

DIFFUSION MODELS FOR OPEN-VOCABULARY SEGMENTATION

Anonymous authors

Paper under double-blind review

ABSTRACT

1 The variety of objects in the real world is unlimited and is thus impossible to
2 capture using models trained on a closed, pre-defined set of categories. Recently,
3 open-vocabulary recognition has garnered significant attention, largely facilitated
4 by advances in large-scale vision-language modelling. In this paper, we present
5 OVDiff, a novel method that leverages the generative properties of text-to-image
6 diffusion models for open-vocabulary segmentation. Specifically, we propose
7 to synthesise support image sets from arbitrary textual categories, creating for
8 each category a set of prototypes representative of both the category itself and
9 its surrounding context (background). Our method relies solely on pre-trained
10 components: segmentation is obtained by simply comparing a target image to the
11 prototypes without further fine-tuning. We show that our method can be used to
12 ground any pre-trained self-supervised feature extractor in natural language and
13 provide explainable predictions by mapping back to regions in the support set. Our
14 approach shows strong performance on a range of open-vocabulary segmentation
15 benchmarks, obtaining a lead of more than 10% over prior work on PASCAL VOC.

16 1 INTRODUCTION

17 Semantic segmentation aims to classify each pixel in an image into a set of categorical labels.
18 Traditionally, this task requires large, densely annotated datasets for training models to predict class
19 assignments at pixel level and relies on the assumption that the set of labels is fixed and predefined.
20 Collecting and annotating such data is not only cumbersome and costly, but it also results in static
21 models that are difficult to extend to new categories.

22 Open-vocabulary semantic segmentation relaxes this restriction by allowing nearly arbitrary free-
23 form text queries as class descriptions. This problem is often approached by extracting image
24 embeddings and matching them to a representation of the text queries. Obtaining these embeddings
25 is challenging as they need to describe the image densely and they must also be compatible with
26 the representation of any possible text query. Prior work addresses this challenge by starting from
27 multi-modal representations (*e.g.*, CLIP (Radford et al., 2021)) to bridge vision and language and
28 further relies on labelled data to fine-tune the representations for the segmentation task. Hence, in line
29 with the zero-shot setting (Bucher et al., 2019), these methods require dense annotations for some
30 known categories, while also extending segmentation to unseen categories by incorporating language.

31 An alternative, which eliminates the need for collecting ad-hoc manual annotations, is to leverage
32 image-text pairs that can be obtained at scale by crawling the Internet. Existing methods (Xu et al.,
33 2022a; Ren et al., 2023; Xu et al., 2023b; Luo et al., 2022; Mukhoti et al., 2022; Cha et al., 2022)
34 observe that large-scale vision-language models such as CLIP have a limited understanding of the
35 positioning of objects within an image and extend these models with additional grouping mechanisms
36 for better localisation using only image-level captions, but no mask supervision. This, however,
37 requires additional contrastive training at scale. Despite yielding promising results, there are some
38 pitfalls to this approach. Firstly, as text might not describe all entities in the image or might mention
39 elements that are not depicted, the training is noisy. Secondly, similar captions may be used to
40 describe a wide range of visual appearances or a similar concept might be described in different ways,
41 though image and language are processed independently. Lastly, most methods resort to heuristics
42 to segment the background (*i.e.*, leave some pixels unlabelled), as it often cannot be described as
43 a textual category. The usual approach is to threshold the similarities to all categories. Finding an

44 appropriate threshold, however, can be challenging and may vary depending on the image, often
45 resulting in imprecise object boundaries. Effectively handling the background remains an open issue.

46 While the field of Computer Vision evolves towards large, pre-trained, general-purpose models, its
47 applications still rely on task-specific approaches, data, and fine-tuning. Thus, in this work, we show
48 that the segmentation problem can be effectively tackled with a combination of frozen “foundation”
49 models without any task-specific adaptation.

50 Specifically, we show that large-scale text-to-image generative models such as StableDiffusion (Rom-
51 bach et al., 2022) open up new avenues for solving this problem, as they are able to bridge the
52 vision-language gap by synthesising data on-the-fly, but also produce latent spaces that are semanti-
53 cally meaningful and well-localised. This also solves a second problem: multi-modal embeddings
54 are difficult to learn and often suffer from ambiguities and differences in detail between modalities.
55 Instead, our approach can use unimodal features for open-vocabulary segmentation, which offers
56 several advantages. Firstly, as text-to-image generators encode a distribution of possible images, this
57 offers a means to deal with intra-class variation and captures the ambiguity in textual descriptions.
58 Secondly, the generative image models encode not only the visual appearance of objects but also
59 provide contextual priors such as backgrounds which can greatly improve the segmentation quality.

60 Given a textual prompt, our method, OVDiff, uses a generative model to produce a support set
61 of visual examples that we then decompose into a set of feature prototypes at different levels of
62 granularity: class, instance, and part prototypes. Prototypes are essentially image features extracted
63 from off-the-shelf unsupervised feature extractors. They can then be used in a simple nearest-
64 neighbour lookup scheme to segment any image. We also propose to leverage the backgrounds from
65 sampled images to encode a set of negative prototypes that enable direct background segmentation.

66 In this work, we present a simple framework that achieves state-of-the-art performance across open-
67 vocabulary segmentation benchmarks. It makes use of several off-the-shelf pre-trained networks
68 and requires no additional data nor fine-tuning. As such, the model can directly benefit from future
69 improvements of its adopted models.

70 2 RELATED WORK

71 **Zero-shot open-vocabulary segmentation.** Open-vocabulary semantic segmentation is a relatively
72 new problem and is typically approached in two different ways. The first line of work poses the
73 problem as a “zero-shot” task, *i.e.*, segmenting unseen classes after training on a set of observed
74 classes with dense annotations. Early approaches (Bucher et al., 2019; Li et al., 2020; Gu et al.,
75 2020; Cheng et al., 2021) explore generative networks to sample features using conditional language
76 embeddings for classes. In (Xian et al., 2019; Li et al., 2021) image encoders are trained to output
77 dense features that can be correlated with word2vec (Mikolov et al., 2013) and CLIP (Radford et al.,
78 2021) text embeddings. Follow-up works (Ghiasi et al., 2022; Liang et al., 2022; Ding et al., 2022;
79 Xu et al., 2022b) approach the problem in two steps, predicting class-agnostic masks and aligning the
80 embeddings of these masks with language. IFSeg (Yun et al., 2023) generates synthetic feature maps
81 by pasting CLIP text embeddings into a known spatial configuration to use as additional supervision.
82 Different from our approach, all these works rely on mask supervision for a set of known classes.

83 The second line of work eliminates the need for mask annotations and instead aims to align image
84 regions with language using only image-text pairs. This is largely enabled by recent advancements in
85 large-scale vision-language models (Radford et al., 2021). Some methods introduce internal grouping
86 mechanisms such as hierarchical grouping (Xu et al., 2022a; Ren et al., 2023), slot-attention (Xu et al.,
87 2023b), or cross-attention to learn cluster centroids (Liu et al., 2022; Luo et al., 2022). Assignment to
88 language queries is performed at group level. An alternative line of work (Zhou et al., 2022; Mukhoti
89 et al., 2022; Cha et al., 2022; Ranasinghe et al., 2022) aims to learn dense features that are better
90 localised when correlated with language embeddings at pixel level. With the exception of (Ranasinghe
91 et al., 2022; Zhou et al., 2022), thresholding is often required to determine the background during
92 inference. Alternatively, Ranasinghe et al. (2022) use a curated list of background prompts.

93 Our method falls into the second category. However, in contrast to prior work, we leverage a
94 generative model to translate language queries to pre-trained image feature extractors without further
95 training. We also segment the background directly, without relying on thresholding or curated list of
96 background prompts. A closely related approach to ours is ReCO (Shin et al., 2022b), where CLIP is

97 used for image retrieval compiling a set of exemplar images from ImageNet for a given language
 98 query, which is then used for co-segmentation. In our method, the shortcoming of an image database
 99 is addressed by synthesising data on-demand. Furthermore, instead of co-segmentation, we leverage
 100 the cross-attention of the generator to extract objects. Instead of similarity of support images, our
 101 method leverages diverse samples and makes use of both foreground and contextual backgrounds.

102 **Diffusion models.** Diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2021)
 103 are a class of generative methods that have seen tremendous success in text-to-image systems such as
 104 DALL-E (Ramesh et al., 2022), Imagen (Saharia et al., 2022), and Stable Diffusion (Rombach et al.,
 105 2022), trained on Internet-scale data such as LAION-5B (Schuhmann et al., 2022). The step-wise
 106 generative process and the language conditioning make pre-trained diffusion models attractive also
 107 for discriminative tasks. They have been recently used in few-shot classification (Zhang et al., 2023),
 108 few-shot segmentation (Baranchuk et al., 2022) and panoptic segmentation (Xu et al., 2023a), and to
 109 generate pairs of images and segmentation masks (Li et al., 2023b). However, these methods rely on
 110 dense manual annotations to associate diffusion features with the desired output.

111 Annotation-free discriminative approaches such as (Li et al., 2023a; Clark & Jaini, 2023) use pre-
 112 trained diffusion models as zero-shot classifiers. DiffuMask (Wu et al., 2023) uses prompt engineering
 113 to synthesise a dataset of “known” and “unseen” categories and trains a closed-set segmenter with
 114 masks obtained from the cross-attention maps of the diffusion model. DiffusionSeg (Ma et al., 2023)
 115 uses DDIM inversion (Song et al., 2021) to obtain feature maps and attention masks of object-centric
 116 images to perform unsupervised object discovery, but relies on ImageNet labels and is not open-
 117 vocabulary. Our approach also leverages the rich semantic information present in diffusion models
 118 for segmentation; unlike these methods, however, it is open-set and does not require further training.

