# DINOv2 Meets Text: A Unified Framework for Image- and Pixel-Level Vision-Language Alignment

Cijo Jose[1]    Théo Moutakanni[1,2]    Dahyun Kang[1,3,*]    Federico Baldassarre[1]

Timothée Darcet[1,4]    Hu Xu[1]    Daniel Li[1]    Marc Szafraniec[1]    Michaël Ramamonjisoa[1]

Maxime Oquab[1]    Oriane Siméoni[1]    Huy V. Vo[1]    Patrick Labatut[1]    Piotr Bojanowski[1]

[1] Meta FAIR    [2] Université Paris-Saclay, CentraleSupélec, MICS
[3] POSTECH    [4] Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP

## Abstract

*Self-supervised visual foundation models produce powerful embeddings that achieve remarkable performance on a wide range of downstream tasks. However, unlike vision-language models such as CLIP [64], self-supervised visual features are not readily aligned with language, hindering their adoption in open-vocabulary tasks. Our method, named* dino.txt, *unlocks this new ability for DINOv2 [60], a widely used self-supervised visual encoder. We build upon the LiT training strategy [92], which trains a text encoder to align with a frozen vision model but leads to unsatisfactory results on dense tasks. We propose several key ingredients to improve performance on both global and dense tasks, such as concatenating the* [CLS] *token with the patch average to train the alignment and curating data using both text and image modalities. With these, we successfully train a CLIP-like model with only a fraction of the computational cost compared to CLIP while achieving state-of-the-art results in zero-shot classification and open-vocabulary semantic segmentation.*
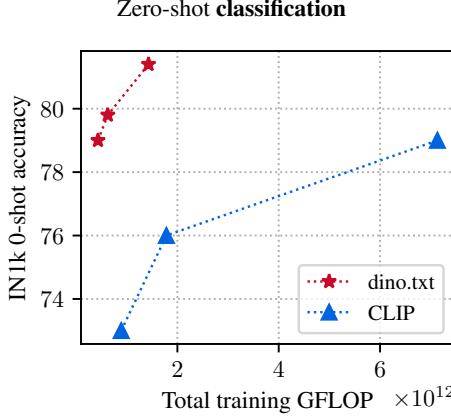
## 1. Introduction

The advent of modern vision foundation models trained in a Self-Supervised Learning (SSL) fashion [4, 11, 14, 37, 60] has resulted in robust, generic features that achieve impressive performance on downstream tasks. These features are typically used *as is*, plugged into a light-weight adapter such as a linear classifier, and deliver strong results without requiring a costly fine-tuning process. As a result, a single strong vision backbone can be used simultaneously for different tasks. DINOv2 [60], in particular, has been popular for its versatility. This self-supervised model, trained to capture both the global context and local information of the image, has led to state-of-the-art performance in tasks that require an overall understanding of the image such as classification and those that necessitate more fine-grained details such as segmentation [50], canopy height prediction [75], object matching [25, 61], object discovery [21] and tracking [77, 91]. However, self-supervised vision models do not provide an interface with language, limiting their use in open-vocabulary scenarios in which multi-modal models [15, 74] that come with built-in language-vision alignment excel. This is a notable weakness in the era of complex and promptable machine learning systems. We aim in this work to equip DINOv2 with a language interface by aligning its feature space with language, which allows us to leverage the strengths of this powerful self-supervised model to tackle open-vocabulary recognition tasks.

Most advanced text-aligned vision foundation models learn with a variation of the CLIP algorithm [64], which trains the visual and textual encoders to align their modality representation in a shared space by exploiting a large-scale, often noisy, paired image-text dataset. They are typically trained from scratch, leading to heavy computational cost. Locked-image text tuning (LiT) [92] is a variant of CLIP that uses a frozen pre-trained vision model as the vision encoder and only train the text encoder to align to the vision encoder's embedding space. This leads to a lower computational cost while retaining desirable properties of the vision encoder. In this work, we argue that given readily available strong vision encoders, we could, and should achieve better vision-language alignment than CLIP at a much smaller cost. To this end, we explore the application of LiT training with DINOv2 as the vision encoder.

