scale. Despite yielding promising results, there are some  
069 additional pitfalls to this approach. Firstly, as the text might  
070 not exhaustively describe all entities in the image or might  
071 mention elements that are not depicted, the training signal  
072 can be noisy. Secondly, similar captions may be used to  
073 describe a wide range of visual appearances, or a similar  
074 concept might be described differently, sometimes even de-  
075 pending on the other context present. There is ambiguity and  
076 a difference in detail between visual and textual data. Lastly,  
077 most methods resort to heuristics to segment the background  
078 (*i.e.*, leave some pixels unlabelled), as it often cannot be  
079 described as a textual category. The usual approach is to  
080 threshold the similarities to all categories. Finding an appro-  
081 priate threshold, however, can be challenging and may vary  
082 depending on the image, often resulting in imprecise object  
083 boundaries. Effectively handling the background remains an  
084 open issue.

085 Our three-step approach departs substantially from both  
086 of these schemes. We show that large-scale text-to-image  
087 generative models, such as StableDiffusion [50], can help  
088 bridge the vision-and-language gap without the need for  
089 annotations or costly training. Furthermore, diffusion models  
090 also produce latent spaces that are semantically meaningful  
091 and well-localised. This solves a second problem: multi-  
092 modal embeddings are difficult to learn and often suffer from  
093 ambiguities and differences in detail between modalities.  
094 Instead, our approach can use unimodal features for open-  
095 vocabulary segmentation, which offers several advantages.  
096 Firstly, as text-to-image generators encode a distribution of  
097 possible images, this offers a means to deal with intra-class  
098 variation and captures the ambiguity in textual descriptions.

099 Secondly, the generative image models encode not only the  
100 visual appearance of objects but also provide contextual  
101 priors, which we use for direct background segmentation.

102 This work presents a simple framework that achieves  
103 state-of-the-art performance across open-vocabulary seg-  
104 mentation benchmarks. It combines several off-the-shelf  
105 pre-trained networks into a segmenter “factory” that seg-  
106 ments images into arbitrary textual categories in three simple  
107 steps. OVDiff requires no additional data, mask supervision,  
108 nor fine-tuning. To summarise, we make the following core  
109 contributions: (1) We introduce a method to use pre-trained  
110 diffusion models for the task of open-vocabulary segmen-  
111 tation, that requires no additional data, mask supervision,  
112 or fine-tuning. (2) We propose a principled way to handle  
113 backgrounds by forming prototypes from contextual priors  
114 built into text-to-image generative models. (3) A set of ad-  
115 ditional techniques for further improving performance, such as  
116 multiple prototypes, category filtering and “stuff” filtering.

## 117 2. Related work

118 **Zero-shot open-vocabulary segmentation.** Open-  
119 vocabulary semantic segmentation is a relatively new  
120 problem and is typically approached in two ways. The first  
121 line of work poses the problem as “zero-shot”, *i.e.*, segmen-  
122 ting unseen classes after training on a set of observed classes  
123 with dense annotations. Early approaches [6, 11, 20, 31]  
124 explore generative networks to sample features using  
125 conditional language embeddings for classes. In [30, 69]  
126 image encoders are trained to output dense features that  
127 can be correlated with word2vec [41] and CLIP [46] text  
128 embeddings. Follow-up works [15, 19, 33, 73] approach  
129 the problem in two steps, predicting class-agnostic masks  
130 and aligning the embeddings of masks with language.  
131 IFSeg [74] generates synthetic feature maps by pasting  
132 CLIP text embeddings into a known spatial configuration to  
133 use as additional supervision. Different from our approach,  
134 all these works rely on mask supervision for a set of known  
135 classes.

136 The second line of work eliminates the need for mask  
137 annotations and instead aims to align image regions with  
138 language using only image-text pairs. This is largely en-  
139 abled by recent advancements in large-scale vision-language  
140 models [46]. Some methods introduce internal group-  
141 ing mechanisms such as hierarchical grouping [49, 70],  
142 slot-attention [71], or cross-attention to learn cluster cen-  
143 troids [35, 37]. Assignment to language queries is performed  
144 at group level. Another line of work [9, 43, 48, 79] aims to  
145 learn dense features that are better localised when correlated  
146 with language embeddings at pixel level. With the exception  
147 of [48, 68, 79], thresholding is often required to determine  
148 the background during inference. Alternatively, a curated  
149 list of background prompts can be used [48].

150 Our method falls into the second category. However,

151 in contrast to prior work, we leverage a generative model  
 152 to translate language queries to pre-trained image feature  
 153 extractors without further training. We also segment the  
 154 background directly, without relying on thresholding or  
 155 curated list of background prompts. A closely related ap-  
 156 proach to ours is ReCO [56], where CLIP is used for im-  
 157 age retrieval compiling a set of exemplar images from Im-  
 158 ageNet for a given language query, which is then used for  
 159 co-segmentation. In our method, the shortcoming of an im-  
 160 age database is addressed by synthesising data on-demand.  
 161 Furthermore, instead of co-segmentation, we leverage the  
 162 cross-attention of the generator to extract objects. Instead  
 163 of similarity of support images, we use diverse samples and  
 164 both foreground and contextual backgrounds.

165 **Diffusion models.** Diffusion models [26, 59, 60] are a class  
 166 of generative methods that have seen tremendous success in  
 167 text-to-image systems such as DALL-E [47], Imagen [52],  
 168 and Stable Diffusion [50], trained on Internet-scale data  
 169 such as LAION-5B [54]. The step-wise generative process  
 170 and the language conditioning make pre-trained diffusion  
 171 models attractive also for discriminative tasks. They have  
 172 been recently used in few-shot classification [77], few-shot  
 173 segmentation [2] and panoptic segmentation [72], and to  
 174 generate pairs of images and segmentation masks [32]. How-  
 175 ever, these methods rely on dense manual annotations to  
 176 associate diffusion features with the desired output.

177 Annotation-free discriminative approaches such as [13,  
 178 29] use pre-trained diffusion models as zero-shot classifiers.  
 179 DiffuMask [67] uses prompt engineering to synthesise a  
 180 dataset of “known” and “unseen” categories and trains a  
 181 closed-set segmenter with masks obtained from the cross-  
 182 attention maps of the diffusion model. DiffusionSeg [38]  
 183 uses DDIM inversion [60] to obtain feature maps and at-  
 184 tention masks of object-centric images to perform unsuper-  
 185 vised object discovery, but relies on ImageNet labels and  
 186 is not open-vocabulary. Our approach also leverages the  
 187 rich semantic information present in diffusion models for  
 188 segmentation; unlike these methods, however, it is open-set  
 189 and does not require further training.

190 **Unsupervised segmentation.** Our work is also related to  
 191 unsupervised segmentation approaches. While early works  
 192 relied on hand-crafted priors [12, 44, 66, 75, 76] later ap-  
 193 proaches leverage feature extractors such as DINO [8] and  
 194 perform further analysis of these methods [21, 39, 55, 57,  
 195 58, 63–65]. Some approaches make use of generative meth-  
 196 ods, usually GANs, to separate images in foreground and  
 197 background layers [3–5, 10] or analyse latent structure to  
 198 induce known foreground-background changes [40, 62] to  
 199 synthesise a training dataset with labels. Largely focused on  
 200 unsupervised saliency prediction, these methods are class-  
 201 agnostic and do not incorporate language.

### 3. Method

We present OVDiff, a method for open-vocabulary segmenta-  
 tion, *i.e.*, semantic segmentation of any category described in  
 natural language. We achieve this goal in three steps: (1) we  
 leverage text-to-image generative models to *generate* a set  
 of images representative of the described category, (2) use  
 these to ground *representations* from off-the-shelf pretrained  
 feature extractors, and (3) *match* these against input image  
 features to perform segmentation.

#### 3.1. OVDiff: Diffusion-based open-vocabulary seg- mentation

Our goal is to devise an algorithm which, given a new vo-  
 cabulary of categories  $c_i \in \mathcal{C}$  formulated as natural language  
 queries, can segment any image against it. Let  $I \in \mathbb{R}^{H \times W \times 3}$   
 be an image to be segmented. Let  $\Phi_v : \mathbb{R}^{H \times W \times 3} \rightarrow$   
 $\mathbb{R}^{H'W' \times D}$  be an off-the-shelf visual feature extractor and  
 $\Phi_t : \mathbb{R}^{d_t} \rightarrow \mathbb{R}^D$  a text encoder. Assuming that image and  
 text encoders are aligned, one can achieve segmentation by  
 simply computing a similarity function, for example, the  
 cosine similarity  $s(\Phi_v(I), \Phi_t(c_i))$ , with  $s(x, y) = \frac{x^T y}{\|x\| \|y\|}$ ,  
 between the encoded image  $\Phi_v(I)$  and an encoding of a  
 class label  $c_i$ . To meaningfully compare different modalities,  
 image and text features must lie in a shared representation  
 space, which is typically learned by jointly training  $\Phi_v$  and  
 $\Phi_t$  using image-text or image-label pairs [46].

We propose two modifications to this approach. First, we  
 observe that it is better to compare representations of the  
*same* modality than across vision and language modalities.  
 We thus replace  $\Phi_t(c_i)$  with a  $D$ -dimensional *visual* repre-  
 sentation  $\bar{P}$  of class  $c_i$ , which we refer to as a *prototype*. In  
 this case, the same feature extractor can be used for both pro-  
 totypes and target images; thus, their comparison becomes  
 straightforward and does not necessitate further training.  
 Second, we propose utilising *multiple* prototypes per cate-  
 gory instead of a single class embedding. This enables us to  
 accommodate intra-class variations in appearance, and, as  
 we explain later, it also allows us to exploit contextual priors,  
 which in turn help to segment the background.

Our approach, thus, proceeds in three steps: (1) a set  
 of support images is sampled based on vocabulary  $\mathcal{C}$ , (2) a  
 set of prototypes  $\mathcal{P}$  is calculated, and (3) a set of images  
 $\{I_1, I_2 \dots\}$  is segmented against these prototypes. We ob-  
 serve that in practical applications, whole image collections  
 are processed using the same vocabulary, as altering the set  
 of target classes for individual images in an informed way  
 would already require some knowledge of their contents.  
 Steps (1) and (2) are, thus, performed very infrequently, and  
 their cost is heavily amortised. Next, we detail each step.

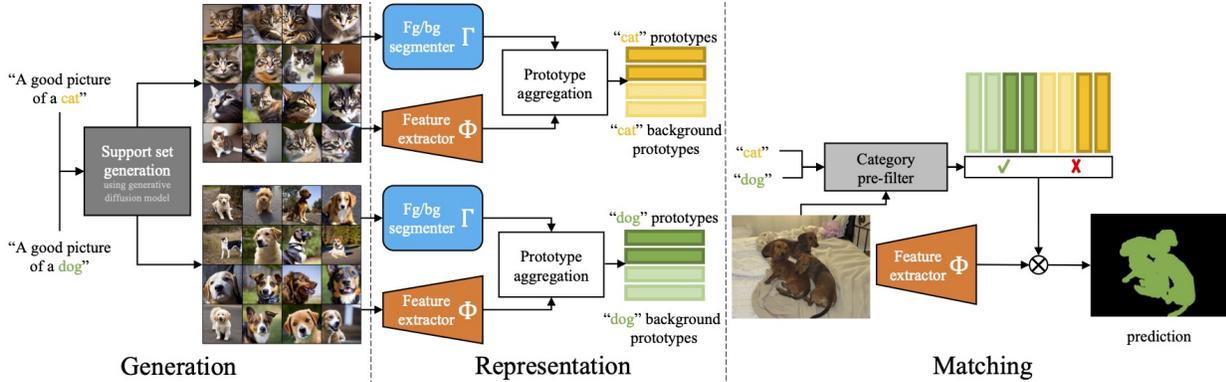


Figure 2. OVDiff overview. Prototype sampling: text queries are used to sample a set of support images which are further processed by a feature extractor and a segmenter forming positive and negative (background) prototypes. Segmentation: image features are compared against prototypes. The CLIP filter removes irrelevant prototypes based on global image contents.