119 **Unsupervised segmentation.** Our work is also related to unsupervised segmentation approaches.
 120 While early works relied on hand-crafted priors (Cheng et al., 2015; Wei et al., 2012; Zhang et al.,
 121 2018; Zeng et al., 2019; Nguyen et al., 2019) later approaches leverage feature extractors such
 122 as DINO (Caron et al., 2021) and perform further analysis of these methods (Wang et al., 2022b;
 123 Melas-Kyriazi et al., 2022a; Siméoni et al., 2021; Siméoni et al., 2022; Hamilton et al., 2022; Shin
 124 et al., 2022a; Wang et al., 2023; 2022a). Some approaches make use of generative methods, usually
 125 GANs, to separate images in foreground and background layers (Bielski & Favaro, 2019; Chen et al.,
 126 2019; Benny & Wolf, 2020; Bielski & Favaro, 2022) or analyse latent structure to induce known
 127 foreground-background changes (Voynov et al., 2021; Melas-Kyriazi et al., 2022b) to synthesise a
 128 training dataset with labels. Largely focused on unsupervised saliency prediction, these methods are
 129 class-agnostic and do not incorporate language.

130 3 METHOD

131 We present OVDiff, a method for open-vocabulary segmentation, *i.e.*, semantic segmentation of any
 132 category described in natural language. To achieve this goal we (1) leverage *text-to-image generative*
 133 *models* to generate a set of images representative of the described category, and (2) use these to ground
 134 off-the-shelf *pretrained feature extractors*. This process does not require further training: it relies
 135 only on pretrained components and does not use additional training data or parameter finetuning.

136 Our goal is to devise an algorithm which, given a new vocabulary of categories $c_i \in \mathcal{C}$ formulated
 137 as natural language queries, can segment any image against it. Let $I \in \mathbb{R}^{H \times W \times 3}$ be an image to
 138 be segmented. Let $\Phi_v : \mathbb{R}^{H \times W \times 3} \mapsto \mathbb{R}^{H' \times W' \times D}$ be an off-the-shelf visual feature extractor and
 139 $\Phi_t : \mathbb{R}^{d_t} \mapsto \mathbb{R}^D$ a text encoder. Assuming that image and text encoders are aligned, one can achieve
 140 zero-shot segmentation by simply computing a similarity function, for example, the cosine similarity
 141 $s(\Phi_v(I), \Phi_t(c_i))$, with $s(x, y) = \frac{x^T y}{\|x\| \|y\|}$, between the encoded image $\Phi_v(I)$ and an encoding of a
 142 class label c_i , which is a simple extension of the zero-shot classification paradigm to dense visual
 143 representations. To meaningfully compare different modalities, image and text features must lie in a
 144 shared representation space, which is typically learned by jointly training Φ_v and Φ_t using image-text
 145 or image-label pairs (Radford et al., 2021).

146 We propose two modifications to this approach. First, we observe that it is better to compare
 147 representations of the *same* modality than across vision and language modalities. We thus replace
 148 $\Phi_t(c_i)$ with a D -dimensional *visual* representation \bar{P} of class c_i , which we refer to as a *prototype*.

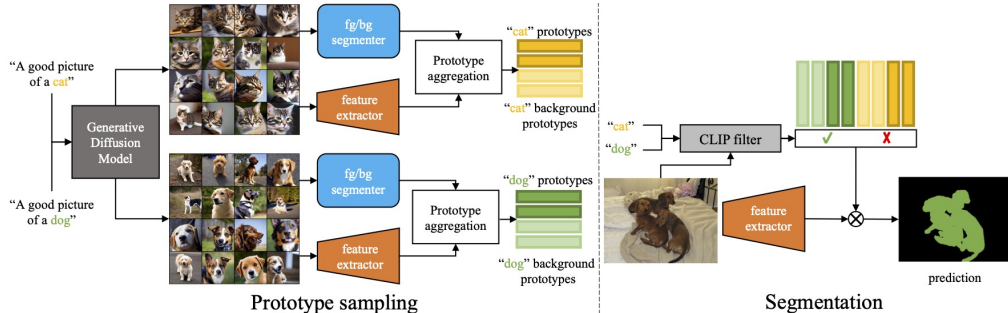


Figure 1: OVDiff overview. Prototype sampling: text queries are used to sample a set of support images which are further processed by a feature extractor and a segmenter forming positive and negative (background) prototypes. Segmentation: image features are compared against prototypes. The CLIP filter removes irrelevant prototypes based on global image contents.

149 In this case, the same feature extractor can be used for both prototypes and target images, thus their
 150 comparison becomes straightforward and does not necessitate further training.

151 Second, we propose utilizing *multiple* prototypes per category instead of a single class embedding.
 152 This enables us to accommodate intra-class variations in appearance and, as we explain later, it also
 153 allows us to exploit contextual priors, which in turn help to effectively segment the background.
 154 Finally, our approach handles the queries c_i independently, allowing for arbitrary changes to the
 155 target vocabulary \mathcal{C} without the need for recomputation.

156 **Support set generation.** To construct a set of prototypes in the visual domain, the first step of
 157 our approach is to sample a support set of images representative of each category c_i . This can be
 158 accomplished by leveraging pretrained text-conditional generative models. Sampling images from
 159 a generative model, as opposed to a curated dataset of real images, aligns well with the goals of
 160 open-vocabulary segmentation as it enables the construction of prototypes for *any* user-specified
 161 category or description, even those for which a manually labelled set may not be readily available
 162 (e.g., $c_i = \text{“donut with chocolate glaze”}$).

163 Specifically, for each query c_i , we define a prompt “A good picture of a $\langle c_i \rangle$ ” and generate
 164 a small batch of N support images $\mathcal{S} = \{S_1, S_2, \dots, S_N \mid S_n \in \mathbb{R}^{h \times w \times 3}\}$ of height h and width w
 165 using Stable Diffusion (Rombach et al., 2022). In its most naïve form, a prototype P could then be
 166 constructed by averaging all features across all images. However, this is unlikely to result in a good
 167 prototype, because not all pixels in the generated image correspond to the category specified by c_i .
 168 To address this issue, we propose to extract the class prototypes as follows.

169 **Class prototypes.** Our approach generates two sets of prototypes, positive and negative, for each
 170 class. Positive prototypes are extracted from image regions that are associated with $\langle c_i \rangle$, while
 171 negative prototypes represent “background” regions. While considering negative or “background”
 172 prototypes is not strictly necessary for segmentation, we found these help to disambiguate objects
 173 from their surroundings by considering contextual priors, which greatly improves performance.

174 Thus, to obtain prototypes, the first step is segmenting the sampled images into foreground (rep-
 175 resenting c_i) and background regions. To identify regions most associated with c_i , we use the
 176 fact that the layout of a generated image is largely dependent on the cross-attention maps of
 177 the diffusion model (Hertz et al., 2022), i.e., pixels attend more strongly to words that describe
 178 them. For a given word or description (in our case c_i), one can generate a set of attribution maps
 179 $\mathcal{A} = \{A_1, A_2, \dots, A_N \mid A_n \in \mathbb{R}^{h \times w}\}$, corresponding to the support set \mathcal{S} , by summing the cross-
 180 attention maps across all layers, heads, and denoising steps of the network (Tang et al., 2022). Yet,
 181 thresholding these attribution maps may not be optimal for foreground/background segmentation, as
 182 they are often coarse or incomplete, and sometimes only parts of objects receive high activation.

183 To address this issue and ensure higher quality masks, we propose to use an unsupervised instance
 184 segmentation method, such as CutLER (Wang et al., 2023). This approach does not use prompts
 185 for object selection and may result in multiple binary object proposals. We denote these as $\mathcal{M}_n =$
 186 $\{M_{nr} \mid M_{nr} \in \{0, 1\}^{h \times w}\}$, where n indexes the support images and r indexes the object masks

187 (including a mask for the background). We thus introduce a promptable extension of CutLER:
 188 for each image, we select from \mathcal{M}_n the mask with the highest (lowest) average attribution as the
 189 foreground (background):

$$M_n^{\text{fg}} = \arg \max_{M \in \mathcal{M}_n} \frac{M^\top A_n}{M^\top M}, \quad M_n^{\text{bg}} = \arg \min_{M \in \mathcal{M}_n} \frac{M^\top A_n}{M^\top M}. \quad (1)$$

190 We can then compute prototypes P_n^g for foreground and background regions ($g \in \{\text{fg}, \text{bg}\}$) as

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

191 where \hat{M}_n^g denotes a resized version of M_n^g that matches the spatial dimensions of $\Phi_v(S_n)$, and
 192 $m_n^g = (\hat{M}_n^g)^\top \hat{M}_n^g$ counts the number of pixels within each mask. In other words, prototypes are
 193 obtained by means of an off-the-shelf pretrained feature extractor and computed as the average feature
 194 within each mask. We refer to these as *instance-level* prototypes, because they are computed from
 195 each image individually and each image in the support set can be viewed as an instance of class c_i .

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

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

203 We consider the union of all these feature prototypes

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

204 and associate all of them with a single category. We note that this process is repeated for each $c_i \in \mathcal{C}$
 205 and we thus refer to \mathcal{P}^{fg} (and \mathcal{P}^{bg}) as $\mathcal{P}_{c_i}^{\text{fg}}$ ($\mathcal{P}_{c_i}^{\text{bg}}$), *i.e.*, as the foreground (background) prototypes of
 206 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
 207 can be dynamically expanded without the need to adapt existing prototypes.