As shown in Table 1, 3, applying LiT on strong DINOv2 encoder is not straightforward as it leads to unsatisfying performance on tasks that require fine-grained details such as semantic segmentation or image-text retrieval. First, it is not trivial to obtain good dense features from a model

---

Figure 1. **Zero-shot classification and open-vocabulary segmentation results with our method `dino.txt` trained with only weak image/caption annotation.** We show that our training strategy leads to state-of-the-art results in zero-shot classification with a fast and efficient training. Also, its produces high quality segmentation results on diverse images showing the quality of the image-to-text alignment.

trained with the CLIP/LiT training paradigm, which contrasts only global image and text representation. Second, the domain gap between visual pre-training data and LiT training data caused by the frozen vision encoder could potentially hinder the alignment of images and their captions. To address these issues, we introduce several improvements to the LiT paradigm. Instead of the commonplace practice of using the `[CLS]` token from the vision encoder to represent the image, we concatenate this token to the average pool of all patch tokens in the image as the vision representation to allow aligning both the global context and the local information of the image to its textual description. Then, we reduce the aforementioned domain gap by adding two learnable vision blocks on top of the frozen vision encoder, thereby allowing vision features to adapt to the new training data. We show the benefit of our training for zero-shot classification and open-vocabulary segmentation[1] in Figure 1.

Moreover, the quality of pre-training data has been shown to strongly influence the model's performance [64, 88] but also the training efficiency [27]. We show that by paying attention to the dataset curation, we can further improve our training procedure. Indeed, we propose to curate the training dataset by balancing the long-tailed distribution of image and text data. A well-balanced data distribution eases the training, allowing us to reach good performance with only a fraction of the computational cost. This in turns allows us to experiment with a wider text encoder, leading to further improvements in performance. Finally, our study not only unlocks text alignment for DI-NOv2 but also reveals limitations in the LiT framework discussed in the error analysis section, pointing toward future directions for more effective and efficient frameworks to achieve language-aligned vision foundation models.

To summarize, the contributions of this work include:

---

<sup>1</sup>We follow the taxonomy in the survey [82] and use the nomenclature 'open-vocabulary' segmentation.

**(i)** a new method `dino.txt`, which unlocks image- and pixel-level text-alignment for DINOv2, **(ii)** key ingredients on top of existing works that allow to train such multi-modal alignment for only a fraction of the usual compute cost, **(iii)** an extensive error analysis, demonstrating the limitations of existing segmentation benchmarks for this task and the different error types with these models.

## 2. Related work

**Self-supervised feature learning.** Visual features from image encoders trained in a self-supervised fashion have been used in many machine learning systems due to their good performance and generalizability. Multiple approaches for learning these models have been developed in recent years. Among these, contrastive learning [59] trains models to pull features of similar images while pushing those of dissimilar images. Notable methods include MoCo [36] which employs a memory bank, BYOL [34] which removes the need for negative pairs, SwAV [10] which contrasts online cluster assignments, or DINO [11] which extends SwAV to Vision Transformers [23]. In contrast, reconstruction-based methods learn by predicting hidden portions of input images such as missing pixels (MAE [37]), patches from quantized code book (BeiT [5]) or patch features in a latent space (I-JEPA [4]). With such an array of viable approaches, a determining factor remains whether the SSL methods can improve with increasing data and model sizes: scaling was explored in multiple works [9, 31–33], with DINOv2 [60] setting the current state of the art for this problem.

**Contrastive text-image pre-training.** The idea of leveraging textual metadata to train image understanding models has a long history in computer vision [24, 28, 35, 46]. In the context of deep neural networks, Joulin *et al.* [42] proposed to use words from image captions as targets to train visual encoders. This core idea has been further im-

proved in CLIP [64]. The authors propose to encode the image and the caption, and train both using a contrastive loss. Deep encoding of captions facilitates robust generalization across sentences, such as by collapsing synonyms, thereby enhancing learning efficiency. Since the original CLIP was trained on a private dataset, open source reproduction attempts have focused on collecting public large-scale datasets (LAION [70]), leading notably to the Open-CLIP [16] family of models. More recently, Xu et al. [88] described a simple procedure for re-balancing image-text web data, reproducing the performance of CLIP. The zero-shot performance was further improved by DFN [27] which proposes filtering the training data to match the distribution of downstream tasks. Even more refined systems have recently been developed: EVA-CLIP [73, 74] or InternVL [15], demonstrated at scales above 5 billions parameters and further narrowed the gap between fully-supervised and zero-shot models on ImageNet.