### 250 3.2. Support set generation

251 To construct a set of prototypes, the first step of our approach  
252 is to sample a support set of images representative of each  
253 category  $c_i$ . This can be accomplished by leveraging pre-  
254 trained text-conditional generative models. Sampling images  
255 from a generative model, as opposed to a curated dataset of  
256 real images, aligns well with the goals of open-vocabulary  
257 segmentation as it enables the construction of prototypes for  
258 any user-specified category or description, even those for  
259 which a manually labelled set may not be readily available  
260 (e.g.,  $c_i = \text{"donut with chocolate glaze"}$ ).

261 Specifically, for each query  $c_i$ , we define a prompt “A  
262 good picture of a  $\langle c_i \rangle$ ” and generate a small batch  
263 of  $N$  support images  $\mathcal{S} = \{S_1, S_2, \dots, S_N \mid S_n \in \mathbb{R}^{h \times w \times 3}\}$   
264 of height  $h$  and width  $w$  using Stable Diffusion [50].

### 265 3.3. Representing categories

266 Naïvely, prototypes  $\bar{P}_{c_i}$  could be constructed by averaging  
267 all features across all images for class  $c_i$ . This is unlikely to  
268 result in good prototypes because not all pixels in the sam-  
269 pled images correspond to the class specified by  $c_i$ . Instead,  
270 we propose to extract the class prototypes as follows.

271 **Class prototypes.** Our approach generates two sets of pro-  
272 totypes, positive and negative, for each class. Positive proto-  
273 types are extracted from image regions that are associated  
274 with  $\langle c_i \rangle$ , while negative prototypes represent “background”  
275 regions. Thus, to obtain prototypes, the first step is segmen-  
276 ting the sampled images into foreground and background. To  
277 identify regions most associated with  $c_i$ , we use the fact that  
278 the layout of a generated image is largely dependent on the  
279 cross-attention maps of the diffusion model [24], i.e., pixels  
280 attend more strongly to words that describe them. For a given  
281 word or description (in our case  $c_i$ ), one can generate a set  
282 of attribution maps  $\mathcal{A} = \{A_1, A_2, \dots, A_N \mid A_n \in \mathbb{R}^{h \times w}\}$ ,  
283 corresponding to the support set  $\mathcal{S}$ , by summing the cross-  
284 attention maps across all layers, heads, and denoising steps

of the network [61].

285 Yet, thresholding these attribution maps may not be op-  
286 timal for segmenting foreground/background, as they are  
287 often coarse or incomplete, and sometimes only parts of  
288 objects receive high activation. To improve segmentation  
289 quality, we propose to optionally leverage an unsupervised  
290 instance segmentation method  $\Gamma$ . Unsupervised segmenters  
291 are not vocabulary-aware and may produce multiple binary  
292 object proposals. We denote these as  $\mathcal{M}_n = \{M_{nr} \mid M_{nr} \in$   
293  $\{0, 1\}^{h \times w}\}$ , where  $n$  indexes the support images and  $r$  in-  
294 dexes the object masks (including a mask for the back-  
295 ground). We thus construct a promptable extension of  $\Gamma$   
296 segmenter to select appropriate proposals for foreground  
297 and background: for each image, we select from  $\mathcal{M}_n$  the  
298 mask with the highest (lowest) average attribution as the  
299 foreground (background):  
300

$$M_n^{\text{fg}} = M_{n \in \mathcal{M}_n} \frac{M^\top A_n}{M^\top M}, \quad M_n^{\text{bg}} = M_{n \in \mathcal{M}_n} \frac{M^\top A_n}{M^\top M}. \quad (1) \quad 301$$

302 **Prototype aggregation.** We can compute prototypes  $P_n^g$  for  
303 foreground and background regions ( $g \in \{\text{fg}, \text{bg}\}$ ) as

$$P_n^g = \frac{(\hat{M}_n^g)^\top \Phi_v(S_n)}{m_n^g} \in \mathbb{R}^D, \quad (2) \quad 304$$

305 where  $\hat{M}_n^g$  denotes a resized version of  $M_n^g$  that matches  
306 the spatial dimensions of  $\Phi_v(S_n)$ , and  $m_n^g = (\hat{M}_n^g)^\top \hat{M}_n^g$   
307 counts the number of pixels within each mask. In other  
308 words, prototypes are obtained by means of an off-the-shelf  
309 pretrained feature extractor and computed as the average  
310 feature within each mask.

311 We refer to these as *instance* prototypes because they are  
312 computed from each image individually, and each image in  
313 the support set can be viewed as an instance of class  $c_i$ .

314 In addition to instance prototypes, we found it helpful  
315 to also compute *class-level* prototypes  $\bar{P}^g$  by averaging  
316 instance prototypes weighted by their mask sizes as  $\bar{P}^g =$   
317  $\sum_{n=1}^N m_n^g P_n^g / \sum_{n=1}^N m_n^g$ .

Finally, we propose to augment the set of class and instance prototypes using  $K$ -Means clustering of the masked features to obtain *part-level* prototypes. We perform spatial clustering separately on foreground and background regions and take each cluster centroid as a prototype  $P_k^g$  with  $1 \leq k \leq K$ . The intuition behind this is to enable segmentation at the level of parts, support greater intra-class variability, and a wider range of feature extractors that might not be scale invariant.

We consider the union of all these feature prototypes:

$$\mathcal{P}^g = \bar{P}^g \cup \{P_n^g \mid 1 \leq n \leq N\} \cup \{P_k^g \mid 1 \leq k \leq K\} \quad (3)$$

for  $g \in \{\text{fg}, \text{bg}\}$ , and associate them with a single category.

We note that this process is repeated for each  $c_i \in \mathcal{C}$  and we hereby refer to  $\mathcal{P}^{\text{fg}}$  (and  $\mathcal{P}^{\text{bg}}$ ) as  $\mathcal{P}_{c_i}^{\text{fg}}$  ( $\mathcal{P}_{c_i}^{\text{bg}}$ ), *i.e.*, as the foreground (background) prototypes of class  $c_i$ .

Since  $\mathcal{P}_{c_i}^{\text{fg}}$  ( $\mathcal{P}_{c_i}^{\text{bg}}$ ) depend only on class  $c_i$ , they can be precomputed, and the set of classes can be dynamically expanded without the need to adapt existing prototypes.

### 3.4. Segmentation via prototype matching

To perform segmentation of any target image  $I$  given a vocabulary  $\mathcal{C}$ , we first extract image features using the same visual encoder  $\Phi_v$  used for the prototypes. The vocabulary is expanded with an additional background class  $\hat{\mathcal{C}} = \{c_{\text{bg}}\} \cup \mathcal{C}$ , for which the positive (*foreground*) prototype is the union of all *background* prototypes in the vocabulary:  $\mathcal{P}_{c_{\text{bg}}}^{\text{fg}} = \bigcup_{c_i \in \mathcal{C}} \mathcal{P}_{c_i}^{\text{bg}}$ . Then, a segmentation map can simply be obtained by matching dense image features to prototypes using cosine similarity. A class with the highest similarity in its prototype set is chosen:

$$M = \max_{c \in \hat{\mathcal{C}}} \max_{P \in \mathcal{P}_c^{\text{fg}}} s(\Phi_v(I), P). \quad (4)$$

**Category pre-filtering.** To limit the impact of spurious correlations that might exist in the feature space of the visual encoder, we introduce a pre-filtering process for the target vocabulary given image  $I$ . Specifically, we leverage CLIP [46] as a strong open-vocabulary classifier but propose to apply it in a multi-label fashion to constrain the segmentation to the subset of categories  $\mathcal{C}' \subseteq \mathcal{C}$  that appear in the target image. First, we encode the target image and each category using CLIP. Any categories that do not score higher than  $1/|\mathcal{C}|$  are removed from consideration, that is we keep the subset  $\{P_{c'}^g \mid c' \in \mathcal{C}'\}$ ,  $g \in \{\text{fg}, \text{bg}\}$ . If more than  $\eta$  categories are present, then the top- $\eta$  are selected. We then form “multi-label” prompts as “ $\langle c_a \rangle$  and  $\langle c_b \rangle$  and . . .” where the categories are selected among the top scoring ones taking into account all  $2^\eta$  combinations. The best-scoring multi-label prompt determines the final list of categories to be used in Equation (4).

Table 1. Open-vocabulary segmentation. Comparison of our approach, OVDiff, to the state of the art (under the mIoU metric). Our results are an average of 5 seeds  $\pm \sigma$ . \*results from [9].

Method	Support Set	Further Training	VOC	Context	Object
ReCo* [56]	Real	✗	25.1	19.9	15.7
ViL-Seg [35]	✗	✓	37.3	18.9	-
MaskCLIP* [79]	✗	✗	38.8	23.6	20.6
TCL [9]	✗	✓	51.2	24.3	30.4
CLIPpy [48]	✗	✓	52.2	-	<u>32.0</u>
GroupViT [70]	✗	✓	52.3	22.4	-
ViewCo [49]	✗	✓	52.4	23.0	23.5
SegCLIP [37]	✗	✓	52.6	<u>24.7</u>	26.5
OVSgmentor [71]	✗	✓	53.8	20.4	25.1
CLIP-DIY [68]	✗	✗	<u>59.9</u>	-	31.0
OVDiff (-CutLER) Synth.	✗	✗	62.8	28.6	34.9
<b>OVDiff</b>	Synth.	✗	<b>66.3 <math>\pm</math> 0.2</b>	<b>29.7 <math>\pm</math> 0.3</b>	<b>34.6 <math>\pm</math> 0.3</b>
<hr/>					
TCL [9] (+PAMR)	✗	✓	<u>55.0</u>	<u>30.4</u>	<u>31.6</u>
<b>OVDiff (+PAMR)</b> Synth.	✗	✗	<b>68.4 <math>\pm</math> 0.2</b>	<b>31.2 <math>\pm</math> 0.4</b>	<b>36.2 <math>\pm</math> 0.4</b>

Table 2. Segmentation performance of OVDiff based on different feature extractors.

Feature Extractor	MAE	DINO	CLIP (token)	CLIP (keys)	SD	SD + DINO + CLIP
<b>VOC</b>	54.9	59.1	51.4	61.8	64.4	66.4

**“Stuff” filtering.** Occasionally,  $c_i$  might not describe a countable object category but an identifiable region in the image, *e.g.*, *sky*, often referred to as a “stuff” class. “Stuff” classes warrant additional consideration as they might appear as background in images of other categories, *e.g.*, *boat* images might often contain regions of *water* and *sky*. As a result, the process outlined above might sample background prototypes for one class that coincide with the foreground prototypes of another. To mitigate this issue, we introduce an additional filtering step to detect and reject such prototypes, when the full vocabulary, *i.e.*, the set of classes under consideration, is known. First, we only consider foreground prototypes for “stuff” classes. Additionally, any negative prototypes of “thing” classes with high cosine similarity with any of the “stuff” class prototypes are simply removed. In our experiments, we use ChatGPT [45] to automatically categorise a set of classes as “thing” or “stuff”.

## 4. Experiments

We evaluate OVDiff on the open-vocabulary semantic segmentation task. First, we consider different feature extractors and investigate how they can be grounded by leveraging our approach. We then compare our method with prior work. We ablate the components of OVDiff, visualize the prototypes, and conclude with a qualitative comparison with prior works on in-the-wild images.