208 **Open-vocabulary segmentation.** To perform segmentation of any target image I given a vocabu-
 209 lary \mathcal{C} , we first extract image features using the same visual encoder Φ_v used for the prototypes.
 210 The vocabulary is expanded with an additional background class $\hat{\mathcal{C}} = \{c_{\text{bg}}\} \cup \mathcal{C}$, for which the
 211 positive (*foreground*) prototype is the union of all *background* prototypes in the vocabulary: $\mathcal{P}_{c_{\text{bg}}}^{\text{fg}}$
 212 $= \bigcup_{c_i \in \mathcal{C}} \mathcal{P}_{c_i}^{\text{bg}}$. Then, a segmentation map can simply be obtained by comparing dense image features
 213 to prototypes using cosine similarity. A class with the highest similarity in its prototype set is chosen:

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

214 **Category pre-filtering.** To limit the impact of spurious correlations that might exist in the feature
 215 space of the visual encoder, we introduce a pre-filtering process for the target vocabulary given image
 216 I . Specifically, we leverage CLIP (Radford et al., 2021) as a strong open-vocabulary classifier but
 217 propose to apply it in a multi-label fashion to constrain the segmentation to the subset of categories
 218 $\mathcal{C}' \subseteq \mathcal{C}$ that appear in the target image. First, we encode the target image and each category using
 219 CLIP. Any categories that do not score higher than $1/|\mathcal{C}'|$ are removed from consideration, that is we
 220 keep the subset $\{P_c^g \mid c' \in \mathcal{C}'\}$, $g \in \{\text{fg}, \text{bg}\}$. If more than η categories are present, then the top- η are
 221 selected. We then form “multi-label” prompts as “ $\langle c_a \rangle$ and $\langle c_b \rangle$ and . . .” where the categories
 222 are selected among the top scoring ones taking into account all 2^η combinations. The best-scoring
 223 multi-label prompt determines the final list of categories to be used in Equation (4).

224 **“Stuff” filtering.** Occasionally, c_i might not describe a countable object category but an identifiable
 225 region in the image, *e.g.*, `sky`, often referred to as a “stuff” class. “Stuff” classes warrant additional
 226 consideration as they might appear as background in images of other categories, *e.g.*, `boat` images
 227 might often contain regions of `water` and `sky`. As a result, the process outlined above might sample

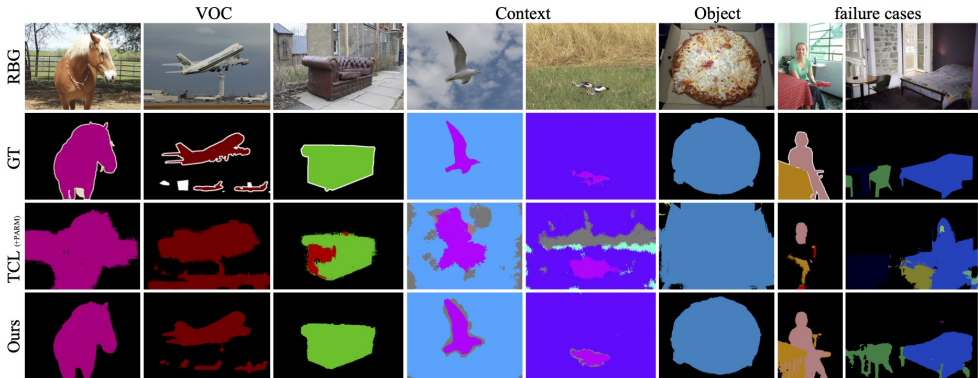


Figure 2: Qualitative results. OVDiff in comparison to TCL (+ PAMR). OVDiff provides more accurate segmentations across a range objects and stuff classes with well defined object boundaries that separate from the background well. Last 2 columns show failure cases. Additional table that appears in the background is segmented. Bed frame legs get misclassified as chairs.

228 background prototypes for one class that coincide with the foreground prototypes of another. To
 229 mitigate this issue, we introduce an additional filtering step to detect and reject such prototypes, when
 230 the full vocabulary, *i.e.*, the set of classes under consideration, is known. First, we only consider
 231 foreground prototypes for “stuff” classes. Additionally, any negative prototypes of “thing” classes
 232 with high cosine similarity with any of the “stuff” class prototypes are simply removed. In our
 233 experiments, we use ChatGPT (OpenAI, 2023) to automatically categorise a set of classes as “thing”
 234 or “stuff”. While this categorisation may contain some errors, this filtering step is still beneficial.

235 4 EXPERIMENTS

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

240 **Datasets and implementation details.** As the approach does not require further training of compo-
 241 nents, we only consider data used for evaluation. Following prior work (Xu et al., 2022a), to assess
 242 the segmentation performance, we report mean Intersection-over-Union (mIoU) on validation splits
 243 of PASCAL VOC (VOC) (Everingham et al., 2012), PASCAL Context (Context) (Mottaghi et al.,
 244 2014) and COCO-Object (Object) (Caesar et al., 2018) datasets, with 20, 59, and 80 foreground
 245 classes, respectively. All datasets have a background class as well. Context also contains both “things”
 246 and “stuff” classes. Similarly to Cha et al. (2022), we employ a sliding window approach. We use
 247 two scales to aid with the limited resolution of off-the-shelf feature extractors with square window
 248 sizes of 448 and 336, and a stride of 224 pixels. We set the size of the support set to $N = 32$. We
 249 detail further specifications of the sampling and other hyper-parameters in Appendix B.5.

250 4.1 GROUNDING FEATURE EXTRACTORS

251 Our method can be used in combination with *any* pretrained visual feature extractor for constructing
 252 prototypes and extracting image features. To verify this quantitatively, we experiment with various
 253 self-supervised ViT feature extractors (Table 2): DINO (Caron et al., 2021), MAE (He et al., 2022),
 254 and CLIP (Radford et al., 2021). We also experiment with SD as a feature extractor. We provide
 255 feature extraction details in Appendix B.2.

256 We find that SD performs the best, though CLIP and DINO also show strong performance based on
 257 our experiments on VOC. MAE shows the weakest performance, which may be attributed to its lack
 258 of semanticity (He et al., 2022); yet it is still competitive with the majority of purposefully trained
 259 networks when employed as part of our approach. We find that taking *keys* of the second to last
 260 layer in CLIP yields better results than using patch tokens (CLIP token). As feature extractors have

Table 1: Open-vocabulary segmentation. Comparison of our approach to the state of the art (under the mIoU metric). Our results are an average of 5 seeds $\pm\sigma$. *results from (Cha et al., 2022).

Method	Support Set	Further Training	VOC	Context	Object
ReCo* (Shin et al., 2022b)	Real	✗	25.1	19.9	15.7
ViL-Seg (Liu et al., 2022)	✗	✓	37.3	18.9	-
MaskCLIP* (Zhou et al., 2022)	✗	✗	38.8	23.6	20.6
TCL (Cha et al., 2022)	✗	✓	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., 2022)	✗	✓	52.6	24.7	26.5
OVSsegmentor (Xu et al., 2023b)	✗	✓	53.8	20.4	25.1
OVDiff (Ours)	Synthetic	✗	66.3 \pm 0.2	29.7 \pm 0.3	34.6 \pm 0.3
TCL (Cha et al., 2022) (+ PAMR)	✗	✓	55.0	30.4	31.6
OVDiff (+ PAMR)	Synthetic	✗	68.4 \pm 0.2	31.2 \pm 0.4	36.2 \pm 0.4

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

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

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

Configuration	VOC	Δ	Context	Δ
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 \bar{P}	62.5	-1.9	28.4	-1.0

261 different training objectives, we hypothesise that their feature spaces might be complementary, thus
 262 we also consider an ensemble approach. In this case, the cosine distances formed between features of
 263 different extractors and respective prototypes are simply averaged. The combination of SD, DINO,
 264 and CLIP performs the best. We adopt this formulation for the main set of experiments.

265 4.2 COMPARISON TO EXISTING METHODS

266 In Table 1, we compare our method with prior work on three datasets: VOC, Context, Object. We
 267 include brief overview of the methods in Appendix B.4. We find that our method compares favourably,
 268 outperforming other methods in all settings. In particular, results on VOC show the largest margin,
 269 with more than 10% improvement over prior work. We hypothesise that this setting is particularly
 270 favourable to our method as it contains scenes where classes take up larger areas of the image.

271 In the same table, we also combine our method with PAMR (Araslanov & Roth, 2020), the post-
 272 processing approach employed by TCL. We find that it improves results for our method though
 273 improvements are less drastic since our method already yields better segmentation and boundaries.

274 Qualitative results are shown in Fig. 2. This figure highlights a key benefit of our approach: the
 275 ability to exploit contextual priors through the use of background prototypes, which in turn allows for
 276 the directly assignment of pixels to a background class. This improves segmentation quality because
 277 it makes it easier to differentiate objects from the background and to delineate their boundaries. In
 278 comparison, TCL predicts very coarse semantic masks and a larger amount of noise.

279 4.3 ABLATIONS

280 Next, we ablate the components of OVDiff on VOC and Context datasets. For these experiments,
 281 only SD is employed as a feature extractor. We remove individual components and measure the

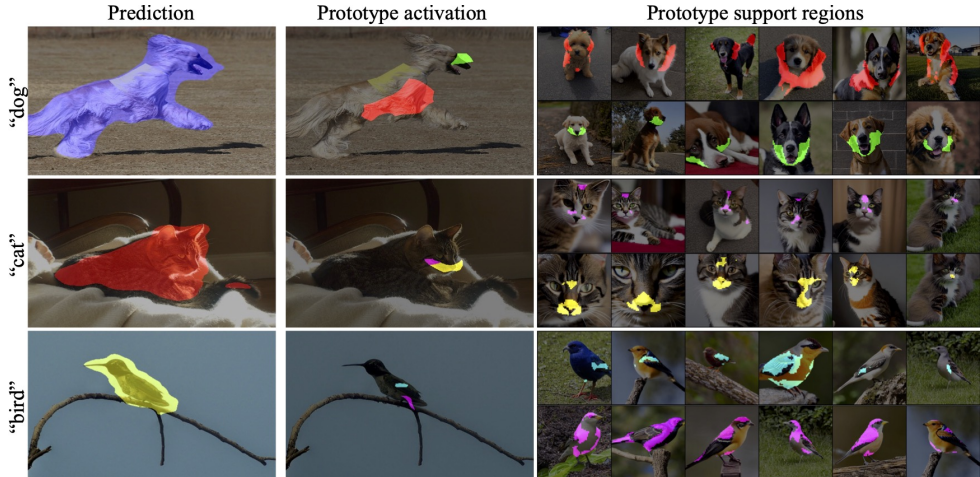


Figure 3: Analysis of the segmentation output by linking regions to samples in the support set. Left: our results for different classes. Middle: select color-coded regions “activated” by different prototypes for the class. Right: regions in the support set images corresponding to these (part-level) prototypes.