Apart from data and model scaling, a few modifications to the initial training algorithm have been proposed. SigLIP [93] considers a binary log loss instead of multinomial cross entropy. LLiP [48] proposes applying the CLIP loss between the text and image register tokens [18] that are conditioned on the text tokens. In contrast, we do not use any improved loss function in our training, and stick to CLIP's original contrastive loss with a frozen image encoder and a learnable text encoder, following the procedure described in locked-image text tuning (LiT [92]).

**Automatic data curation at scale.** Automatic data curation plays a crucial role in the training of foundation models with massive-scale web-crawled datasets, which typically reach over hundreds of millions [29, 69, 88] of data. At this scale, manual annotation becomes infeasible, so these datasets are often collected without supervision from the internet. Such in-the-wild data inherently exhibit a long-tailed distribution of data categories [53] which limits a model's ability to efficiently learn to cover broad concepts. To address this issue, related works on foundation models often construct balanced training datasets by suppressing head (frequent) concepts and boosting tail (rare) ones. For example, CLIP's [64] data preparation pipeline collects five hundred thousand frequent words and queries each word in the raw dataset to retrieve a balanced number of (image, caption) pairs. The unrevealed details of this pipeline are later reproduced and formalized by MetaCLIP [88]. DFN [27] trains a data evaluation model that assesses the "quality" of data to sample the top samples among the raw data pool. SemDeDup [2] prunes it by removing duplicated data points detected with clustering. Finally, Vo et al. [79] balances data distribution by sampling uniformly over the data support. These methods apply data curation to either images or captions, while we balance both distributions, leading to better performance and more efficient training.

**Open-vocabulary segmentation.** CLIP models can be adapted to produce patch-level features aligned with text by performing several forward passes on different views of the image [3, 41, 45, 85] or producing code books [44, 71, 72] of prototypes per concept of interest. MaskCLIP [97], which can be applied to most vision-language models (VLMs), adapts the model by removing the final global pooling and applying the final projection to the value embedding of the last attention layer, achieving dense features in CLIP space. Such features can be refined with improved attention mechanisms [7, 80], or using an SSL model as a guide [43, 47, 86]. Such efforts are orthogonal to our work as they can be applied to any dense CLIP-like features.

Improved patch-level alignment can be obtained by fine-tuning or training from scratch a CLIP-like model with dedicated objectives using pixel-level annotations [22, 30, 49, 51, 66, 81] or coarse image/caption annotations [12, 30, 51, 52, 54, 58, 65, 89, 90, 94]. In this work, we focus on the latter. ViewCO [68] leverages multi-view consistency and CLIPSelf [84] uses a teacher-student learning strategy to produce dense features aligned with those obtained from crops. GroupViT [89] integrates learnable tokens that are trained using grouping blocks, which CoCu [87] further improves by using image retrieval on image-caption pairs to create concept banks and use them as training data. PACL [58] trains patch-to-text affinity with a dedicated module, TCL [12] proposes a local contrastive objective to align well-selected patches to the text, and CoDe [83] uses a word-region local contrastive objective to match regions of the image to segments of the text. Closer to us, CLIPpy [65] fine-tunes an SSL vision backbone and pre-trained text encoder to produce aligned features (using the average patch), however, at the cost of worse classification results. In this work, we show that it is possible to train a model with both image- and pixel-level alignment with a simple loss.

# 3. Proposed approach: `dino.txt`

In this work, we demonstrate the simplicity and efficiency of aligning a text encoder to a self-supervised visual foundation model. We show that we can directly use the foundation model's embedding space to perform both zero-shot classification and open-vocabulary semantic segmentation. The first section defines the model architecture and training objective. We then describe the data curation that establishes our training dataset, followed by our inference protocol. An overview of our method is shown in Figure 2.