**Datasets and implementation details.** As the approach does not require further training of components, we only

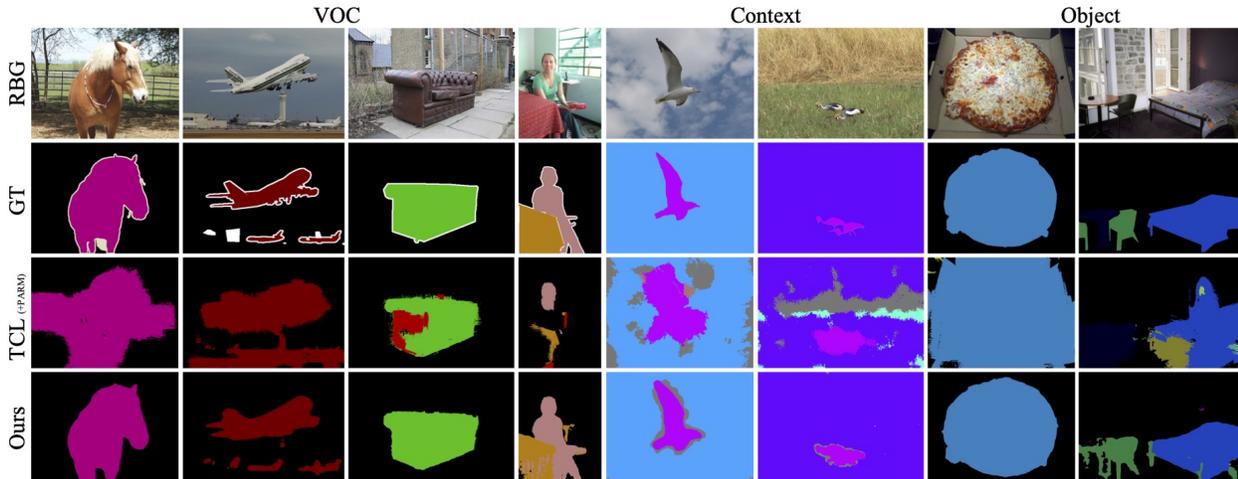


Figure 3. Qualitative results. OVDiff in comparison to TCL (+ PAMR). OVDiff provides more accurate segmentations across a range objects and stuff classes with well defined object boundaries that separate from the background well.

392 consider data for evaluation. Following prior work [70],  
 393 to assess the segmentation performance, we report mean  
 394 Intersection-over-Union (mIoU) on validation splits of PAS-  
 395 CAL VOC (VOC) [18], PASCAL Context (Context) [42] and  
 396 COCO-Object (Object) [7] datasets, with 20, 59, and 80 fore-  
 397 ground classes, respectively. These datasets include a back-  
 398 ground class to reflect a realistic setting of non-exhaustive  
 399 vocabularies. Context also contains both “things” and “stuff”  
 400 classes. We also evaluate without background on VOC, Con-  
 401 text, ADE20K [78], COCO-Stuff [7] and Cityscapes [14],  
 402 with 20, 59, 150, 171, and 19 classes, respectively, but do not  
 403 consider this a realistic setting as it relies on knowing which  
 404 pixels cannot be described by a set of categories. Thus we  
 405 leave such evaluation to Appendix A.3. Similar to [9, 70, 71],  
 406 we employ a sliding window approach. We use two scales to  
 407 aid with the limited resolution of off-the-shelf feature extrac-  
 408 tors with square window sizes of 448 and 336 and a stride  
 409 of 224 pixels. We set the size of the support set to  $N = 32$ .  
 410 For the diffusion model, we use Stable Diffusion v1.5; for  
 411 unsupervised segmenter  $\Gamma$ , we employ CutLER [64].

#### 412 4.1. Grounding feature extractors

413 Our method can be combined with *any* pretrained visual  
 414 feature extractor for constructing prototypes and extracting  
 415 image features. To verify this quantitatively, we experiment  
 416 with various self-supervised ViT feature extractors (Tab. 2):  
 417 DINO [8], MAE [23], and CLIP [46]. We also use SD as a  
 418 feature extractor.

419 We find that SD performs the best, though CLIP and  
 420 DINO also show strong performance based on our experi-  
 421 ments on VOC. MAE shows the weakest performance, which  
 422 may be attributed to its lack of semanticity [23]; yet it is still  
 423 competitive with the majority of purposefully trained net-  
 424 works when employed as part of our approach. We find that  
 425 taking keys of the second to last layer in CLIP yields better

426 results than using patch tokens (CLIP token). As feature  
 427 extractors have different training objectives, we hypothesise  
 428 that their feature spaces might be complementary. Thus, we  
 429 also consider an ensemble approach. In this case, the cosine  
 430 distances formed between features of different extractors  
 431 and respective prototypes are averaged. The combination  
 432 of SD, DINO, and CLIP performs the best. We adopt this  
 433 formulation for the main set of experiments.

#### 434 4.2. Comparison to existing methods

435 In Tab. 1, we compare our method with prior work that does  
 436 not rely on manual mask annotation on three datasets: VOC,  
 437 Context, Object. We include a brief overview of the meth-  
 438 ods in the supplement. We find that our method compares  
 439 favourably, outperforming other methods in all settings. In  
 440 particular, results on VOC show the largest margin, with  
 441 more than 5% improvement over prior work.

442 We also consider a version of our method, OVDiff (-  
 443 CutLER), that does not rely on an additional unsupervised  
 444 segmenter  $\Gamma$ . Instead, the attention masks are thresholded.  
 445 We observe that such a version of OVDiff has strong per-  
 446 formance, outperforming prior work as well. CutLER is  
 447 helpful, but not a critical component, and OVDiff performs  
 448 strongly without it.

449 In the same table, we also combine our method with  
 450 PAMR [1], the post-processing approach employed by TCL.  
 451 We find that it improves results for our method, though im-  
 452 provements are less drastic since our method already yields  
 453 better segmentation and boundaries.

454 Qualitative results are shown in Fig. 3. This figure high-  
 455 lights a key benefit of our approach: the ability to exploit  
 456 contextual priors through the use of background prototypes,  
 457 which in turn allows for the direct assignment of pixels to  
 458 a background class. This improves segmentation quality  
 459 because it makes it easier to differentiate objects from the

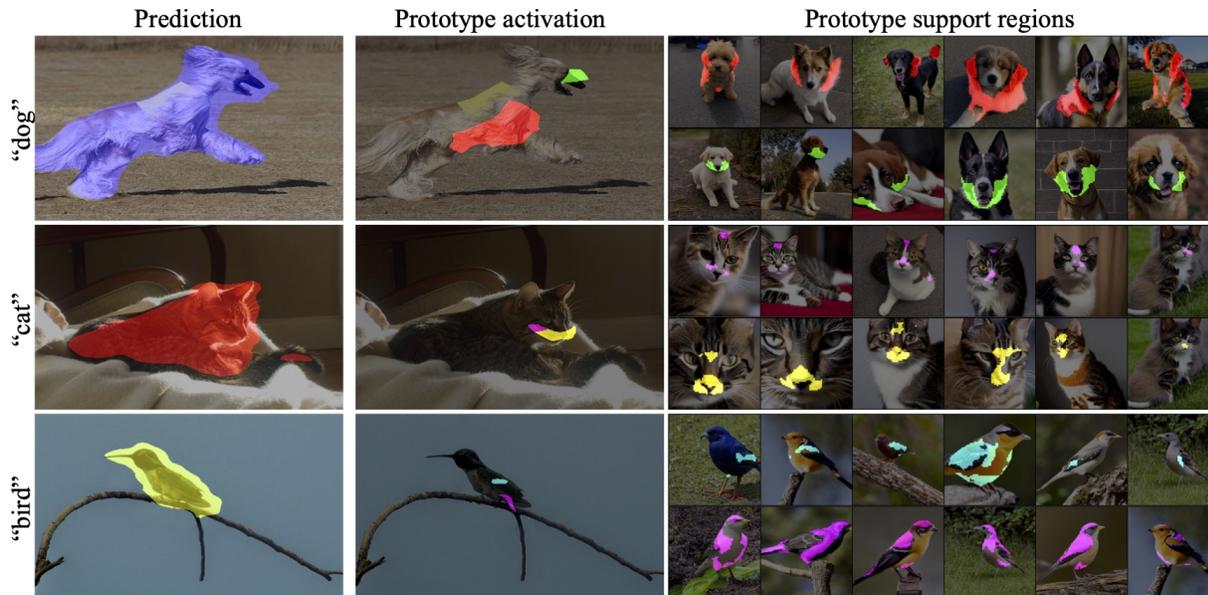


Figure 4. Analysis of the segmentation output by linking regions to samples in the support set. Left: our results for different classes. Middle: select color-coded regions “activated” by different prototypes for the class. Right: regions in the support set images corresponding to these (part-level) prototypes.

Table 3. Ablation of different components. Each component is removed in isolation, measuring the drop ( $\Delta$ ) in mIoU on VOC and Context datasets. Using SD features.

Configuration	VOC	$\Delta$	Context	$\Delta$
Full	64.4		29.4	
w/o bg prototypes	53.2	-11.2	28.9	-0.5
w/o category filter	54.4	-10.0	25.2	-4.2
w/o “stuff” filter	n/a		26.9	-2.5
w/o CutLER	60.4	-4.0	27.6	-1.8
w/o sliding window	62.2	-2.2	28.6	-0.8
only average $P$	62.5	-1.9	28.4	-1.0

460 background and to delineate their boundaries. In comparison,  
461 TCL predictions are very coarse and contain more noise.

462 **Computation cost.** We focus on a construction of a method  
463 to show that existing foundational diffusion models can be  
464 used for segmentation with great efficacy without further  
465 training. OVDiff requires computing prototypes instead.  
466 With our unoptimized implementation, we measure around  
467  $110 \pm 10$ s to calculate prototypes using SD for a single  
468 class, or around 1.14 TFLOP/s-hours of compute. While the  
469 focus of this study is not computational efficiency, we can  
470 compare prototype sampling to the cost of additional training  
471 of other methods: TCL requires 2688, GroupViT 10752, and  
472 OVSegmentor 624 TFLOP/s-hours.<sup>1</sup> While training has an  
473 upfront compute cost and requires special infrastructure (e.g.  
474 OVSegmentor uses  $16 \times A100$ s), OVDiff’s prototype set can  
475 be grown progressively as needed, while showing better  
476 performance.

<sup>1</sup>Estimated as training time  $\times$  num. GPUs  $\times$  theoretical peak TFLOP/s for GPU type.

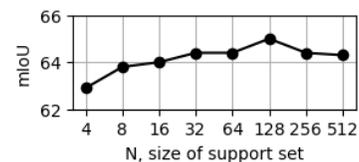


Figure 5. PascalVOC results with increasing support size  $N$ .

### 4.3. Ablations

477  
478 Next, we ablate the components of OVDiff on VOC and Con-  
479 text datasets. For these experiments, only SD is employed  
480 as a feature extractor. We remove individual components  
481 and measure the change in segmentation performance, sum-  
482 marising the results in Tab. 3. Our first observation is that  
483 background prototypes have a major impact on performance.  
484 When removing them from consideration, we instead thresh-  
485 old the similarity scores of the images with the foreground  
486 prototypes (set to 0.72, determined via grid search); in this  
487 case, the performance drops significantly, which again high-  
488 lights the importance of leveraging contextual priors.  
489 On Context, the impact is less significant, likely due to the  
490 fact that the dataset contains “stuff” categories. Remov-  
491 ing the *instance*- and *part-level* prototypes also negatively  
492 affects performance. Additionally, removing the category  
493 pre-filtering has a major impact. We hypothesize that this  
494 introduces spurious correlations between prototypes of dif-  
495 ferent classes. On Context, “stuff” filtering is also important.  
496 We again consider the importance of using an unsupervised  
497 segmenter, CutLER, for prototype mask extractions, using  
498 thresholding instead. We find this slightly reduces perfor-  
499 mance in this setting as well. Overall, background prototypes  
500 and pre-filtering contribute the most.