282 change in segmentation performance, summarising the results in Table 3. Our first observation
 283 is that background prototypes have a major impact on performance. When removing them from
 284 consideration, we instead threshold the similarity scores of the images with the foreground proto-
 285 types (set to 0.72, determined via grid search); in this case, the performance drops significantly,
 286 which again highlights the importance of leveraging contextual priors. On Context, the impact is
 287 less significant, likely due to the fact that the dataset contains “stuff” categories. Removing the
 288 *instance-* and *part-level* prototypes also negatively affects performance. Additionally, removing
 289 the category pre-filtering has a major impact. We hypothesize that this introduces spurious cor-
 290 relations between prototypes of different classes. On Context, “stuff” filtering is also important.
 291 Next, we evaluate the importance of using CutLER to obtain foreground/background prototypes.
 292 Instead of a segmentation method, one can threshold the attribution
 293 maps obtained directly through the diffusion process. However, we
 294 find that this slightly reduces performance. Overall, background
 295 prototypes and pre-filtering contribute the most.

296 Finally, we measure the effect of varying the size of the support
 297 set N in Fig. 4. We find that even at a low number of samples for
 298 each query, our method already shows strong performance. With
 299 increasing the number of samples, the performance improves, satu-
 300 rating at around 32, which is what we use in our main experiments.
 301 We leave additional ablations for Appendix C.2.

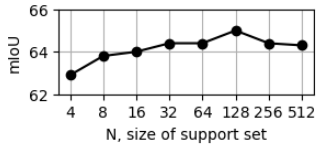


Figure 4: PascalVOC results with increasing support size N .

302 4.4 EXPLAINING SEGMENTATIONS

303 We inspect how our method segments certain regions by considering which prototype from $\mathcal{P}_c^{\text{fg}}$ was
 304 used to assign a class c to a pixel. Prototypes have a mapping to regions in the support set from
 305 where they were aggregated, *e.g.*, instances prototypes are associated with foreground masks M_n^{fg} and part
 306 prototypes with centroids/clusters.

307 By following these mappings, a set of support image regions can be retrieved for each segmentation
 308 decision providing a degree of explainability. Fig. 3 illustrates this for examples of dog, cat, and
 309 bird classes. For visualisation purposes, select prototypes and corresponding regions are shown.
 310 On the left, we show the full segmentation result of each image. In the middle, we select regions that
 311 correlated best with certain prototypes of the class. On the right, we retrieve images from the support
 312 set and highlight where each prototype emerged. We find that meaningful part segmentation merges
 313 due to clustering the support image features, and similar regions are segmented by corresponding
 314 prototypes. Though sometimes region covered in the input image will not fully align with whole

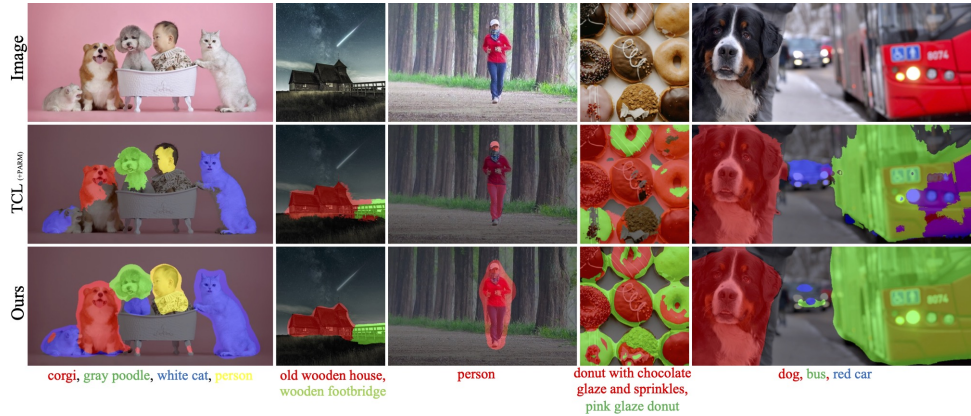


Figure 5: Qualitative comparison on in-the-wild images with TCL, which struggles with object boundaries, missing parts of objects, or including surroundings. Our method has more appropriate boundaries but does produce small halo effect around objects due to upscaling of feature extractors.

315 prototype (e.g. cat’s face around the eyes or lower belly/tail of bird). This shows how each
 316 segmentation produced by OVDiff is explained by precise regions in a small set of support images.

317 4.5 IN-THE-WILD

318 In Fig. 5, we investigate OVDiff on in-the-wild images containing simple and complex backgrounds.
 319 We compare with TCL+PAMR. In the first three images, both methods correctly detect the objects
 320 identified by the queries. TCL however misses parts of the objects, such as most of the person, and
 321 parts of animal bodies. The distinction between the house and the bridge in the second image is also
 322 better with OVDiff. We note that our segmentations sometimes have halos around objects. This is
 323 caused by the upscaling of the low-resolution feature extractor (SD in this case). The last two images
 324 contain difficult scenarios where both approaches struggle. The fourth image only contains similar
 325 objects of the same type. Both methods incorrectly identify plain donuts as either of the specified
 326 queries. OVDiff however correctly identifies chocolate donuts with varied sprinkles and separates
 327 all donuts from the background. In the final picture, the query “red car” is added, although no such
 328 object is present. The extra query causes TCL to incorrectly identify parts of the red bus as a car.
 329 Both methods incorrectly segment the gray car in the distance. However, overall, our method is more
 330 robust and delineates objects better despite the lack of specialized training or post-processing.

331 4.6 LIMITATIONS

332 As OVDiff relies on pretrained components, it inherits some of their limitations. OVDiff works with
 333 the limited resolution of feature extractors, due to which it might miss tiny objects. While this can
 334 be partially mitigated with a sliding window, employing high-resolution feature extractors is one
 335 direction of future improvements. Furthermore, OVDiff cannot segment what the generator cannot
 336 generate. For example, current diffusion models struggle with producing legible text. Finally, one
 337 limitation comes from the computational overhead of sampling support images. We observe that in
 338 practice often whole image collections are segmented by the same set of queries, amortising this cost.

339 5 CONCLUSION

340 We introduce OVDiff, an open-vocabulary segmentation method that operates in two stages. First,
 341 given queries, support images are sampled and their features are extracted to create class prototypes.
 342 These prototypes are then compared to features from an inference image. This approach offers
 343 multiple advantages without task-specific adaptation of its pre-trained components: diverse prototypes
 344 accommodating various visual appearances and negative prototypes for background localisation.
 345 OVDiff outperforms prior work on benchmarks, exhibiting fewer errors, effectively separating objects
 346 from background, and providing explainability through segmentation mapping to support set regions.

347 REFERENCES

- 348 Nikita Araslanov and Stefan Roth. Single-stage semantic segmentation from image labels. In *Proceedings of the*
349 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4253–4262, 2020. 7
- 350 Dmitry Baranchuk, Andrey Voynov, Ivan Rubachev, Valentin Khrukov, and Artem Babenko. Label-efficient
351 semantic segmentation with diffusion models. In *International Conference on Learning Representations*,
352 2022. 3
- 353 Yaniv Benny and Lior Wolf. Onegan: Simultaneous unsupervised learning of conditional image generation,
354 foreground segmentation, and fine-grained clustering. In *European Conference on Computer Vision*, pp.
355 514–530. Springer, 2020. 3
- 356 Adam Bielski and Paolo Favaro. Emergence of object segmentation in perturbed generative models. *Advances*
357 *in Neural Information Processing Systems*, 32, 2019. 3
- 358 Adam Bielski and Paolo Favaro. Move: Unsupervised movable object segmentation and detection. In *Advances*
359 *in Neural Information Processing Systems*, 2022. 3
- 360 Maxime Bucher, Tuan-Hung Vu, Matthieu Cord, and Patrick Pérez. Zero-shot semantic segmentation. *Advances*
361 *in Neural Information Processing Systems*, 32, 2019. 1, 2
- 362 Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *Computer*
363 *vision and pattern recognition (CVPR), 2018 IEEE conference on*. IEEE, 2018. 6, 15
- 364 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin.
365 Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international*
366 *conference on computer vision*, pp. 9650–9660, 2021. 3, 6, 15
- 367 Junbum Cha, Jonghwan Mun, and Byungseok Roh. Learning to generate text-grounded mask for open-world
368 semantic segmentation from only image-text pairs. *arXiv preprint arXiv:2212.00785*, 2022. 1, 2, 6, 7, 15, 16,
369 18, 19
- 370 Mickaël Chen, Thierry Artières, and Ludovic Denoyer. Unsupervised object segmentation by redrawing.
371 *Advances in neural information processing systems*, 32, 2019. 3
- 372 Jiaxin Cheng, Soumyaroop Nandi, Prem Natarajan, and Wael Abd-Almageed. Sign: Spatial-information
373 incorporated generative network for generalized zero-shot semantic segmentation. In *Proceedings of the*
374 *IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 9556–9566, October 2021. 2
- 375 Ming-Ming Cheng, Niloy J. Mitra, Xiaolei Huang, Philip H. S. Torr, and Shi-Min Hu. Global contrast based
376 salient region detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(3):569–582,
377 2015. 3
- 378 Kevin Clark and Priyank Jaini. Text-to-image diffusion models are zero-shot classifiers. *arXiv preprint*
379 *arXiv:2303.15233*, 2023. 3
- 380 Jian Ding, Nan Xue, Gui-Song Xia, and Dengxin Dai. Decoupling zero-shot semantic segmentation. In
381 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11583–11592,
382 2022. 2
- 383 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,
384 Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words:
385 Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
386 14
- 387 M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes
388 (voc) challenge. *International Journal of Computer Vision*, 88(2):303–338, June 2010. 15
- 389 M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PAS-
390 CAL Visual Object Classes Challenge 2012 (VOC2012) Results. [http://www.pascal-](http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html)
391 [network.org/challenges/VOC/voc2012/workshop/index.html](http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html), 2012. 6, 15
- 392 Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Scaling open-vocabulary image segmentation with
393 image-level labels. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October*
394 *23–27, 2022, Proceedings, Part XXXVI*, pp. 540–557. Springer, 2022. 2
- 395 Zhangxuan Gu, Siyuan Zhou, Li Niu, Zihan Zhao, and Liqing Zhang. Context-aware feature generation for
396 zero-shot semantic segmentation. In *Proceedings of the 28th ACM International Conference on Multimedia*,
397 pp. 1921–1929, 2020. 2