## 3.1. Locked-image text alignment

**Image and text encoders.** For the image encoder, we use a frozen ViT model [23], denoted $\phi_V(\cdot)$, trained in a self-supervised fashion following DINOv2 [60]. The encoder takes as input an image $\mathbf{I} \in \mathbb{R}^{3 \times H \times W}$ and divides it into a sequence of $N$ patch tokens to which is prepended a
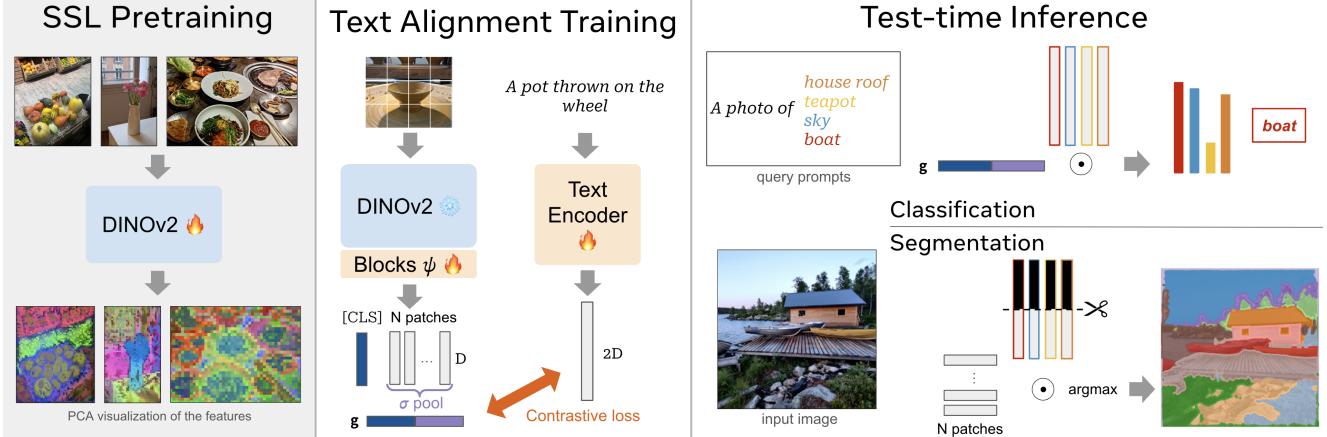
Figure 2. **Overview of our method `dino.txt`.** We first show the localization quality of the self-supervised features (left). We then present our training strategy (middle) which consists in aligning the frozen SSL backbone to a text-encoder trained from scratch. We additionally add a light vision block on top of the visual encoder in order to better align with the text. We train our model for just 50k iterations and achieve SoTA results on both zero-shot classification and open-vocabulary segmentation (right).

learned `[CLS]` token, thus giving $\phi_{\mathrm{V}}(\mathbf{I}) = [\mathbf{c}, \mathbf{f}_1, \cdots, \mathbf{f}_N]$ in $\mathbb{R}^{(1+N)\times D}$. Here, $\mathbf{c} \in \mathbb{R}^D$ represents the `[CLS]` output embedding and $\mathbf{f}_p \in \mathbb{R}^D$ denotes the output embedding for a patch $p \in [1, \cdots, N]$. We discard any potential register tokens [18] as they are not used.

The text encoder consists of a series of Transformer [78] blocks and a single linear layer on top that maps features to the image embedding size. All parameters of the text encoder are trained from scratch following LiT [92]. We follow [64, 92] and align the output `[EOS]` text token to the image embedding, therefore obtaining global alignment between the corresponding sentences and images.

**Improved image representation for text alignment.** Here we describe our choice of image representation for the image-text alignment. We aim to improve the global-level text alignment used for classification and retrieval tasks, as well as the patch-level alignment for segmentation, using *a single learning objective that does not require any pixel-level supervision*. Previous works have proposed approaches designed specifically for a task: for image-level classification, the `[CLS]` embedding $\mathbf{c}$ is the predominant choice for text alignment [64], while for segmentation, the final patch embeddings can be pooled [58, 65], *e.g.*, using max or average pooling, enforcing gradient to the patches but unfortunately hurting classification performance [65]. Instead, we aim here to enforce *both* global and local alignment with text. To do so, we concatenate the `[CLS]` embedding to the average-pooled patch embeddings.