501 Finally, we measure the effect of varying the size of the



Figure 6. Qualitative comparison on challenging in-the-wild images with TCL, which struggles with object boundaries, missing parts of objects, or including surroundings. Our method has more appropriate boundaries and makes fewer errors overall, but does produce a small halo effect around objects due to the upscaling of feature extractors.

support set  $N$  in Fig. 5. We find that OVDiff already shows strong performance even at a low number of samples for each query. With increasing the number of samples, the performance improves, saturating at around  $N = 32$ , which we use in our main experiments.

#### 4.4. Explaining segmentations

We inspect how our method segments certain regions by considering which prototype from  $\mathcal{P}_c^{\text{fg}}$  was used to assign a class  $c$  to a pixel. Prototypes map to regions in the support set from where they were aggregated, e.g., instances of prototypes are associated with foreground masks  $M_n^{\text{fg}}$  and part prototypes with centroids/clusters. By following these mappings, a set of support image regions can be retrieved for each segmentation decision, providing a degree of explainability. Fig. 4 illustrates this for examples of dog, cat, and bird classes. For visualisation purposes, selected prototypes and corresponding regions are shown. On the left, we show the full segmentation result of each image. In the middle, we select regions that correlate best with certain class prototypes. On the right, we retrieve images from the support set and highlight where each prototype emerged. We find that meaningful part segmentation merges due to clustering the support image features, and similar regions are segmented by corresponding prototypes. However, sometimes region covered in the input image will not fully align with the whole prototype (e.g. cat’s face around the eyes or lower belly/tail of bird). Each segmentation is explained by precise regions in a small support set.

#### 4.5. In-the-wild

In Fig. 6, we investigate OVDiff on challenging in-the-wild images with simple and complex backgrounds. We compare with TCL+PAMR. In the first three images, both methods

correctly detect the objects identified by the queries. OVDiff has small false positive “corgi” patches. TCL however misses large parts of the objects, such as most of the person, and parts of animal bodies. The distinction between the house and the bridge in the second image is also better with OVDiff. We also note that our segmentations sometimes have halos around objects. This is caused by upscaling the low-resolution feature extractor (SD in this case). The last two images contain challenging scenarios where both approaches struggle. The fourth image only contains similar objects of the same type. Both methods incorrectly identify plain donuts as either of the specified queries. OVDiff however correctly identifies chocolate donuts with varied sprinkles and separates all donuts from the background. In the final picture, the query “red car” is added, although no such object is present. The extra query causes TCL to incorrectly identify parts of the red bus as a car. Both methods incorrectly segment the gray car in the distance. However, overall, our method is more robust and delineates objects better despite the lack of specialized training or post-processing.

## 5. Conclusion

We introduce OVDiff, an open-vocabulary segmentation method that operates in two stages. First, given queries, support images are sampled and their features are extracted to create class prototypes. These prototypes are then compared to features from an inference image. This approach offers multiple advantages: diverse prototypes accommodating various visual appearances and negative prototypes for background localisation. OVDiff outperforms prior work on benchmarks, exhibiting fewer errors, effectively separating objects from background, and providing explainability through segmentation mapping to support set regions.

566 **References**

- 567 [1] Nikita Araslanov and Stefan Roth. Single-stage seman-  
568 tic segmentation from image labels. In *Proceedings of*  
569 *the IEEE/CVF Conference on Computer Vision and Pattern*  
570 *Recognition*, pages 4253–4262, 2020. 6
- 571 [2] Dmitry Baranchuk, Andrey Voynov, Ivan Rubachev, Valentin  
572 Khrukov, and Artem Babenko. Label-efficient semantic seg-  
573 mentation with diffusion models. In *International Conference*  
574 *on Learning Representations*, 2022. 3
- 575 [3] Yaniv Benny and Lior Wolf. Onegan: Simultaneous unsuper-  
576 vised learning of conditional image generation, foreground  
577 segmentation, and fine-grained clustering. In *European Con-*  
578 *ference on Computer Vision*, pages 514–530. Springer, 2020.  
579 3
- 580 [4] Adam Bielski and Paolo Favaro. Emergence of object segmen-  
581 tation in perturbed generative models. *Advances in Neural*  
582 *Information Processing Systems*, 32, 2019.
- 583 [5] Adam Bielski and Paolo Favaro. Move: Unsupervised mov-  
584 able object segmentation and detection. In *Advances in Neural*  
585 *Information Processing Systems*, 2022. 3
- 586 [6] Maxime Bucher, Tuan-Hung Vu, Matthieu Cord, and Patrick  
587 Pérez. Zero-shot semantic segmentation. *Advances in Neural*  
588 *Information Processing Systems*, 32, 2019. 2
- 589 [7] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-  
590 stuff: Thing and stuff classes in context. In *Computer vision*  
591 *and pattern recognition (CVPR), 2018 IEEE conference on*.  
592 IEEE, 2018. 6, 19
- 593 [8] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou,  
594 Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerg-  
595 ing properties in self-supervised vision transformers. In *Pro-*  
596 *ceedings of the IEEE/CVF international conference on com-*  
597 *puter vision*, pages 9650–9660, 2021. 3, 6, 18
- 598 [9] Junbum Cha, Jonghwan Mun, and Byungseok Roh. Learn-  
599 ing to generate text-grounded mask for open-world semantic  
600 segmentation from only image-text pairs. *arXiv preprint*  
601 *arXiv:2212.00785*, 2022. 2, 5, 6, 13, 15, 19
- 602 [10] Mickaël Chen, Thierry Artières, and Ludovic Denoyer. Un-  
603 supervised object segmentation by redrawing. *Advances in*  
604 *neural information processing systems*, 32, 2019. 3
- 605 [11] Jiaxin Cheng, Soumyaroop Nandi, Prem Natarajan, and Wael  
606 Abd-Almageed. Sign: Spatial-information incorporated gen-  
607 erative network for generalized zero-shot semantic segmen-  
608 tation. In *Proceedings of the IEEE/CVF International Con-*  
609 *ference on Computer Vision (ICCV)*, pages 9556–9566, 2021.  
610 2
- 611 [12] Ming-Ming Cheng, Niloy J. Mitra, Xiaolei Huang, Philip  
612 H. S. Torr, and Shi-Min Hu. Global contrast based salient  
613 region detection. *IEEE Transactions on Pattern Analysis and*  
614 *Machine Intelligence*, 37(3):569–582, 2015. 3
- 615 [13] Kevin Clark and Priyank Jaini. Text-to-image diffusion mod-  
616 els are zero-shot classifiers. *arXiv preprint arXiv:2303.15233*,  
617 2023. 3
- 618 [14] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo  
619 Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke,  
620 Stefan Roth, and Bernt Schiele. The cityscapes dataset for  
621 semantic urban scene understanding. In *Proceedings of the*  
*IEEE Conference on Computer Vision and Pattern Recogni-*  
*tion (CVPR)*, 2016. 6
- [15] Jian Ding, Nan Xue, Gui-Song Xia, and Dengxin Dai. De-  
coupling zero-shot semantic segmentation. In *Proceedings of*  
*the IEEE/CVF Conference on Computer Vision and Pattern*  
*Recognition*, pages 11583–11592, 2022. 2
- [16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov,  
Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,  
Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl-  
vain Gelly, et al. An image is worth 16x16 words: Transform-  
ers for image recognition at scale. In *International Conference*  
*on Learning Representations*, 2021. 18
- [17] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn,  
and A. Zisserman. The pascal visual object classes (voc)  
challenge. *International Journal of Computer Vision*, 88(2):  
303–338, 2010. 19
- [18] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn,  
and A. Zisserman. The PASCAL Visual Object Classes  
Challenge 2012 (VOC2012) Results. [http://www.pascal-](http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html)  
[network.org/challenges/VOC/voc2012/workshop/index.html](http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html),  
2012. 6, 19
- [19] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Scaling  
open-vocabulary image segmentation with image-level labels.  
In *Computer Vision—ECCV 2022: 17th European Conference,*  
*Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part*  
*XXXVI*, pages 540–557. Springer, 2022. 2
- [20] Zhangxuan Gu, Siyuan Zhou, Li Niu, Zihan Zhao, and Liqing  
Zhang. Context-aware feature generation for zero-shot se-  
mantic segmentation. In *Proceedings of the 28th ACM Inter-*  
*national Conference on Multimedia*, pages 1921–1929, 2020.  
2
- [21] Mark Hamilton, Zhoutong Zhang, Bharath Hariharan, Noah  
Snavely, and William T Freeman. Unsupervised semantic  
segmentation by distilling feature correspondences. In *Inter-*  
*national Conference on Learning Representations*, 2022.  
3
- [22] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Gir-  
shick. Mask r-cnn. In *Proceedings of the IEEE international*  
*conference on computer vision*, pages 2961–2969, 2017. 18
- [23] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr  
Dollár, and Ross Girshick. Masked autoencoders are scalable  
vision learners. In *Proceedings of the IEEE/CVF Conference*  
*on Computer Vision and Pattern Recognition*, pages 16000–  
16009, 2022. 6
- [24] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman,  
Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt im-  
age editing with cross attention control. *arXiv preprint*  
*arXiv:2208.01626*, 2022. 4
- [25] Jonathan Ho and Tim Salimans. Classifier-free diffusion  
guidance. In *NeurIPS 2021 Workshop on Deep Generative*  
*Models and Downstream Applications*. 19
- [26] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffu-  
sion probabilistic models. *Advances in Neural Information*  
*Processing Systems*, 33:6840–6851, 2020. 3
- [27] Ronghang Hu, Shoubhik Debnath, Saining Xie, and Xinlei  
Chen. Exploring long-sequence masked autoencoders. *arXiv*  
*preprint arXiv:2210.07224*, 2022. 18