- 398 Mark Hamilton, Zhoutong Zhang, Bharath Hariharan, Noah Snavely, and William T Freeman. Unsupervised
399 semantic segmentation by distilling feature correspondences. In *International Conference on Learning*
400 *Representations*, 2022. 3
- 401 Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE*
402 *international conference on computer vision*, pp. 2961–2969, 2017. 15
- 403 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders
404 are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
405 *Recognition*, pp. 16000–16009, 2022. 6
- 406 Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt
407 image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 4
- 408 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep*
409 *Generative Models and Downstream Applications*. 16
- 410 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural*
411 *Information Processing Systems*, 33:6840–6851, 2020. 3
- 412 Ronghang Hu, Shoubhik Debnath, Saining Xie, and Xinlei Chen. Exploring long-sequence masked autoencoders.
413 *arXiv preprint arXiv:2210.07224*, 2022. 15
- 414 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. 2014. 14
- 415 Alexander C Li, Mihir Prabhudesai, Shivam Duggal, Ellis Brown, and Deepak Pathak. Your diffusion model is
416 secretly a zero-shot classifier. *arXiv preprint arXiv:2303.16203*, 2023a. 3
- 417 Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and Rene Ranftl. Language-driven semantic
418 segmentation. In *International Conference on Learning Representations*, 2021. 2
- 419 Peike Li, Yunchao Wei, and Yi Yang. Consistent structural relation learning for zero-shot segmentation. *Advances*
420 *in Neural Information Processing Systems*, 33:10317–10327, 2020. 2
- 421 Ziyi Li, Qinye Zhou, Xiaoyun Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Guiding text-to-image diffusion
422 model towards grounded generation. *arXiv:2301.05221*, 2023b. 3
- 423 Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda,
424 and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. *arXiv preprint*
425 *arXiv:2210.04150*, 2022. 2
- 426 Quande Liu, Youpeng Wen, Jianhua Han, Chunjing Xu, Hang Xu, and Xiaodan Liang. Open-world semantic
427 segmentation via contrasting and clustering vision-language embedding. In *Computer Vision–ECCV 2022:*
428 *17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XX*, pp. 275–292.
429 Springer, 2022. 2, 7, 15, 16
- 430 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for
431 guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022. 16
- 432 Huaishao Luo, Junwei Bao, Youzheng Wu, Xiaodong He, and Tianrui Li. SegCLIP: Patch aggregation with
433 learnable centers for open-vocabulary semantic segmentation. *arXiv preprint arXiv:2211.14813*, 2022. 1, 2,
434 7, 15, 16
- 435 Chaofan Ma, Yuhuan Yang, Chen Ju, Fei Zhang, Jinxiang Liu, Yu Wang, Ya Zhang, and Yanfeng Wang.
436 Diffusionseg: Adapting diffusion towards unsupervised object discovery. *arXiv preprint arXiv:2303.09813*,
437 2023. 3
- 438 Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. Deep spectral methods: A surprisingly
439 strong baseline for unsupervised semantic segmentation and localization. In *Proceedings of the IEEE/CVF*
440 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8364–8375, June 2022a. 3
- 441 Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. Finding an unsupervised image
442 segmenter in each of your deep generative models. In *International Conference on Learning Representations*,
443 2022b. 3
- 444 Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words
445 and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013. 2

- 446 Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun,
447 and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *Proceedings*
448 *of the IEEE conference on computer vision and pattern recognition*, pp. 891–898, 2014. 6, 15
- 449 Jishnu Mukhoti, Tsung-Yu Lin, Omid Poursaeed, Rui Wang, Ashish Shah, Philip HS Torr, and Ser-Nam
450 Lim. Open vocabulary semantic segmentation with patch aligned contrastive learning. *arXiv preprint*
451 *arXiv:2212.04994*, 2022. 1, 2
- 452 Tam Nguyen, Maximilian Dax, Chaithanya Kumar Mummadi, Nhung Ngo, Thi Hoai Phuong Nguyen, Zhongyu
453 Lou, and Thomas Brox. Deepusps: Deep robust unsupervised saliency prediction via self-supervision. In
454 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural*
455 *Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. 3
- 456 OpenAI. Introducing chatgpt. <https://openai.com/blog/chatgpt>, 2023. 6
- 457 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
458 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language
459 supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021. 1, 2, 3, 5, 6, 15
- 460 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional
461 image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 3
- 462 Kanchana Ranasinghe, Brandon McKinzie, Sachin Ravi, Yinfei Yang, Alexander Toshev, and Jonathon Shlens.
463 Perceptual grouping in vision-language models. *arXiv preprint arXiv:2210.09996*, 2022. 2, 7, 15, 16
- 464 Pengzhen Ren, Changlin Li, Hang Xu, Yi Zhu, Guangrun Wang, Jianzhuang Liu, Xiaojun Chang, and Xiaodan
465 Liang. Viewco: Discovering text-supervised segmentation masks via multi-view semantic consistency. *arXiv*
466 *preprint arXiv:2302.10307*, 2023. 1, 2, 7, 15, 16
- 467 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image
468 synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
469 *Pattern Recognition*, pp. 10684–10695, 2022. 2, 3, 4, 14
- 470 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image
471 segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th*
472 *International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pp. 234–241.
473 Springer, 2015. 14
- 474 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour,
475 Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion
476 models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–
477 36494, 2022. 3
- 478 Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating
479 inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer*
480 *Vision and Pattern Recognition (CVPR)*, pp. 22522–22531, June 2023. 14
- 481 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo
482 Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for
483 training next generation image-text models. *arXiv preprint arXiv:2210.08402*, 2022. 3
- 484 Gyungin Shin, Samuel Albanie, and Weidi Xie. Unsupervised salient object detection with spectral cluster
485 voting. In *CVPRW*, 2022a. 3
- 486 Gyungin Shin, Weidi Xie, and Samuel Albanie. Reco: Retrieve and co-segment for zero-shot transfer. In
487 *Advances in Neural Information Processing Systems (NeurIPS)*, 2022b. 2, 7, 16
- 488 Oriane Siméoni, Gilles Puy, Huy V. Vo, Simon Roburin, Spyros Gidaris, Andrei Bursuc, Patrick Pérez, Renaud
489 Marlet, and Jean Ponce. Localizing objects with self-supervised transformers and no labels. November 2021.
490 3
- 491 Oriane Siméoni, Chloé Sekkat, Gilles Puy, Antonin Vobecky, Éloi Zablocki, and Patrick Pérez. Unsupervised
492 object localization: Observing the background to discover objects. *arXiv preprint arXiv:2212.07834*, 2022. 3
- 493 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning
494 using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pp. 2256–2265.
495 PMLR, 2015. 3

- 496 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole.
497 Score-based generative modeling through stochastic differential equations. In *International Conference on*
498 *Learning Representations*, 2021. 3
- 499 Raphael Tang, Akshat Pandey, Zhiying Jiang, Gefei Yang, Karun Kumar, Jimmy Lin, and Ferhan Ture. What the
500 daam: Interpreting stable diffusion using cross attention. *arXiv preprint arXiv:2210.04885*, 2022. 4
- 501 Andrey Voynov, Stanislav Morozov, and Artem Babenko. Object segmentation without labels with large-scale
502 generative models. In *International Conference on Machine Learning*, pp. 10596–10606. PMLR, 2021. 3
- 503 Xinlong Wang, Zhiding Yu, Shalini De Mello, Jan Kautz, Anima Anandkumar, Chunhua Shen, and Jose M
504 Alvarez. Freesolo: Learning to segment objects without annotations. In *Proceedings of the IEEE/CVF*
505 *Conference on Computer Vision and Pattern Recognition*, pp. 14176–14186, 2022a. 3
- 506 Xudong Wang, Rohit Girdhar, Stella X Yu, and Ishan Misra. Cut and learn for unsupervised object detection and
507 instance segmentation. *arXiv preprint arXiv:2301.11320*, 2023. 3, 4, 16
- 508 Yangtao Wang, Xi Shen, Shell Xu Hu, Yuan Yuan, James L. Crowley, and Dominique Vaufreydaz. Self-
509 supervised transformers for unsupervised object discovery using normalized cut. In *Proceedings of the*
510 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14543–14553, June 2022b.
511 3
- 512 Yichen Wei, Fang Wen, Wangjiang Zhu, and Jian Sun. Geodesic saliency using background priors. In *ECCV*,
513 2012. 3
- 514 Weijia Wu, Yuzhong Zhao, Mike Zheng Shou, Hong Zhou, and Chunhua Shen. Diffumask: Synthesizing
515 images with pixel-level annotations for semantic segmentation using diffusion models. *arXiv preprint*
516 *arXiv:2303.11681*, 2023. 3
- 517 Yongqin Xian, Subhabrata Choudhury, Yang He, Bernt Schiele, and Zeynep Akata. Semantic projection network
518 for zero-and few-label semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer*
519 *Vision and Pattern Recognition*, pp. 8256–8265, 2019. 2
- 520 Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit:
521 Semantic segmentation emerges from text supervision. In *Proceedings of the IEEE/CVF Conference on*
522 *Computer Vision and Pattern Recognition*, pp. 18134–18144, 2022a. 1, 2, 6, 7, 15, 16
- 523 Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong Wang, and Shalini De Mello. Open-vocabulary
524 panoptic segmentation with text-to-image diffusion models. *arXiv preprint arXiv:2303.04803*, 2023a. 3
- 525 Jilan Xu, Junlin Hou, Yuejie Zhang, Rui Feng, Yi Wang, Yu Qiao, and Weidi Xie. Learning open-vocabulary
526 semantic segmentation models from natural language supervision. *arXiv preprint arXiv:2301.09121*, 2023b.
527 1, 2, 7, 15, 16, 18
- 528 Mengde Xu, Zheng Zhang, Fangyun Wei, Yutong Lin, Yue Cao, Han Hu, and Xiang Bai. A simple baseline for
529 open-vocabulary semantic segmentation with pre-trained vision-language model. In *European Conference on*
530 *Computer Vision*, pp. 736–753, 2022b. 2
- 531 Sukmin Yun, Seong Hyeon Park, Paul Hongsuck Seo, and Jinwoo Shin. Ifseg: Image-free semantic segmentation
532 via vision-language model. *arXiv preprint arXiv:2303.14396*, 2023. 2
- 533 Yu Zeng, Yunzhi Zhuge, Huchuan Lu, Lihe Zhang, Mingyang Qian, and Yizhou Yu. Multi-source weak
534 supervision for saliency detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 3
- 535 Jing Zhang, T. Zhang, Yuchao Dai, Mehrtash Harandi, and Richard I. Hartley. Deep unsupervised saliency
536 detection: A multiple noisy labeling perspective. *2018 IEEE/CVF Conference on Computer Vision and*
537 *Pattern Recognition*, pp. 9029–9038, 2018. 3
- 538 Renrui Zhang, Xiangfei Hu, Bohao Li, Siyuan Huang, Hanqiu Deng, Hongsheng Li, Yu Qiao, and Peng Gao.
539 Prompt, generate, then cache: Cascade of foundation models makes strong few-shot learners. *arXiv preprint*
540 *arXiv:2303.02151*, 2023. 3
- 541 Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through
542 ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
543 633–641, 2017. 20
- 544 Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from clip. In *Computer Vision–ECCV*
545 *2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVIII*, pp.
546 696–712. Springer, 2022. 2, 7, 15, 16