We further improve the alignment to the text modality by adding two trainable transformer blocks, noted $\psi$, on top of the frozen vision backbone, which we refer to as "vision blocks" throughout the paper. We use the outputs of the blocks $\psi$, which preserves the dimensionality $D$ of the

descriptors, following:

$$\psi([\mathbf{c}, \mathbf{f_1}, \cdots, \mathbf{f_N}]) = [\mathbf{c}', \mathbf{f_1}', \cdots, \mathbf{f_N}'], \tag{1}$$

to produce a representation for the image. Specifically, we concatenate the updated class token $\mathbf{c}'$ with the pooled patch tokens, obtained by applying a pooling operation $\sigma$ over the patches $[\mathbf{f_1}', \cdots, \mathbf{f_N}']$ and produce the global descriptor $\mathbf{g}$ of dimension $\mathbb{R}^{2D}$:

$$\mathbf{g} = [\mathbf{c}'; \sigma([\mathbf{f_1}', \cdots, \mathbf{f_N}'])], \tag{2}$$

with ';' denoting the concatenation. We found that using average pooling for $\sigma$ yields the best results. By propagating the gradient through the average of patch tokens, each token can learn to contribute to the final descriptor, enabling a more granular alignment. We discuss the importance of this representation for dense tasks in Table 2. Interestingly, we observe that this joint representation improves alignment for downstream classification and segmentation tasks.

**Contrastive locked-image text alignment.** The image-text alignment objective encourages the text representation to be close to its paired image while simultaneously repelling it from non-corresponding images. As discussed, we use the image descriptor $\mathbf{g}$ as the alignment target to the text embedding. With the image backbone frozen, the trainable parts consist of the additional vision blocks and the text encoder. We train with the contrastive learning objective [64] on a dataset of image-caption pairs, which is automatically curated without any supervision.

### 3.2. Text- and image-based data curation

Training data plays a crucial role in machine learning model performance [79, 88]. CLIP-style VLM training requires good quality image-text pairs, however existing data curation methods for this rely only on the text modality [64, 88]

4

- click to enlarge
- ~product.metadata.name~
- Certified pre-owned 2018

Figure 3. **Examples of poor, ambiguous or too generic captions** found in our data pool.

which we later show is suboptimal. For example, recent MetaCLIP [88] curates CLIP training data from a large pool of image-text pairs collected from the Internet. It first constructs a set of text queries based on WordNet [55] synsets and Wikipedia. Next, a mapping is establishes from each query to the set of image-text pairs whose caption contains the query. Finally, pairs are sub-sampled from each of these sets to form the final dataset. This approach results in a more balanced distribution of concepts within the data pool.

Although the pipeline significantly improves the performance of CLIP, it overlooks the distribution of visual concepts appearing in the data pool. This would not be an issue if there was a perfect alignment of concepts between the images and their captions because selecting data based on text would lead to a balanced selection in the visual space. However, captions automatically collected from the Internet are often noisy and do not exactly describe what is depicted in the corresponding images (see Figure 3), therefore, ignoring the image distribution is suboptimal.

In this work, we propose to balance both the caption and image distributions by combining text and image based curation. We use [88] for text curation while for image curation, we use the clustering-based method of Vo et al. [79]. The latter divides the data pool into coherent clusters from which data are sampled to form a curated dataset. As opposed to [88], this method does not require a pre-defined set of concepts to perform clustering, which is not trivial to construct in the visual space. Instead it builds clusters using hierarchical $k$-means. Clusters obtained this way distribute evenly over the space and their size naturally follow a long-tailed distribution. The curated dataset is then formed by sub-sampling from the clusters to diminish the impact of head clusters and thus balance the concepts. In this work, we propose applying this curation method on images and the pipeline of [88] on captions. We then take the intersection of these results to form the final selection.

### 3.3. Inference

At inference, we consider a set of text queries $\mathcal{Q}$ which we want to compare either with the image representation (*e.g.*, for classification and retrieval tasks) or to the image pixels for dense tasks. In both cases, each text query is encoded by the trained text encoder as $T_q = \phi_T(q) \in \mathbb{R}^{2D}$ for $q \in \mathcal{Q}$.

In order to perform open-vocabulary zero-shot classification and retrieval, we extract a global image descriptor $\mathbf{g}$ which we compare with each of the text queries using cosine

similarity. In the case of a dense task which requires pixel-level features, with our model, there is no need to adapt the output specifically to the task, *e.g.*, as done in MaskCLIP [97]. Instead, we extract for each patch $p \in [1, \cdots, N]$ the final representation $\mathbf{f}'_{\mathbf{p}}$ in $\mathbb{R}^D$ outputted by the model. We then compare, using cosine similarity, each patch representation with the part of the text embedding aligned during training with the average patch. We then upsample the logits to fit to the image resolution. Doing so, we obtain good quality predictions without needing any model adaptation [97]. This is possible because the patch space benefits both from the SSL localization quality and the alignment to the text learned with our objective.