- [28] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. 2014. 18
- [29] Alexander C Li, Mihir Prabhudesai, Shivam Duggal, Ellis Brown, and Deepak Pathak. Your diffusion model is secretly a zero-shot classifier. *arXiv preprint arXiv:2303.16203*, 2023. 3
- [30] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and Rene Ranftl. Language-driven semantic segmentation. In *International Conference on Learning Representations*, 2021. 2
- [31] Peike Li, Yunchao Wei, and Yi Yang. Consistent structural relation learning for zero-shot segmentation. *Advances in Neural Information Processing Systems*, 33:10317–10327, 2020. 2
- [32] Ziyi Li, Qinye Zhou, Xiaoyun Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Guiding text-to-image diffusion model towards grounded generation. *arXiv:2301.05221*, 2023. 3
- [33] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. *arXiv preprint arXiv:2210.04150*, 2022. 2
- [34] Ping-Sung Liao, Tse-Sheng Chen, Pau-Choo Chung, et al. A fast algorithm for multilevel thresholding. *J. Inf. Sci. Eng.*, 17(5):713–727, 2001. 19
- [35] Quande Liu, Youpeng Wen, Jianhua Han, Chunjing Xu, Hang Xu, and Xiaodan Liang. Open-world semantic segmentation via contrasting and clustering vision-language embedding. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XX*, pages 275–292. Springer, 2022. 2, 5, 19
- [36] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022. 19
- [37] Huaishao Luo, Junwei Bao, Youzheng Wu, Xiaodong He, and Tianrui Li. SegCLIP: Patch aggregation with learnable centers for open-vocabulary semantic segmentation. *arXiv preprint arXiv:2211.14813*, 2022. 2, 5, 19
- [38] Chaofan Ma, Yuhuan Yang, Chen Ju, Fei Zhang, Jinxiang Liu, Yu Wang, Ya Zhang, and Yanfeng Wang. Diffusionseg: Adapting diffusion towards unsupervised object discovery. *arXiv preprint arXiv:2303.09813*, 2023. 3
- [39] Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. Deep spectral methods: A surprisingly strong baseline for unsupervised semantic segmentation and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8364–8375, 2022. 3
- [40] Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. Finding an unsupervised image segmenter in each of your deep generative models. In *International Conference on Learning Representations*, 2022. 3
- [41] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013. 2
- [42] 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. 6, 19
- [43] 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. *arXiv preprint arXiv:2212.04994*, 2022. 2
- [44] Tam Nguyen, Maximilian Dax, Chaithanya Kumar Mummadi, Nhung Ngo, Thi Hoai Phuong Nguyen, Zhongyu Lou, and Thomas Brox. Deepusps: Deep robust unsupervised saliency prediction via self-supervision. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2019. 3
- [45] OpenAI. Introducing chatgpt. <https://openai.com/blog/chatgpt>, 2023. 5
- [46] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2, 3, 5, 6, 18
- [47] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 3
- [48] Kanchana Ranasinghe, Brandon McKinzie, Sachin Ravi, Yinfei Yang, Alexander Toshev, and Jonathon Shlens. Perceptual grouping in vision-language models. *arXiv preprint arXiv:2210.09996*, 2022. 2, 5, 13, 19
- [49] Pengzhen Ren, Changlin Li, Hang Xu, Yi Zhu, Guangrun Wang, Jianzhuang Liu, Xiaojun Chang, and Xiaodan Liang. Viewco: Discovering text-supervised segmentation masks via multi-view semantic consistency. *arXiv preprint arXiv:2302.10307*, 2023. 2, 5, 19
- [50] 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. 1, 2, 3, 4, 18
- [51] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015. 18
- [52] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 3
- [53] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of* 736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792

- 793 *the IEEE/CVF Conference on Computer Vision and Pattern*  
794 *Recognition (CVPR)*, pages 22522–22531, 2023. 15
- [54] Christoph Schuhmann, Romain Beaumont, Richard Vencu,  
795 Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes,  
796 Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al.  
797 Laion-5b: An open large-scale dataset for training next generation  
798 image-text models. *arXiv preprint arXiv:2210.08402*,  
799 2022. 3
- [55] Gyungin Shin, Samuel Albanie, and Weidi Xie. Unsuper-  
800 vised salient object detection with spectral cluster voting. In  
801 *CVPRW*, 2022. 3
- [56] Gyungin Shin, Weidi Xie, and Samuel Albanie. Reco: Re-  
802 trieve and co-segment for zero-shot transfer. In *Advances in*  
803 *Neural Information Processing Systems (NeurIPS)*, 2022. 3,  
804 5, 19
- [57] Oriane Siméoni, Gilles Puy, Huy V. Vo, Simon Roburin, Spy-  
805 ros Gidaris, Andrei Bursuc, Patrick Pérez, Renaud Marlet,  
806 and Jean Ponce. Localizing objects with self-supervised trans-  
807 formers and no labels. 2021. 3
- [58] Oriane Siméoni, Chloé Sekkat, Gilles Puy, Antonin Vobecky,  
808 Éloi Zablocki, and Patrick Pérez. Unsupervised object local-  
809 ization: Observing the background to discover objects. *arXiv*  
810 *preprint arXiv:2212.07834*, 2022. 3
- [59] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan,  
811 and Surya Ganguli. Deep unsupervised learning using  
812 nonequilibrium thermodynamics. In *International Confer-*  
813 *ence on Machine Learning*, pages 2256–2265. PMLR, 2015.  
814 3
- [60] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Ab-  
815 hishek Kumar, Stefano Ermon, and Ben Poole. Score-based  
816 generative modeling through stochastic differential equations.  
817 In *International Conference on Learning Representations*,  
818 2021. 3
- [61] Raphael Tang, Akshat Pandey, Zhiying Jiang, Gefei Yang,  
819 Karun Kumar, Jimmy Lin, and Ferhan Ture. What the daam:  
820 Interpreting stable diffusion using cross attention. *arXiv*  
821 *preprint arXiv:2210.04885*, 2022. 4
- [62] Andrey Voynov, Stanislav Morozov, and Artem Babenko. Ob-  
822 ject segmentation without labels with large-scale generative  
823 models. In *International Conference on Machine Learning*,  
824 pages 10596–10606. PMLR, 2021. 3
- [63] Xinlong Wang, Zhiding Yu, Shalini De Mello, Jan Kautz,  
825 Anima Anandkumar, Chunhua Shen, and Jose M Alvarez.  
826 Freesolo: Learning to segment objects without annotations.  
827 In *Proceedings of the IEEE/CVF Conference on Computer*  
828 *Vision and Pattern Recognition*, pages 14176–14186, 2022. 3
- [64] Xudong Wang, Rohit Girdhar, Stella X Yu, and Ishan Misra.  
829 Cut and learn for unsupervised object detection and instance  
830 segmentation. *arXiv preprint arXiv:2301.11320*, 2023. 6, 13
- [65] Yangtao Wang, Xi Shen, Shell Xu Hu, Yuan Yuan, James L.  
831 Crowley, and Dominique Vaufreydaz. Self-supervised trans-  
832 formers for unsupervised object discovery using normalized  
833 cut. In *Proceedings of the IEEE/CVF Conference on Com-*  
834 *puter Vision and Pattern Recognition (CVPR)*, pages 14543–  
835 14553, 2022. 3
- [66] Yichen Wei, Fang Wen, Wangjiang Zhu, and Jian Sun.  
836 Geodesic saliency using background priors. In *ECCV*, 2012.  
837 3
- [67] Weijia Wu, Yuzhong Zhao, Mike Zheng Shou, Hong Zhou,  
838 and Chunhua Shen. Diffumask: Synthesizing images with  
839 pixel-level annotations for semantic segmentation using diffu-  
840 sion models. *arXiv preprint arXiv:2303.11681*, 2023. 3
- [68] Monika Wysockańska, Michaël Ramamonjisoa, Tomasz Trz-  
841 ciński, and Oriane Siméoni. Clip-diy: Clip dense inference  
842 yields open-vocabulary semantic segmentation for-free. In  
843 *Proceedings of the IEEE/CVF Winter Conference on Appli-*  
844 *cations of Computer Vision*, pages 1403–1413, 2024. 2, 5,  
845 19
- [69] Yongqin Xian, Subhabrata Choudhury, Yang He, Bernt  
846 Schiele, and Zeynep Akata. Semantic projection network  
847 for zero-and few-label semantic segmentation. In *Proceed-*  
848 *ings of the IEEE/CVF Conference on Computer Vision and*  
849 *Pattern Recognition*, pages 8256–8265, 2019. 2
- [70] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon,  
850 Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit:  
851 Semantic segmentation emerges from text supervision. In  
852 *Proceedings of the IEEE/CVF Conference on Computer Vi-*  
853 *sion and Pattern Recognition*, pages 18134–18144, 2022. 2,  
854 5, 6, 19
- [71] Jilan Xu, Junlin Hou, Yuejie Zhang, Rui Feng, Yi Wang, Yu  
855 Qiao, and Weidi Xie. Learning open-vocabulary semantic seg-  
856 mentation models from natural language supervision. *arXiv*  
857 *preprint arXiv:2301.09121*, 2023. 2, 5, 6, 13, 19
- [72] Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong  
858 Wang, and Shalini De Mello. Open-vocabulary panoptic  
859 segmentation with text-to-image diffusion models. *arXiv*  
860 *preprint arXiv:2303.04803*, 2023. 3
- [73] Mengde Xu, Zheng Zhang, Fangyun Wei, Yutong Lin, Yue  
861 Cao, Han Hu, and Xiang Bai. A simple baseline for open-  
862 vocabulary semantic segmentation with pre-trained vision-  
863 language model. In *European Conference on Computer Vi-*  
864 *sion*, pages 736–753, 2022. 2
- [74] Sukmin Yun, Seong Hyeon Park, Paul Hongsuck Seo, and  
865 Jinwoo Shin. Ifseg: Image-free semantic segmentation via  
866 vision-language model. *arXiv preprint arXiv:2303.14396*,  
867 2023. 2
- [75] Yu Zeng, Yunzhi Zhuge, Huchuan Lu, Lihe Zhang, Mingyang  
868 Qian, and Yizhou Yu. Multi-source weak supervision for  
869 saliency detection. In *IEEE Conference on Computer Vision*  
870 *and Pattern Recognition*, 2019. 3
- [76] Jing Zhang, T. Zhang, Yuchao Dai, Mehrtash Harandi, and  
871 Richard I. Hartley. Deep unsupervised saliency detection: A  
872 multiple noisy labeling perspective. *2018 IEEE/CVF Con-*  
873 *ference on Computer Vision and Pattern Recognition*, pages  
874 9029–9038, 2018. 3
- [77] Renrui Zhang, Xiangfei Hu, Bohao Li, Siyuan Huang, Han-  
875 qiu Deng, Hongsheng Li, Yu Qiao, and Peng Gao. Prompt,  
876 generate, then cache: Cascade of foundation models makes  
877 strong few-shot learners. *arXiv preprint arXiv:2303.02151*,  
878 2023. 3
- [78] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Bar-  
879 riuoso, and Antonio Torralba. Scene parsing through ade20k  
880 dataset. In *Proceedings of the IEEE conference on computer*  
881 *vision and pattern recognition*, pages 633–641, 2017. 6, 15
- [79] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free  
882 dense labels from clip. In *Computer Vision–ECCV 2022: 17th*  
883 907  
884 908

909 *European Conference, Tel Aviv, Israel, October 23–27, 2022,*  
910 *Proceedings, Part XXVIII*, pages 696–712. Springer, 2022. 2,  
911 5, 19

912 **Supplementary Material**

913 In this supplementary material, we provide additional exper-  
914 imental results, including further ablations and qualitative  
915 comparisons (Appendix A), consider the limitations and  
916 broader impacts of our work (Appendix B), and conclude  
917 with additional details concerning the implementation (Ap-  
918 pendix C).

919 **A. Additional experiments**

920 This section provides additional experimental results of  
921 OVDiff.

922 **A.1. Additional Comparisons**

923 **Category filter.** To ensure that the category pre-filtering  
924 does not give our approach an unfair advantage, we augment  
925 two methods (TCL [9] and OVSegmentor [71], which are  
926 the closest baselines with code and checkpoints available)  
927 with our category pre-filtering. We evaluate on the Pascal  
928 VOC dataset (where the category filter shows a significant  
929 impact; see Table 3) and report the results in Tab. A.1. We  
930 observe that TCL improves by 0.6, while the performance  
931 of OVSegmentor drops by 0.1. On the contrary, our method  
932 benefits substantially from this component, but it still shows  
933 stronger performance without the filter than baselines with.  
934 **Influence of  $\Gamma$  segmentation method.** We also further in-  
935 vestigate the use of CutLER [64] to obtain segmentation  
936 masks. We also provide example results of segmentation in  
937 Fig. C.4. In Tab. A.2, we devise a baseline where CutLER-  
938 predicted masks are used to average the CLIP image en-  
939 coder’s final spatial tokens after projection. Averaged tokens  
940 are compared with CLIP text embeddings to assign a class.  
941 While relying on pre-trained components (like ours), this  
942 avoids support set generation. In the same table, we also con-  
943 sider whether the objectness prior provided by CutLER could  
944 be beneficial to other methods as well. We consider a version  
945 of TCL [9] and OVSegmentor [71] which we augment with  
946 CutLER. That is, after methods assign class probabilities to  
947 each pixel/patch, a majority voting for a class is performed in  
948 every region predicted by CutLER. This combines CutLER’s  
949 understanding of objects and their boundaries, aspects where  
950 prior methods struggle, with open-vocabulary segmentation.  
951 However, we observe that this negatively impacts the perfor-  
952 mance of these methods, which we attribute to only a limited  
953 performance of CutLER in complex scenes present in the  
954 datasets. Finally, we also include a version of OVDiff that  
955 does not rely on CutLER for mask extractions, instead using  
956 thresholded masks. We observe that such a version of our  
957 method also has strong performance.