547 SUPPLEMENTARY MATERIAL

548 In this supplementary material, we consider the broader impacts of our work (Appendix A), provide
 549 additional details concerning the implementation (Appendix B), and conclude with additional results
 550 (Appendix C).

551 A BROADER IMPACT

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

559 Because OVDiff does not require further training, it is more versatile, but also, inherits the weaknesses
 560 of the components it is built on. For example, it might contain the biases (e.g., gender bias) of its
 561 components, in particular Stable Diffusion (Schramowski et al., 2023), which is used for generating
 562 support images for any given category/description. Thus it should not be exposed without further
 563 filtering and detection of, e.g., NSFW material in the sampled support set. Finally, OVDiff is also
 564 bound by the licenses of its components.

565 B OVDIFF: FURTHER DETAILS

566 In this section, we provide additional details concerning the implementation of OVDiff. We begin
 567 with a brief overview of the attention mechanism and diffusion models central to extracting features
 568 and sampling images. We review different feature extractors used. We specify the hyperparameter
 569 setting for all our experiments and provide an overview of the exchange with ChatGPT used to
 570 categorise classes into “thing” and “stuff”.

571 B.1 PRELIMINARIES

572 **Attention.** In this work, we make use of pre-trained ViT (Dosovitskiy et al., 2021) networks as
 573 feature extractors, which repeatedly apply multi-headed attention layers. In an attention layer, input
 574 sequences $X \in \mathbb{R}^{l_x \times d}$ and $Y \in \mathbb{R}^{l_y \times d}$ are linearly project to forms *keys*, *queries*, and *values*:
 575 $K = W_k Y$, $Q = W_q X$, $V = W_v X$. In self-attention, $X = Y$. Attention is calculated as
 576 $A = \text{softmax}(\frac{1}{\sqrt{d}} Q K^T)$, and softmax is applied along the sequence dimension l_y . The layer outputs
 577 an update $Z = X + A \cdot V$. ViTs use multiple heads, replicating the above process in parallel with
 578 different projection matrices W_k, W_q, W_v . In this work, we consider *queries* and *keys* of attention
 579 layers as points where useful features that form meaningful inner-products can be extracted. As
 580 we detail later (Appendix B.2), we use the *keys* from attention layers of ViT feature extractors
 581 (DINO/MAE/CLIP), concatenating multiple heads if present.

582 **Text-to-image diffusion models.** Diffusion models are a class of generative models that form
 583 samples starting with noise and gradually denoising it. We focus on latent diffusion models (Rombach
 584 et al., 2022) which operate in the latent space of an image VAE (Kingma & Welling, 2014) forming
 585 powerful conditional image generators. During training, an image is encoded into VAE latent space
 586 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
 587 timestep $\tau \in \{1 \dots T\}$, where α_τ are variance values that define a noise schedule such that the
 588 resulting z_T is approximately unit normal. A conditional UNet (Ronneberger et al., 2015), $\epsilon_\theta(z_t, t, c)$,
 589 is trained to predict the injected noise, minimising the mean squared error $\mathbb{E}_t (\alpha_t \|\epsilon_\theta(z_t, t, c) - z_0\|_2)$
 590 for some caption c and additional constants a_t . The network forms new samples by reversing the noise-
 591 injecting chain. Starting from $\hat{z}_T \sim \mathcal{N}(\hat{z}_T; 0, I)$, one iterates $\hat{z}_{t-1} = \frac{1}{\sqrt{1 - \alpha_t}} (\hat{z}_t + \alpha_t \epsilon_\theta(\hat{z}_t, t, c)) +$
 592 $\sqrt{\alpha_t} \hat{z}_t$ until \hat{z}_0 is formed and decoded into image space using the VAE decoder. The conditional UNet
 593 uses cross-attention layers between image patches and language (CLIP) embeddings to condition on
 594 text c and achieve text-to-image generation.

595 B.2 FEATURE EXTRACTORS

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

- 598 • **DINO** (Caron et al., 2021) is a self-supervised method that trains networks by exploring
599 alignment between multiple views using an exponential moving average teacher network.
600 We use the ViT-B/8 model pre-trained on ImageNet¹ and extract features from the *keys* of
601 the last attention layer.
- 602 • **MAE** (He et al., 2017) is a self-supervised method that uses masked image inpainting as
603 a learning objective, where a portion of image patches are dropped and the network seeks
604 to reconstruct the full input. We use the ViT-L/16 model pre-trained on ImageNet at a
605 resolution of 448 (Hu et al., 2022).² The *keys* of the last layer of the *encoder* network are
606 used. No masking is performed.
- 607 • **CLIP** (Radford et al., 2021) is trained using image-text pairs on an internal dataset WIT-
608 400M. We use ViT-B/16 model³. We consider two locations to obtain dense features:
609 *keys* from a self-attention layer of the image encoder and *tokens* which are the outputs of
610 transformer layers. We find that *keys* of the second-to-last layer give better performance.
- 611 • We also consider **Stable Diffusion**⁴ (v1.5) itself as a feature extractor. To that end, we use
612 the *queries* from the cross-attention layers in the UNet denoiser, which correspond to the
613 image modality. Its UNet is organised into 3 downsampling blocks, a middle block, and 3
614 upsampling blocks. We observe that the middle layers have the most semantic content, so
615 we consider the middle block, 1st and 2nd upsampling blocks and aggregate features from
616 all three cross-attention layers in each block. As the features are quite low in resolution,
617 we include the first downsampling cross-attention layer and the last upsampling cross-
618 attention layer as well. The feature maps are bilinearly upsampled to resolution 64×64 and
619 concatenated. A noise appropriate for $\tau = 200$ timesteps is added to the input. For feature
620 extraction, we run SD in *unconditional* mode, supplying an empty string for text caption.

621 B.3 DATASETS

622 We evaluate on validation splits of PASCAL VOC (VOC), Pascal Context (Context) and COCO-Object
623 (Object) datasets. PASCAL VOC (Everingham et al., 2010; 2012) has 21 classes: 20 foreground
624 plus a background class. For Pascal Context (Mottaghi et al., 2014), we use the common variant
625 with 59 foreground classes and 1 background class. It contains both “things” and “stuff” classes.
626 The COCO-Object is a variant of COCO-Stuff Caesar et al. (2018) with 80 “thing” classes and
627 one class for the background. Textual class names are used as natural language specification of
628 names. We renamed or specified certain class names to fix errors (e.g. pottedplant → potted
629 plant), resolve ambiguity better (e.g. mouse → computer mouse) or change to more common
630 spelling/word (e.g. aeroplane → airplane), resulting in 14 fixes. We experiment and measure
631 an impact of this in Appendix C.1 for our and prior work.

632 B.4 COMPARATIVE BASELINES

633 We briefly review the prior work in Table 1. Most prior work (Liu et al., 2022; Cha et al., 2022; Xu
634 et al., 2022a; Ren et al., 2023; Luo et al., 2022; Xu et al., 2023b) trains image and text encoders on
635 large image-text datasets with a contrastive loss. The methods mainly differ in their architecture
636 and use of grouping mechanisms to ground image-level text on regions. ViL-Seg (Liu et al., 2022)
637 uses online clustering, GroupViT (Xu et al., 2022a) and ViewCo (Ren et al., 2023) employ group
638 tokens. OVSegmentor (Xu et al., 2023b) uses slot-attention and SegCLIP Luo et al. (2022) a grouping
639 mechanism with learnable centers. CLIPPy (Ranasinghe et al., 2022), TCL (Cha et al., 2022), and
640 MaskCLIP (Zhou et al., 2022) predict classes for each image patch: Ranasinghe et al. (2022) use
641 max-pooling aggregation, Cha et al. (2022) self-masking, and Zhou et al. (2022) modify CLIP

¹Model and code available at <https://github.com/facebookresearch/dino>.

²Model and code from https://github.com/facebookresearch/long_seq_mae.

³Model and code from <https://github.com/openai/CLIP>.

⁴We use implementation from <https://github.com/huggingface/diffusers>.

642 for dense predictions. To assign a background label (Liu et al., 2022; Cha et al., 2022; Xu et al.,
 643 2022a; Ren et al., 2023; Luo et al., 2022) use thresholding while Ranasinghe et al. (2022) uses
 644 dataset-specific prompts. ReCO (Shin et al., 2022b) is closer in spirit to our approach as it uses a
 645 support set for each prompt; this set, however, is CLIP-retrieved from curated image collections,
 646 which may not be applicable for any category in-the-wild.