## 4. Experiments

### 4.1. Tasks and metrics

**Zero-shot classification.** We evaluate zero-shot classification using the protocol described in CLIP on ImageNet-1K [19] (IN1K), ImageNet-v2 [67] (IN-v2), ImageNet-Adversarial [38] (IN-A), ObjectNet [6] (ObjNet), iNaturalist2021 [40] (iNat21) and Places205 [95] (PL205). At test time, we feed the class names to the text encoder to retrieve text vectors, and measure their cosine similarity with the global descriptor produced by the image encoder.

**Image-text retrieval.** We evaluate image-text retrieval on the standard cross-modal retrieval benchmarks: COCO2017 [76] and Flickr30K [63]. These datasets comprise pairs of images and their corresponding descriptive captions. The task involves finding the most similar image based on a text query. We use the metric Recall@1, which equals 1 if the nearest image matches the ground-truth pair, and 0 otherwise.

**Open-vocabulary segmentation.** We evaluate the results of dino.txt on the task of open-vocabulary segmentation on the datasets ADE20K [96](ADE), Cityscapes [17] (City.), COCO-Stuff [8](Stuff), PASCAL Context [56] (C59), and PASCAL VOC20 [26](VOC). We employ the mIoU metric (mean intersection-over-union). In order to generate pixel-level features, we use the inference procedure detailed in Section 3.3 and additionally follow the sliding window protocol of TCL [12]. If not stated otherwise we use the final patch token embeddings produced by our model. However, most general-purpose image-text encoders [27, 64, 73, 74, 93] do not apply supervisory signals on the final patch embeddings leading to poor output patch quality. Therefore, we employ the well-known MaskCLIP [97] strategy to evaluate such methods on the segmentation task. It forwards the value embeddings in the last attention layer (bypassing the final attention) resulting in patch embeddings in the aligned space.

5

## 4.2. Implementation details

**Training.** We implement our training framework in PyTorch [62]. We follow the implementation of the CLIP loss function from the OpenCLIP library [39]. We employ the `torch.compile` feature of PyTorch for maximally efficient training on Nvidia A100 GPUs with 80 GB VRAM. The DINOv2 vision encoder initialized from [18] is kept frozen which saves compute and memory, allowing larger batch sizes compared to CLIP, which is important as shown in Table 3. For numbers in Table 1,2 we follow the CLIP recipe and set the batch size to 32K. In other tables we use 65K for better results. We also observe that good results can be obtained by training for 50K iterations which corresponds to 1.6 and 3.2 billion image-text pairs seen at 32K and 65K batch size respectively. We chose this setup as default and discuss more hyper-parameter in Section 4.3.

**Training dataset.** We apply the text and image curation process described in Section 3.2 to an initial data pool derived from CommonCrawl [1], consisting of 2.3 billion image-text pairs. We sample 650 million pairs per-epoch using our curation strategy. For the text-based curation part, text frequencies in the data pool are precomputed offline, and data with frequent texts are stochastically dropped following [88]. This process keeps 900 millions pairs per epoch. For image-based curation, we use pre-trained DINOv2 ViT-L/14 to extract embeddings for an offline 3-level hierarchical $k$-means [79] with 20M, 800K and 80K centroids respectively on each level. We similarly drop pairs whose images appear in large clusters, resulting in 1.5 billion pairs per epoch. Our final training dataset for the given epoch consists of the text-image pairs kept in both the text- and image-based curation process, hereafter noted as LVTD-2.3B which stands for Large Vision Text Dataset.

**High-resolution inference.** A typical segmentation protocol, popularized by TCL [12], consists of applying a sliding window strategy and aggregating the segmentation results in a single prediction map. We extend this strategy to a high-resolution windowing procedure in which we sample crops of various sizes (1%, 10%, 100% of the total area) in a dense sliding window manner, and add noise to the coordinates, such that the crops correspond to non-rectangular quadrilaterals. We distort the crops into squares, extract features, then project the features back onto the dense pixel grid with interpolation, and average all contributions. We cluster features using $k$-means with $k$=32, then run the zero-shot classifier on the centroids. For our results using this procedure, each pixel is visited on average 40 times, for a total of approximately 800 crops processed by the vision model in 10 seconds on an A100 GPU. This approach showcases the features at finer scales, and improving the procedure is a direction for future work. We provide results in Table 6 (last row) and visualization in Figure 4.