958 We additionally experiment with stronger segmenters to  
959 understand the influence of FG/BG mask quality. We replace  
960 our FG/BG segmentation approach with strong supervised  
961 models: with SAM, we achieve 67.1 on VOC, and with

Table A.1. Use of category filter component. OVDiff without category filter outperforms prior work with cat. filter.

Model	Category filter	
	✗	✓
OVSegmentor	53.8	53.7
TCL	51.2	51.8
TCL (+PAMR)	55.0	56.0
OVDiff	<b>56.2</b>	<b>66.4</b>

Table A.2. Application of CutLER. Prior work does not benefit from using CutLER during inference, while OVDiff shows strong results without it.

Model	CutLER	VOC	Context	Object
CLIP	✓	33.0	11.6	11.1
OVSegmentor		53.8	20.4	25.1
OVSegmentor	✓	38.7	14.4	16.8
TCL		51.2	24.3	30.4
TCL	✓	43.1	20.5	22.7
OVDiff		62.8	28.6	34.9
OVDiff	✓	<b>66.3 ± 0.2</b>	<b>29.7 ± 0.3</b>	<b>34.6 ± 0.3</b>

962 Grounded SAM, 68.5. This slightly improves results from  
963 66.3 of our configuration with CutLER, but the performance  
964 gain is not large and thus not critical.

965 **Class prompts.** We additionally consider whether correc-  
966 tions introduced to class prompts might have similarly pro-  
967 vided additional benefits to our approach (see Appendix C.3  
968 for details). To that end, we also evaluate TCL and OVSeg-  
969 menter (methods that do not rely on additional prompt cu-  
970 ration) with our corrected prompts and consider a version  
971 of our method without such corrections in Tab. A.3. We  
972 observe only marginal to no impact on the performance.

973 **Prompt template** Finally, we consider the prompt tem-  
974 plate employed when sampling support image set: “A good  
975 picture of a  $\langle c_i \rangle$ ” for class prompt  $c_i$ . This template  
976 is generic and broadly applicable to virtually any natural  
977 language specification of a target class. While prior work  
978 adopts prompt expansion by considering a list of synonyms  
979 and subcategories, it is not entirely clear how such a strat-  
980 egy could be systematically performed for any in-the-wild  
981 prompts, such as a “chocolate glazed donut”. We experiment  
982 with a list of synonyms and subclasses, as employed by [48],  
983 on VOC datasets measuring 66.4 mIoU, which is similar to  
984 our single prompt performance  $66.3 \pm 0.2$ . Curating such  
985 lists automatically is an interesting future scaling direction.

986 **A.2. Additional ablations**

987 **Prototype combinations.** In Tab. A.6, we consider the three  
988 different types of prototypes described in Section 3 and test  
989 their performance individually and in various combinations.  
990 We find that the “part” prototypes obtained by  $K$ -means

Table A.3. Using corrected prompts. We consider if corrected class names benefit prior work. We observe negligible to no effect.

Model	Correction	VOC	Context	Object
OVSegmentor		53.8	20.4	25.1
OVSegmentor	✓	53.9	20.4	25.1
TCL		51.2	24.3	30.4
TCL	✓	50.6	24.3	30.4
OVDiff		66.1	29.5	34.9
OVDiff	✓	<b>66.3 ± 0.2</b>	<b>29.7 ± 0.3</b>	<b>34.6 ± 0.3</b>

Table A.4. Choice of  $K$  for number of centroids.

K	VOC	Context
8	63.8	29.2
16	64.0	29.3
32	64.4	29.4
64	64.3	28.0

Table A.5. Ablation of different SD feature configurations. Removing first and last cross attention *layers*, mid, 1<sup>st</sup> and 2<sup>nd</sup> upsampling *blocks* (all layers in the block) has a negative effect.

1st layer	Mid block	Up-1 block	Up-2 block	Last layer	Context
✓	✓	✓	✓	✓	29.4
	✓	✓	✓	✓	29.4
✓		✓	✓	✓	29.2
✓	✓		✓	✓	27.3
✓	✓	✓		✓	28.9
✓	✓	✓	✓		29.3

Table A.6. Ablation of various configurations for prototypes. We consider average  $\bar{P}$ , instance  $P_n$ , and part  $P_k$  prototypes individually and in various combinations on VOC and Context datasets. Combination of all three types of prototypes shows strongest results.

$\bar{P}$	$P_n$	$P_k$	VOC	Context
✓	✓	✓	64.4	29.4
✓		✓	61.7	29.3
✓	✓		63.5	29.4
		✓	62.5	28.4
		✓	63.7	28.8
	✓		60.0	29.0
✓			62.5	28.4

991 clustering show strong performance when considered indi-  
 992 vidually on VOC. Instance prototypes show strong individual  
 993 performance on Context, as well as in combination with the  
 994 average category prototype. The combination of all three  
 995 types shows the strongest results across the two datasets,  
 996 which is what we adopt in our main set of experiments.

997 We also consider the treatment of prototypes under the

stuff filter. We investigate the impact of not excluding back-  
 ground prototypes for “stuff” classes. In this setting, we  
 measure 29.1 on Context, which is a slight reduction in per-  
 formance. We also investigate the benefit of categorisation  
 into “things” and “stuff” used in the stuff filter component.  
 Instead, we filter all background prototypes using all fore-  
 ground prototypes. In this configuration, we measure 27.6  
 on Context. Both configurations show a reduction from 29.4,  
 measuring using the stuff filter with categorisation in “stuff”  
 and “things”, as used in our main experiments. Finally,  
 we experiment by removing part-level prototypes for “stuff”  
 classes, which also results in a performance drop to 28.0.

**$K$  - number of clusters.** In Tab. A.4, we investigate the  
 sensitivity of the method to the choice of  $K$  for the number  
 of “part” prototypes extracted using  $K$ -means clustering.  
 Although our setting  $K = 32$  obtains slightly better results  
 on Context and VOC, other values result in comparable  
 segmentation performance suggesting that OVDiff is not  
 sensitive to the choice of  $K$  and a range of values is viable.

**SD features.** When using Stable Diffusion as a feature ex-  
 tractor, we consider various combinations of layers/blocks  
 in the UNet architecture. We follow the nomenclature used  
 in the Stable Diffusion implementation where consecutive  
 layers of Unet are organised into *blocks*. There are 3 down-  
 sampling blocks with 2 cross-attention layers each, a mid-  
 block with a single cross-attention, and 3 up-sampling blocks  
 with 3 cross-attention layers each. We report our findings in  
 Tab. A.5. Including the first and last cross-attention layers in  
 the feature extraction process has a small positive impact on  
 segmentation performance, which we attribute to the high  
 feature resolution. We also consider excluding features from  
 the middle block of the network due to small  $8 \times 8$  resolu-  
 tion but observe a small negative impact on performance on  
 the Context dataset. We also investigate whether including  
 the first (Up-1) and the second upsampling (Up-2) blocks  
 are necessary. Without them, the performance drops the  
 most out of the configurations considered. Thus, we use a  
 concatenation of features from the middle, first and second  
 upsampling blocks and the first and last layers in our main  
 experiments.

### A.3. Evaluation without background

One of the notable advantages of our approach is the ability  
 to represent background regions via (negative) prototypes,  
 leading to improved segmentation performance. Neverthe-  
 less, we hereby also evaluate our method under a differ-  
 ent evaluation protocol adopted in prior work, which ex-  
 cludes the *background* class from the evaluation. We note  
 that prior work often requires additional considerations to  
 handle background, such as thresholding. In this setting,  
 however, the background class is *not* predicted, and the  
 set of categories, thus, must be exhaustive. As in practice,  
 this is not the case, and datasets contain unlabelled pixels



Figure A.1. Qualitative comparison on in-the-wild images. OVDiff performs significantly better than prior state-of-the-art, TCL, on wildlife images containing multiple instances, studio photos with simple backgrounds, images containing multiple categories and an image containing a rare instance of a class.

Table A.7. Comparison with methods when background is excluded (decided by ground truth). OVDiff shows comparable performance to prior works despite only relying on pretrained feature extractors. \* result from [9].

Method	VOC-20	Context-59	ADE	Stuff	City
CLIPpy	–	–	13.5	–	–
OVSegmentor	–	–	5.6	–	–
GroupViT*	79.7	23.4	9.2	15.3	11.1
MaskCLIP*	74.9	26.4	9.8	16.4	12.6
ReCo*	57.5	22.3	11.2	14.8	21.1
TCL	77.5	30.3	<b>14.9</b>	19.6	23.1
<b>OVDiff</b>	<b>80.9</b>	<b>32.9</b>	<u>14.1</u>	<b>20.3</b>	<b>23.4</b>

(or simply a background label), such image areas are removed from consideration. Consequently, less emphasis is placed on object boundaries in this setting. As in this setting the background prediction is invalid, we do not consider negative prototypes. For this setting, we benchmark on 5 datasets following [9]: PascalVOC without background, termed VOC-20, Pascal Context without background, termed Context-59, and ADE20k [78], which contains 150 foreground classes, termed ADE-150, COCO-Stuff, termed Stuff, and Cityscapes, termed City. This setting tests the ability of various methods to discriminate between different classes, which for OVDiff is inherent to the choice of feature extractors. Despite this, our method shows competitive performance across wide range of benchmarks Tab. A.7.

#### A.4. Qualitative results

We include additional qualitative results from the benchmark datasets in Fig. A.2. Our method achieves high-quality segmentation across all examples without any post-processing or refinement steps. In Fig. A.3, we show examples of support images sampled for some things, and stuff categories. In Fig. C.5, we show examples of support set images sampled for rare *pikachu* class.

#### B. Broader impact

Semantic segmentation is a component in a vast and diverse spectrum of applications in healthcare, image processing, computer graphics, surveillance and more. As for any foundational technology, applications can be good or bad. OVDiff is similarly widely applicable. It also makes it easier to use semantic segmentation in new applications by leveraging existing and new pre-trained models. This is a bonus for inclusivity, affordability, and, potentially, environmental impact (as it requires no additional training, which is usually computationally intensive); however, these features also mean that it is easier for bad actors to use the technology.