647 We also note that prior work builds on top of similar pre-trained components such as CLIP in (Shin
 648 et al., 2022b; Zhou et al., 2022; Cha et al., 2022; Luo et al., 2022), DINO + T5/RoBERTa in (Ranas-
 649 inghe et al., 2022; Xu et al., 2023b). We additionally make use of StableDiffusion, which is trained
 650 on a larger dataset (3B, compared to 400M of CLIP). OVDiff is, however, fundamentally different to
 651 all prior work, as (a) it generates a support set of synthetic images given a class description, and (b) it
 652 does not rely on additional training data and further training for learning to segment.

653 B.5 HYPERPARAMETERS

654 OVDiff has relatively few hyperparameters and we use the same set in all experiments. Unless
 655 otherwise specified, $N = 32$ images are sampled using classifier-free guidance scale (Ho & Salimans)
 656 of 8.0 and 30 denoising steps. We employ DPM-Solver scheduler (Lu et al., 2022). When sampling
 657 images for the support sets we also use a negative prompt “text, low quality, blurry, cartoon, meme,
 658 low resolution, bad, poor, faded”. If/when CutLER fails to extract any components in a sampled
 659 image, a fallback of $M_n^{\text{fb}} = A_n > 0.5$ and $M_n^{\text{bg}} = A_n < 0.2$ is used instead. During inference
 660 we set $\eta = 10$, which results in 1024 text prompts processed in parallel, a choice made mainly
 661 to due computational constraints. We set the thresholds for the “stuff” filter between background
 662 prototypes for “things” classes and the foreground of “stuff” at 0.85 for all feature extractors. When
 663 sampling, a seed is set for each category individually to aid reproducibility. With our unoptimized
 664 implementation, we measure around 154 ± 10 s to calculate prototypes for a single category, or 78 ± 4 s
 665 without clustering.

666 B.6 INTERACTION WITH CHATGPT

667 We interact with ChatGPT to categorise classes into “stuff” and “things” for stuff filter component.
 668 Due to input limits, the categories are processed in blocks. Specifically, we input “In semantic
 669 segmentation, there are “stuff” or “thing” classes. Please indicate whether the following class
 670 prompts should be considered “stuff” or “things”:.”. We show the output in Table 4. Note there are
 671 several errors in the response, e.g. glass, blanket, and trade name are actually instances of
 672 tableware, bedding and signage, respectively, so should more appropriately be treated as “things”.
 673 Similarly, land and sand might be more appropriately handled as “stuff”, same as snow and
 674 ground. Despite this, We find ChatGPT contains sufficient knowledge when prompted with “in
 675 semantic segmentation”. We have estimated the accuracy of ChatGPT in thing/stuff classification
 676 using the categories of COCO-Stuff, which are defined as 80 “things” and 91 “stuff” categories.
 677 ChatGPT achieves an accuracy rate of 88.9% in this case.

678 C ADDITIONAL EXPERIMENTS

679 In this section, we provide additional experimental results of OVDiff.

680 C.1 ADDITIONAL COMPARISONS

681 **Category filter.** To ensure that the category pre-filtering does not give our approach an unfair
 682 advantage, we augment two methods (TCL (Cha et al., 2022) and OVSegmentor (Xu et al., 2023b),
 683 which are the closest baselines with code and checkpoints available) with our category pre-filtering.
 684 We evaluate on the Pascal VOC dataset (where the category filter shows a significant impact, see
 685 Table 3) and report the results in Table 5. We observe that TCL improves by 0.6, while the performance
 686 of OVSegmentor drops by 0.1. On the contrary, our method benefits substantially from this component,
 687 but it still shows stronger performance without the filter than baselines with.

688 **CutLER (Wang et al., 2023) baseline.** We also further investigate the use of CutLER to obtain
 689 segmentation masks. In Table 6, we devise a baseline where CutLER-predicted masks are used to
 690 average the CLIP image encoder’s final spatial tokens after projection. Averaged tokens are compared

Table 4: **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

Table 5: 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	56.2	66.4

Table 6: 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	✓	66.3 ± 0.2	29.7 ± 0.3	34.6 ± 0.3

Table 7: 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	✓	66.3 ± 0.2	29.7 ± 0.3	34.6 ± 0.3

Table 8: 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

691 with CLIP text embeddings to assign a class. While relying on pre-trained components (like ours),
 692 this avoids support set generation. In the same table, we also consider whether the objectness prior
 693 provided by CutLER could be beneficial to other methods as well. We consider a version of TCL (Cha
 694 et al., 2022) and OVSegmentor (Xu et al., 2023b) which we augment with CutLER. That is, after
 695 methods assign class probabilities to each pixel/patch, a majority voting for class is performed in
 696 every region predicted by CutLER. This combines CutLER’s understanding of objects and their
 697 boundaries, aspects where prior methods struggle, with open-vocabulary segmentation. However, we
 698 observe that this negatively impacts the performance of these methods, which we attribute to only
 699 limited performance of CutLER in complex scenes present in the datasets. Finally, we also include
 700 a version of OVDiff that does not rely on CutLER for mask extractions, instead using thresholded
 701 masks. We observe that such version of our method also has strong performance, showing that
 702 CutLER is helpful but not a critical component and OVDiff performs strongly without it as well.

703 **Class prompts.** We additionally consider whether corrections introduced to class prompts might
 704 have similarly provided additional benefits to our approach. To that end, we also evaluate TCL and
 705 OVSegmentor (methods that do not rely on additional prompt curation) with our corrected prompts
 706 and consider a version of our method without such corrections in Table 7. We observe only marginal
 707 to no impact to the performance.

708 C.2 ADDITIONAL ABLATIONS

709 **Prototype combinations.** In Table 9, we consider the three different types of prototypes described
 710 in Section 3 and test their performance individually and in various combinations. We find that
 711 the “part” prototypes obtained by K -means clustering show strong performance when considered
 712 individually on VOC. Instance prototypes show strong individual performance on Context, as well as
 713 in combination with the average category prototype. The combination of all three types shows the
 714 strongest results across the two datasets, which is what we adopt in our main set of experiments.

715 We also consider the treatment of prototypes under the stuff filter. We investigate the impact of not
 716 excluding background prototypes for “stuff” classes. In this setting, we measure 29.1 on Context,
 717 which is a slight reduction in performance. We also investigate the benefit of categorisation into
 718 “things” and “stuff” used in the stuff filter component. We instead filter all background prototypes using
 719 all foreground prototypes. In this configuration, we measure 27.6 on Context. Both configurations

Table 9: 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

Table 10: Ablation of different SD feature configurations. Removing first and last cross attention layers, mid, 1st and 2nd 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 11: Comparison with methods when background is *not* considered. We compare OVDiff with prior work on VOC-20, Context-59 and ADE datasets in a setting that considers only the foreground pixels (decided by ground truth). Our method shows comparable performance to prior works despite only relying on pretrained feature extractors. Our results are an average of 5 seeds $\pm\sigma$. * result from (Cha et al., 2022).

Method	VOC-20	Context-59	ADE-150
GroupViT*	79.7	23.4	9.2
MaskCLIP*	74.9	26.4	9.8
ReCo*	57.5	22.3	11.2
PACL	72.3	50.1	31.4
TCL	77.5	30.3	14.9
OVDiff	80.2 \pm 0.6	33.0 \pm 0.2	14.1 \pm 0.2

720 show a reduction from 29.4, measuring using the stuff filter with categorisation in “stuff” and “things”,
721 as used in our main experiments.

722 **K - number of clusters.** In Table 8, we investigate the sensitivity of the method to the choice of
723 K for the number of “part” prototypes extracted using K -means clustering. Although our setting
724 $K = 32$ obtains slightly better results on Context and VOC, other values result in comparable
725 segmentation performance suggesting that OVDiff is not sensitive to the choice of K and a range of
726 values are viable.

727 **SD features.** When using Stable Diffusion as a feature extractor, we consider various combinations
728 of layers/blocks in the UNet architecture. We follow the nomenclature used in the Stable Diffusion
729 implementation where consecutive layers of Unet are organised into *blocks*. There are 3 down-
730 sampling blocks with 2 cross-attention layers each, a mid-block with a single cross-attention, and 3
731 up-sampling blocks with 3 cross-attention layers each. We report our findings in Table 10. Including
732 the first and last cross-attention layers in the feature extraction process has a small positive impact
733 on segmentation performance, which we attribute to the high feature resolution. We also consider
734 excluding features from the middle block of the network due to small 8×8 resolution but observe a
735 small negative impact on performance on Context dataset and substantial decrease on VOC. We also
736 investigate whether including the first (Up-1) and the second upsampling (Up-2) blocks are necessary.
737 Without them, the performance drops the most out of the configurations considered. Thus, we use a
738 concatenation of features from the middle, first and second upsampling blocks and the first and last
739 layers in our main experiments.

740 C.3 EVALUATION WITHOUT BACKGROUND

741 One of the notable advantages of our approach is the ability to represent background regions via
742 (negative) prototypes, leading to improved segmentation performance. Nevertheless, we hereby also
743 evaluate our method under a different evaluation protocol, adopted in prior work, which excludes the

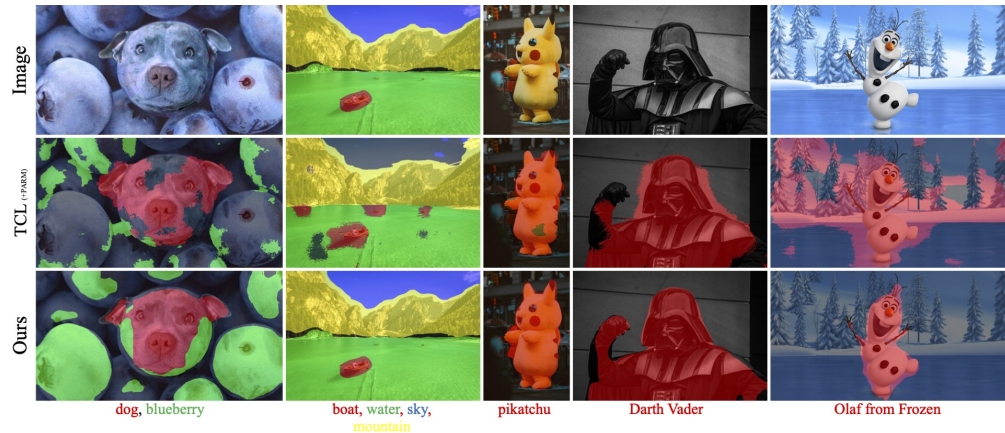


Figure 6: Qualitative comparison on in-the-wild images. OVDiff performs significantly better than prior state-of-the-art, TCL, on a confusing composite (photoshopped) image, a scenery photo, and realistic and cartoon images containing popular characters.