Because OVDiff does not require further training, it is more versatile but also inherits the weaknesses of the components it is built on. For example, it might contain the biases (e.g., gender bias) of its components, in particular Stable Diffusion [53], which is used for generating support images for any given category/description. Thus, it should not be exposed without further filtering and detection of, e.g., NSFW material in the sampled support set. Finally, OVDiff

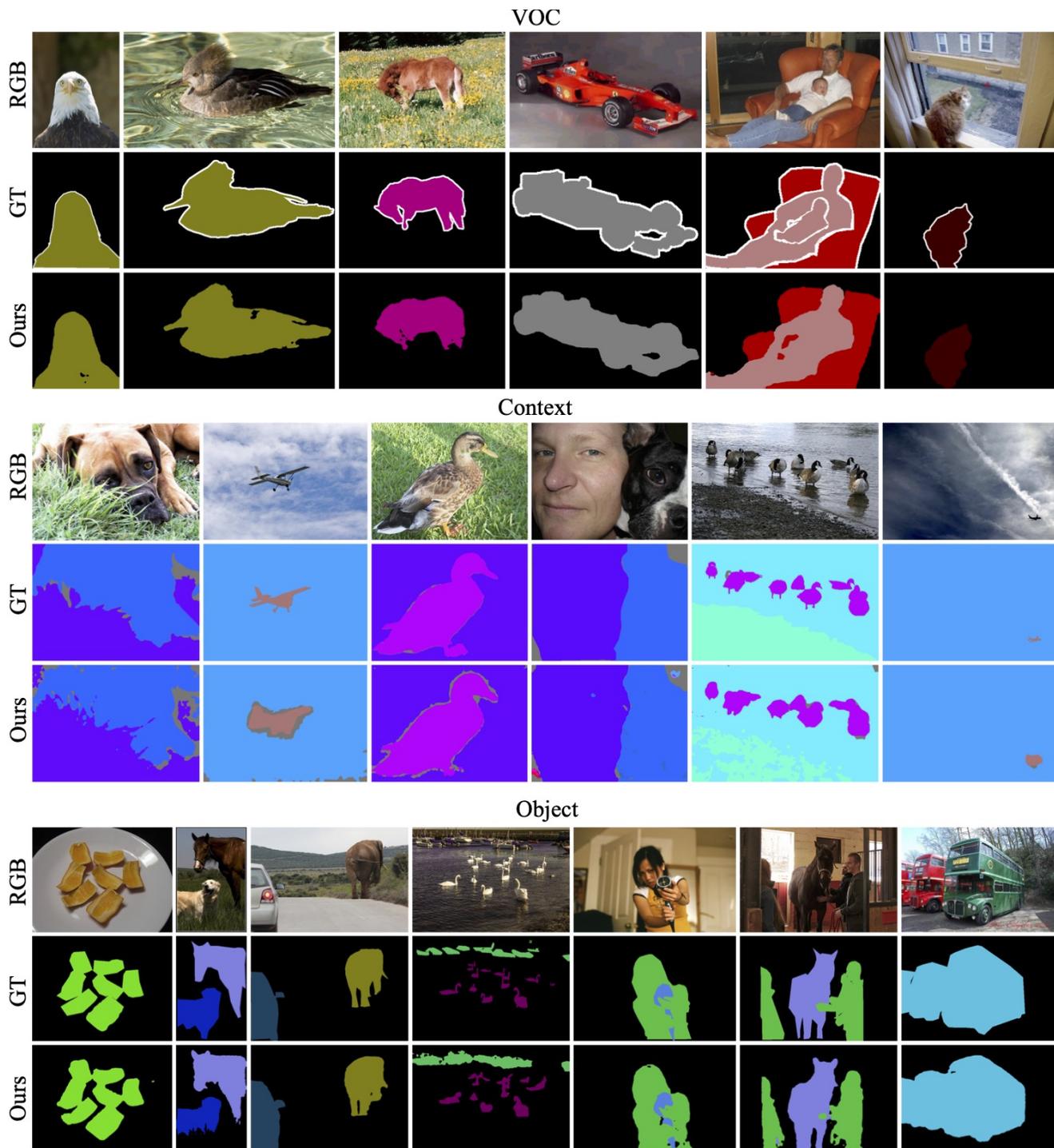


Figure A.2. Additional qualitative results. Images from Pascal VOC (top), Pascal Context (middle), and COCO Object (bottom).

1092 is also bound by the licenses of its components.

### 1093 B.1. Limitations

1094 As OVDiff relies on pretrained components, it inherits some  
1095 of their limitations. OVDiff works with the limited resolution

1096 of feature extractors, due to which it might occasionally  
1097 miss tiny objects. Furthermore, OVDiff cannot segment  
1098 what the generator cannot generate. For example, current  
1099 diffusion models struggle with producing legible text, which  
1100 can make it difficult to segment specific words. Furthermore,

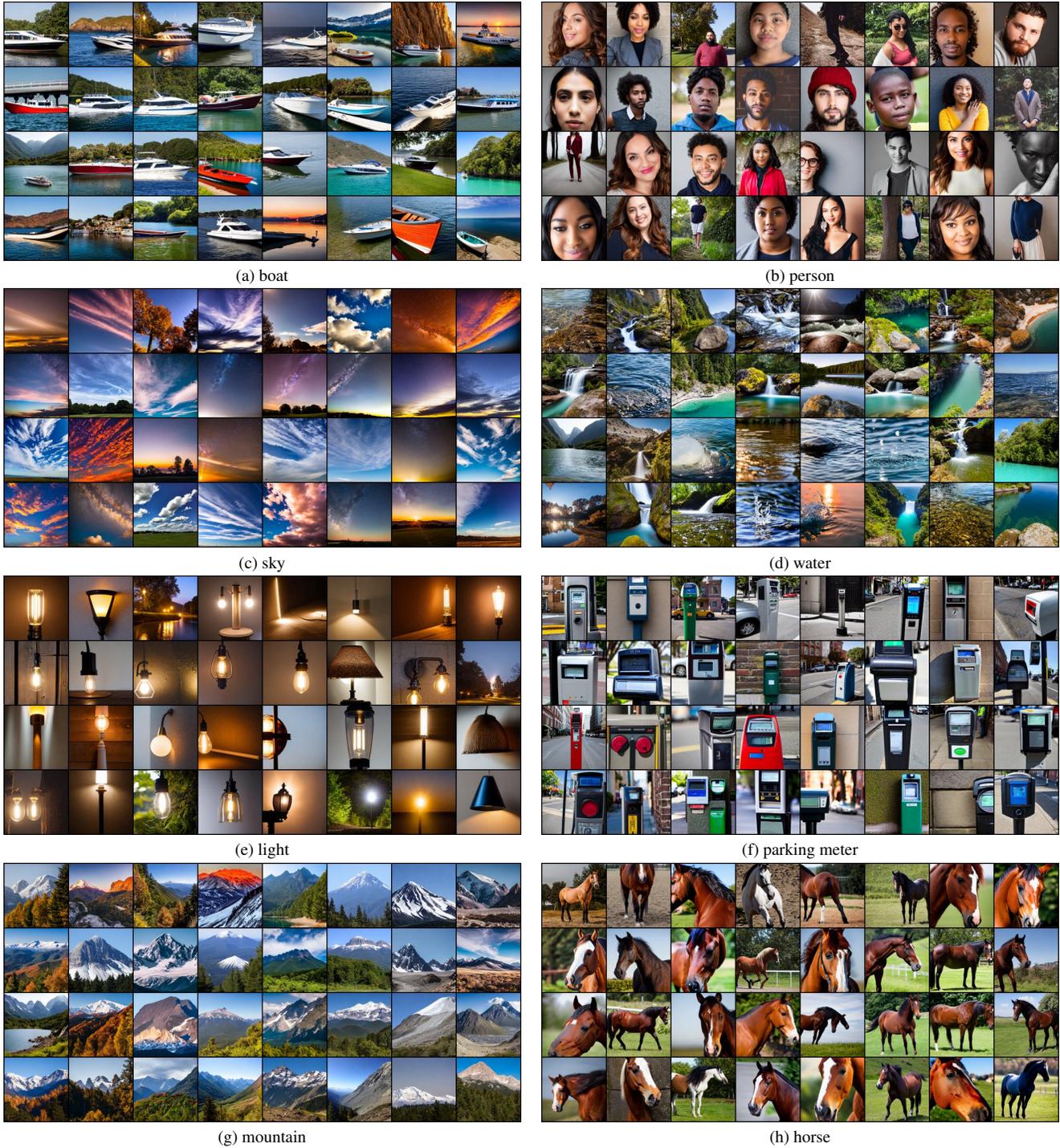


Figure A.3. Images sampled for a support set of some categories.

1101 applications in domains far from the generator’s training data  
1102 (e.g. medical imaging) are unlikely to work out of the box.

### C. OVDiff: Further details

In this section, we provide additional details concerning the implementation of OVDiff. We begin with a brief overview of the attention mechanism and diffusion models central to

1103

1104

1105

1106

extracting features and sampling images. We review different feature extractors used. We specify the hyperparameter setting for all our experiments and provide an overview of the exchange with ChatGPT used to categorise classes into “thing” and “stuff”.

### C.1. Preliminaries

**Attention.** In this work, we make use of pre-trained ViT [16] networks as feature extractors, which repeatedly apply multi-headed attention layers. In an attention layer, input sequences  $X \in \mathbb{R}^{l_x \times d}$  and  $Y \in \mathbb{R}^{l_y \times d}$  are linearly projected to forms *keys*, *queries*, and *values*:  $K = W_k Y$ ,  $Q = W_q X$ ,  $V = W_v X$ . In self-attention,  $X = Y$ . Attention is calculated as  $A = \text{softmax}(\frac{1}{\sqrt{d}} Q K^\top)$ , and softmax is applied along the sequence dimension  $l_y$ . The layer outputs an update  $Z = X + A \cdot V$ . ViTs use multiple heads, replicating the above process in parallel with different projection matrices  $W_k, W_q, W_v$ . In this work, we consider *queries* and *keys* of attention layers as points where useful features that form meaningful inner products can be extracted. As we detail later (Appendix C.2), we use the *keys* from attention layers of ViT feature extractors (DINO/MAE/CLIP), concatenating multiple heads if present.

**Text-to-image diffusion models.** Diffusion models are a class of generative models that form samples starting with noise and gradually denoising it. We focus on latent diffusion models [50] which operate in the latent space of an image VAE [28] forming powerful conditional image generators. During training, an image is encoded into VAE latent space, forming a latent vector  $z_0$ . A noise is injected forming a sample  $z_\tau \sim \mathcal{N}(z_\tau; \sqrt{1 - \alpha_\tau} z_0, \alpha_\tau I)$  for timestep  $\tau \in \{1 \dots T\}$ , where  $\alpha_\tau$  are variance values that define a noise schedule such that the resulting  $z_T$  is approximately unit normal. A conditional UNet [51],  $\epsilon_\theta(z_t, t, c)$ , is trained to predict the injected noise, minimising the mean squared error  $\mathbb{E}_t(\alpha_t \|\epsilon_\theta(z_t, t, c) - z_0\|_2)$  for some caption  $c$  and additional constants  $\alpha_t$ . The network forms new samples by reversing the noise-injecting chain. Starting from  $\hat{z}_T \sim \mathcal{N}(\hat{z}_T; 0, I)$ , one iterates  $\hat{z}_{t-1} = \frac{1}{\sqrt{1 - \alpha_t}}(\hat{z}_t + \alpha_t \epsilon_\theta(\hat{z}_t, t, c)) + \sqrt{\alpha_t} \hat{z}_t$  until  $\hat{z}_0$  is formed and decoded into image space using the VAE decoder. The conditional UNet uses cross-attention layers between image patches and language (CLIP) embeddings to condition on text  $c$  and achieve text-to-image generation.

### C.2. Feature extractors

OVDiff is buildable on top of any pre-trained feature extractor. In our experiments, we have considered several networks as feature extractors with various self-supervised training regimes:

- **DINO** [8] is a self-supervised method that trains networks by exploring alignment between multiple views using an exponential moving average teacher network. We use

the ViT-B/8 model pre-trained on ImageNet<sup>2</sup> and extract features from the *keys* of the last attention layer.

- **MAE** [22] is a self-supervised method that uses masked image inpainting as a learning objective, where a portion of image patches are dropped, and the network seeks to reconstruct the full input. We use the ViT-L/16 model pre-trained on ImageNet at a resolution of 448 [27].<sup>3</sup> The *keys* of the last layer of the *encoder* network are used. No masking is performed.
- **CLIP** [46] is trained using image-text pairs on an internal dataset WIT-400M. We use ViT-B/16 model<sup>4</sup>. We consider two locations to obtain dense features: *keys* from a self-attention layer of the image encoder and *tokens* which are the outputs of transformer layers. We find that *keys* of the second-to-last layer give better performance.
- We also consider **Stable Diffusion**<sup>5</sup> (v1.5) itself as a feature extractor. To that end, we use the *queries* from the cross-attention layers in the UNet denoiser, which correspond to the image modality. Its UNet is organised into three downsampling blocks, a middle block, and three upsampling blocks. We observe that the middle layers have the most semantic content, so we consider the middle block, 1st and 2nd upsampling blocks and aggregate features from all three cross-attention layers in each block. As the features are quite low in resolution, we include the first downsampling cross-attention layer and the last upsampling cross-attention layer as well. The feature maps are bilinearly upsampled to resolution  $64 \times 64$  and concatenated. A noise appropriate for  $\tau = 200$  timesteps is added to the input. For feature extraction, we run SD in *unconditional* mode, supplying an empty string for text caption.



Figure C.4. FG/BG segmentation of classes of *water*, *snow* and *grass*. The foreground is in red, while the background is shown in blue.

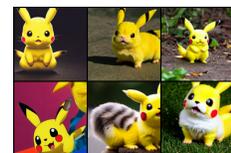


Figure C.5. Example images from the support set of a rare *pikachu* class.

<sup>2</sup>Model and code available at <https://github.com/facebookresearch/dino>.

<sup>3</sup>Model and code from [https://github.com/facebookresearch/long\\_seq\\_mae](https://github.com/facebookresearch/long_seq_mae).

<sup>4</sup>Model and code from <https://github.com/openai/CLIP>.

<sup>5</sup>We use implementation from <https://github.com/huggingface/diffusers>.

1189	<b>C.3. Datasets</b>		
1190	We evaluate on validation splits of PASCAL VOC (VOC),	OVDiff has relatively few hyperparameters and we use the	1241
1191	Pascal Context (Context) and COCO-Object (Object)	same set in all experiments. Unless otherwise specified,	1242
1192	datasets. PASCAL VOC [17, 18] has 21 classes: 20 fore-	$N = 32$ images are sampled using classifier-free guid-	1243
1193	ground plus a background class. For Pascal Context [42],	ance scale [25] of 8.0 and 30 denoising steps. We employ	1244
1194	we use the common variant with 59 foreground classes and	DPM-Solver scheduler [36]. When sampling images for	1245
1195	1 background class. It contains both “things” and “stuff”	the support sets, we also use a negative prompt “ <i>text, low</i>	1246
1196	classes. The COCO-Object is a variant of COCO-Stuff [7]	<i>quality, blurry, cartoon, meme, low resolution, bad, poor,</i>	1247
1197	with 80 “thing” classes and one class for the background.	<i>faded</i> ”. If/when segmenter $\Gamma$ fails to extract any components	1248
1198	Textual class names are used as natural language specifica-	in a sampled image, a fallback of adaptive thresholding of	1249
1199	tions of names. We renamed or specified certain class names	$A_n$ is used, following [34]. During inference, we set $\eta = 10$ ,	1250
1200	to fix errors ( <i>e.g.</i> pottedplant $\rightarrow$ potted plant),	which results in 1024 text prompts processed in parallel, a	1251
1201	resolve ambiguity better ( <i>e.g.</i> mouse $\rightarrow$ computer	choice made mainly due to computational constraints. We	1252
1202	mouse) or change to more common spelling/word ( <i>e.g.</i>	set the thresholds for the “stuff” filter between background	1253
1203	aeroplane $\rightarrow$ airplane), resulting in 14 fixes. We	prototypes for “things” classes and the foreground of “stuff”	1254
1204	experiment and measure the impact of this in Appendix A.1	at 0.85 for all feature extractors. When sampling, a seed	1255
1205	for our and prior work.	is set for each category individually to aid reproducibility.	1256
1206	<b>C.4. Comparative baselines</b>	With our unoptimized implementation, we measure around	1257
1207	We briefly review the prior work in used in our experi-	$110 \pm 10$ s to calculate prototypes (sample images, extract fea-	1258
1208	ments, mainly in Table 1. We consider baselines that do	tures and aggregate) for a single category or $50.2 \pm 2$ s without	1259
1209	not rely on mask annotations and have code and check-	clustering using SD. Using CLIP, we measure $49.2 \pm 0.2$ s	1260
1210	points available or detail their evaluation protocol that	with clustering and $47.7 \pm 0.2$ s without. We note that sam-	1261
1211	matches that used in other prior works [9, 70, 71]. Most	pling time grows linearly: we measure 55s for 16, 110s for	1262
1212	prior work [9, 35, 37, 49, 70, 71] trains image and text	32, and 213s for 64 images per class. The prototype storage	1263
1213	encoders on large image-text datasets with a contrastive	requirements are 0.39MB using CLIP/DINO for each class.	1264
1214	loss. The methods mainly differ in their architecture and	We additionally measure the speed of inference at 0.6s	1265
1215	use of grouping mechanisms to ground image-level text on	per image, which is slightly slower but comparable to 0.2s	1266
1216	regions. ViL-Seg [35] uses online clustering, GroupViT [70]	for TCL and 0.08s for OVSegmentor. We performed infer-	1267
1217	and ViewCo [49] employ group tokens. OVSegmentor [71]	ence measurements using SD on the same machine with a	1268
1218	uses slot-attention and SegCLIP [37] a grouping mecha-	2080Ti GPU using 21 classes and the same resolution/sliding	1269
1219	nism with learnable centers. CLIPPy [48], TCL [9], and	window settings for all methods.	1270
1220	MaskCLIP [79] predict classes for each image patch: [48]	<b>C.6. Interaction with ChatGPT</b>	1271
1221	use max-pooling aggregation, [9] self-masking, and [79]	We interact with ChatGPT to categorise classes into “stuff”	1272
1222	modify CLIP for dense predictions. To assign a background	and “things” for the stuff filter component. Due to input lim-	1273
1223	label [9, 35, 37, 49, 70] use thresholding while [48] uses	its, the categories are processed in blocks. Specifically, we	1274
1224	dataset-specific prompts. CLIP-DIY [68] leverages CLIP	input “ <i>In semantic segmentation, there are “stuff” or “thing”</i>	1275
1225	as a zero-shot classifier and applies it on multiple scales to	<i>classes. Please indicate whether the following class prompts</i>	1276
1226	form a dense segmentation. ReCO [56] is closer in spirit to	<i>should be considered “stuff” or “things”:</i> ”. We show the out-	1277
1227	our approach as it uses a support set for each prompt; this set,	put in Tab. C.8. Note there are several errors in the response,	1278
1228	however, is CLIP-retrieved from curated image collections,	<i>e.g.</i> glass, blanket, and trade name are actually in-	1279
1229	which may not be applicable for any category in-the-wild.	stances of tableware, bedding and signage, respectively, so	1280
1230	We also note that prior work builds on top of similar	should more appropriately be treated as “things”. Similarly,	1281
1231	pre-trained components such as CLIP in [9, 37, 56, 79],	land and sand might be more appropriately handled as	1282
1232	OpenCLIP in [68], DINO + T5/RoBERTa in [48, 71]. We	“stuff”, same as snow and ground. Despite this, We find	1283
1233	additionally make use of StableDiffusion, which is trained	ChatGPT contains sufficient knowledge when prompted with	1284
1234	on a larger dataset (3B, compared to 400M of CLIP or 2B or	“in semantic segmentation”. We have estimated the accuracy	1285
1235	OpenCLIP). OVDiff is, however, fundamentally different to	of ChatGPT in thing/stuff classification using the categories	1286
1236	all prior work, as (a) it generates a support set of synthetic	of COCO-Stuff, which are defined as 80 “things” and 91	1287
1237	images given a class description, and (b) it does not rely on	“stuff” categories. ChatGPT achieves an accuracy rate of	1288
1238	additional training data and further training for learning to	88.9% in this case. We also measure the impact the potential	1289
1239	segment.	errors have on our performance by providing “oracle” an-	1290
		swers on the Context dataset. We measure 29.6 mIoU, which	1291

Table C.8. **Response from interaction with ChatGPT.** We used ChatGPT model to automatically categorise classes in “stuff” or “things”.

airplane:	thing	window:	thing	awning:	thing
bag:	thing	wood:	stuff	streetlight:	thing
bed:	thing	windowpane:	thing	booth:	thing
bedclothes:	stuff	earth:	thing	television receiver:	thing
bench:	thing	painting:	thing	dirt track:	thing
bicycle:	thing	shelf:	thing	apparel:	thing
bird:	thing	house:	thing	pole:	thing
boat:	thing	sea:	thing	land:	thing
book:	thing	mirror:	thing	bannister:	thing
bottle:	thing	rug:	thing	escalator:	thing
building:	thing	field:	thing	ottoman:	thing
bus:	thing	armchair:	thing	buffet:	thing
cabinet:	thing	seat:	thing	poster:	thing
car:	thing	desk:	thing	stage:	thing
cat:	thing	wardrobe:	thing	van:	thing
ceiling:	stuff	lamp:	thing	ship:	thing
chair:	thing	bathtub:	thing	fountain:	thing
cloth:	stuff	railing:	thing	conveyer belt:	thing
computer:	thing	cushion:	thing	canopy:	thing
cow:	thing	base:	thing	washer:	thing
cup:	thing	box:	thing	plaything:	thing
curtain:	stuff	column:	thing	swimming pool:	thing
dog:	thing	signboard:	thing	stool:	thing
door:	thing	chest of drawers:	thing	barrel:	thing
fence:	stuff	counter:	thing	basket:	thing
floor:	stuff	sand:	thing	waterfall:	thing
flower:	thing	sink:	thing	tent:	thing
food:	thing	skyscraper:	thing	minibike:	thing
grass:	stuff	fireplace:	thing	cradle:	thing
ground:	stuff	refrigerator:	thing	oven:	thing
horse:	thing	grandstand:	thing	ball:	thing
keyboard:	thing	path:	thing	step:	stuff
light:	thing	stairs:	thing	tank:	thing
motorbike:	thing	runway:	thing	trade name:	stuff
mountain:	stuff	case:	thing	microwave:	thing
mouse:	thing	pool table:	thing	pot:	thing
person:	thing	pillow:	thing	animal:	thing
plate:	thing	screen door:	thing	lake:	stuff
platform:	stuff	stairway:	thing	dishwasher:	thing
plant:	thing	river:	thing	screen:	thing
road:	stuff	bridge:	thing	blanket:	stuff
rock:	stuff	bookcase:	thing	sculpture:	thing
sheep:	thing	blind:	thing	hood:	thing
shelves:	thing	coffee table:	thing	sconce:	thing
sidewalk:	stuff	toilet:	thing	vase:	thing
sign:	thing	hill:	thing	traffic light:	thing
sky:	stuff	countertop:	thing	tray:	stuff
snow:	stuff	stove:	thing	ashcan:	thing
sofa:	thing	palm:	thing	fan:	thing
table:	thing	kitchen island:	thing	pier:	thing
track:	stuff	swivel chair:	thing	crt screen:	thing
train:	thing	bar:	thing	bulletin board:	thing
tree:	thing	arcade machine:	thing	shower:	thing
truck:	thing	hovel:	thing	radiator:	thing
monitor:	thing	towel:	thing	glass:	stuff
wall:	stuff	tower:	thing	clock:	thing
water:	stuff	chandelier:	thing	flag:	thing

1292 is similar to  $29.7 \pm 0.3$  of using ChatGPT, showing that small  
 1293 errors do not drastically affect the method, however, enable  
 1294 using “stuff” filter component, which improves performance  
 1295 (see Table 3).