744 *background* class from the evaluation. We note that prior work often requires additional considerations
 745 to handle background, such as thresholding. In this setting, however, the background class is *not*
 746 predicted, and the set of categories, thus, must be exhaustive. As in practice this is not the case, and
 747 datasets contain unlabelled pixels (or simply a background label), such image areas are removed from
 748 consideration. Consequently, less emphasis is placed on object boundaries in this setting. We test our
 749 method on three datasets: PascalVOC without background termed VOC-20, Pascal Context without
 750 background termed Context-59, and ADE20k (Zhou et al., 2017) which contains 150 foreground
 751 classes. As in this setting the background prediction is invalid, we do not consider negative prototypes.
 752 This setting tests the ability of various methods to discriminate between different classes, which for
 753 OVDiff is inherent to the choice of feature extractors. Despite this, our method shows competitive
 754 performance. There exists a notable gap between PACL and other works, including ours, on Context-
 755 59 and ADE-150. In the case of OVDiff, we attribute this to the limited resolution of our feature
 756 extractors, especially on ADE-150 where a variety of tiny objects is present. PACL, on the other
 757 hand, proposes a method to increase the resolution of their trained network 4 times during inference.

758 C.4 QUALITATIVE RESULTS

759 We include additional qualitative results from the benchmark datasets in Fig. 7. Our method achieves
 760 high-quality segmentation across all examples, without any post-processing or refinement steps.
 761 Finally, in Fig. 8, we show examples of support images sampled for some thing, and stuff categories.

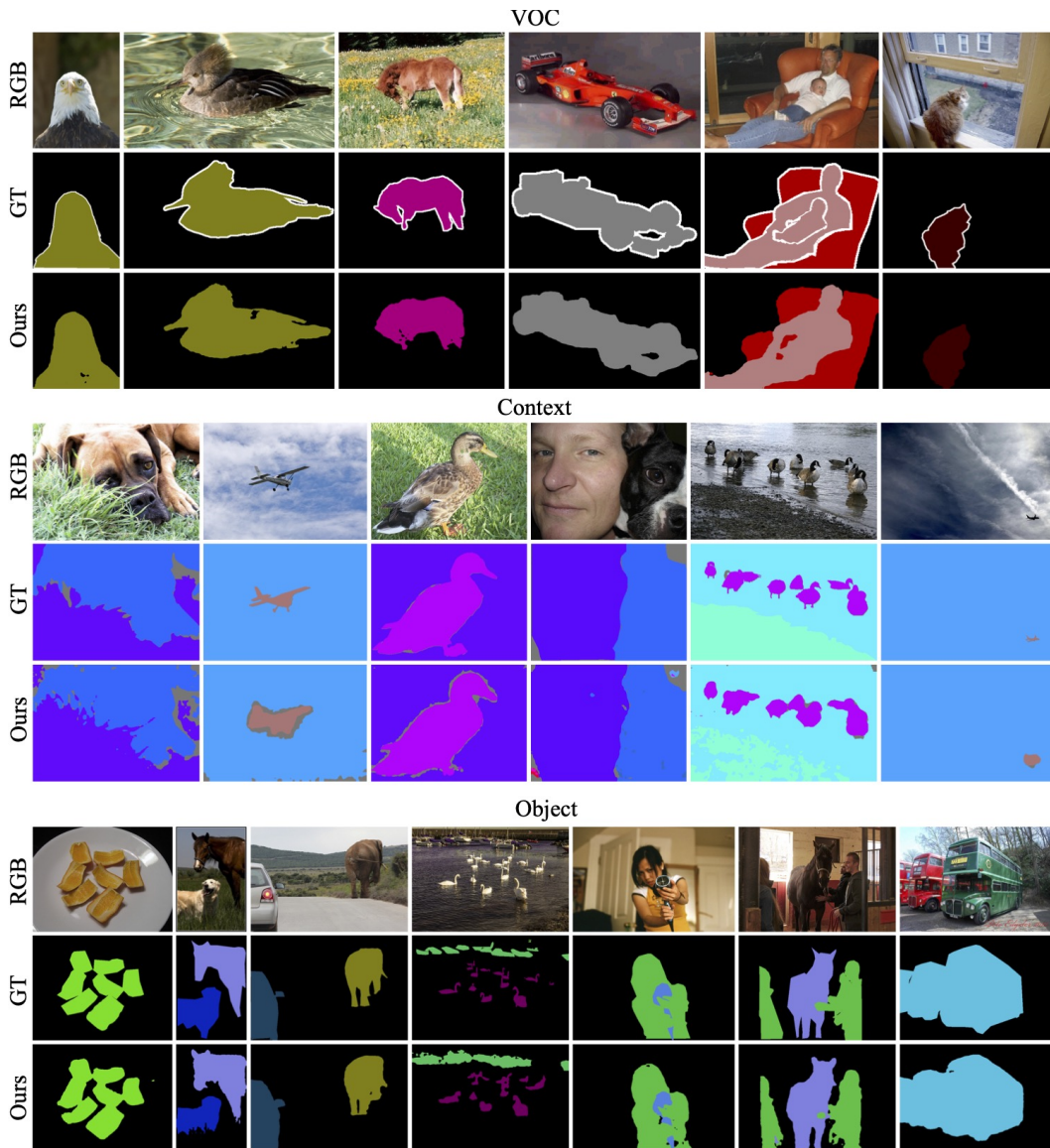


Figure 7: Additional qualitative results. Images from Pascal VOC (top), Pascal Context (middle), and COCO Object (bottom).

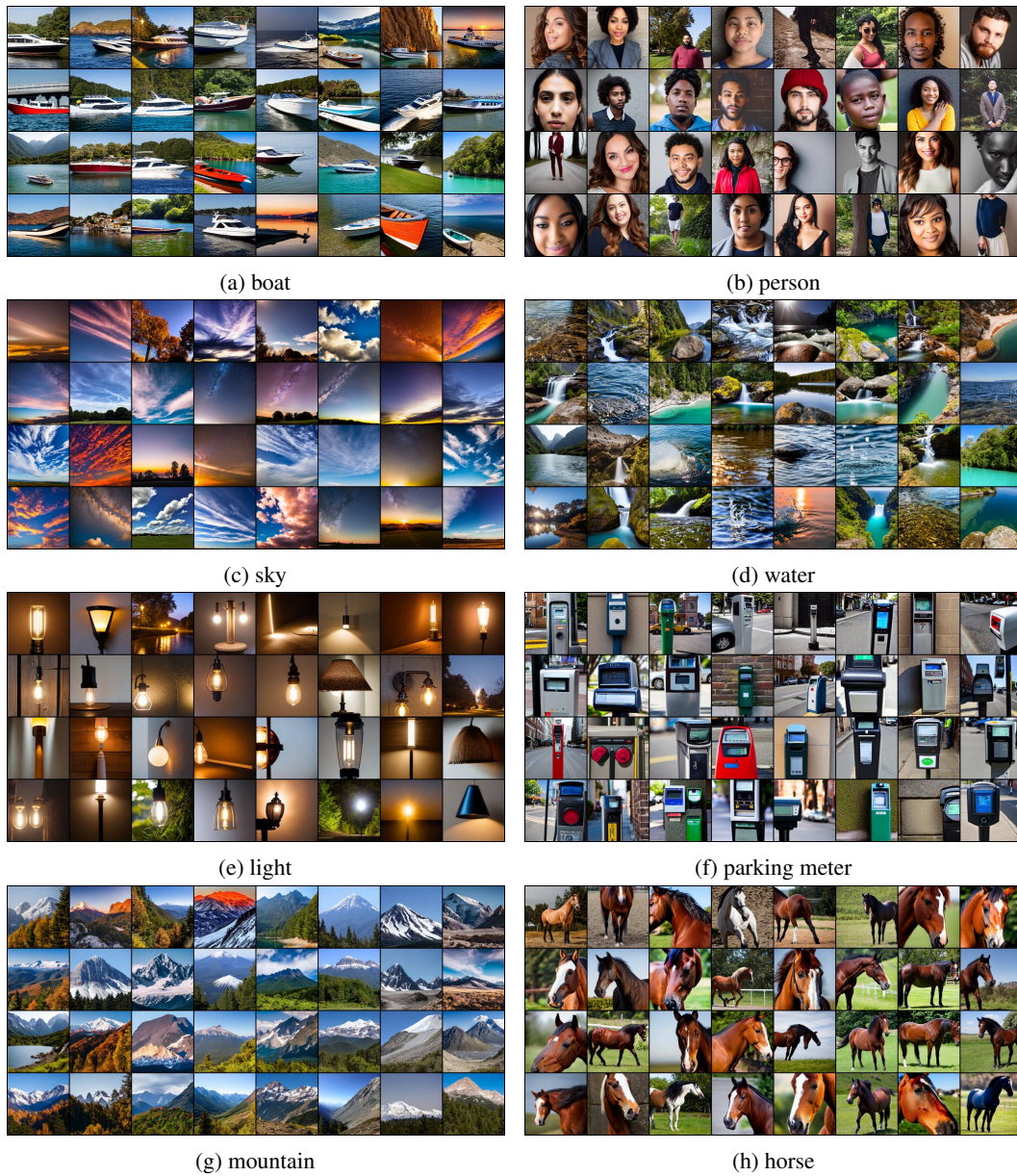


Figure 8: Images sampled for a support set of some categories.