## 4.3. Ablation study of our method `dino.txt`

We study in this section the impact of different components of `dino.txt`. In all experiments, we train models on the text-based curated dataset obtained from LVTD-2.3B following Xu et al. [88] unless otherwise specified.

| Model | Vision backbone | Arch. | *class.* IN1K | *retrieval* COCO |
|---|---|---|---|---|
| CLIP | scratch | ViT-L/14 | 73.0 | **38.0** |
| LiT | MAE [37] | ViT-L/14 | 52.3 | 13.6 |
| | I-JEPA [4] | ViT-H/14 | 67.7 | 20.1 |
| | DINO [11] | ViT-B/8 | 71.3 | 26.5 |
| | DINOv2 [60] | ViT-L/14 | **78.8** | 30.2 |

Table 1. **Comparison of trainable vision backbone (CLIP) and frozen pre-trained SSL backbones (LiT).** We produce CLIP results using exactly CLIP recipe in our setup. All the models are trained for 50K iterations.

**LiT with SSL is not obvious.** In order to train a text encoder to align with DINOv2 features, we resort to LiT. Our preliminary LiT experiments with DINOv2 ViT-L/14 vision encoder and a pre-trained BERT-base [20] led us to 70.0 zero-shot accuracy on IN1K when training on CC12M [13]. For comparison, the original LiT paper reports 67.6 with a ViT-L/16 vision encoder and pre-trained BERT-large text encoder when training on the same dataset. We next train our models on the larger dataset LVTD-2.3B, and present in Table 1 a comparison between CLIP and LiT with different pre-trained vision encoders. We observe that among considered vision backbones, DINOv2 leads to the best performance. It enables LiT to achieve good results in classification. However, there is a drop in retrieval performance compared to CLIP, likely due to the frozen vision encoder not being able to adapt to new training data. These results suggest that we need a new strategy to align a frozen backbone encoder to text that can generalize for different tasks, such as our proposed `dino.txt`.

| $\sigma$ pooling | *class.* IN1K | *retrieval* COCO | *seg.* ADE |
|---|---|---|---|
| `[CLS]` | 78.8 | 30.2 | 8.3 |
| `[avg]` | 74.7 | 32.7 | 13.3 |
| `[max]` | 70.2 | 25.7 | 18.0 |
| `[CLS max]` | 78.2 | 31.9 | 16.8 |
| `[CLS avg]` | **79.2** | **34.7** | **18.2** |

Table 2. **Impact of the pooling operation** $\sigma$ used at training in `dino.txt` on zero-shot performance. Results of the first row are produced with MaskCLIP strategy, others with the output patches. The experiment corresponds to the first row in Table 3.

**Impact of the pooling operation** $\sigma$ **at training.** We evaluate in Table 2 the impact of the choice of the pooling operation, applied during training, to the performance on downstream tasks. Typically, methods have used max [65] or average pooling [58] in order to align patch embeddings to the text, but this hurts classification results. We observe the same phenomenon in our experiments where max or average pooling alone degrades classification performance. However, when using our proposed concatenation pooling `[CLS avg]`, we obtain a significant boost for *both* classification and dense tasks, showing that there is no need to choose between one task or the other.

| | *class.* | *retrieval* | *seg.* |
|---|---|---|---|
| Ablation | IN1K | COCO | ADE |
| Reference | 78.8 | 30.2 | 8.3 |
| + `avg-pool` | 79.2 | 34.7 | 18.2 |
| + 65K batch size | 79.8 | 35.1 | 18.2 |
| + 1 vision block | 79.8 | 40.8 | 20.5 |
| + 2 vision blocks | 79.7 | 42.1 | 20.4 |
| + text-large (768 → 1280) | 80.8 | 43.9 | 20.5 |
| + img-based curation | **81.4** | **45.4** | **20.6** |

Table 3. **Exploration of the `dino.txt` recipe.** We start from a 'Reference' configuration which consists of training LiT with a frozen DINOv2-ViT-L/14 vision encoder and a BERT-base sized text encoder trained from scratch; we then add modifications progressively. The first row is evaluated following MaskCLIP whilst the next ones use the output patch tokens.

**Improved training recipe.** We evaluate the impact of the different components of our training in Table 3. Beside our `pool` strategy, we observe that using a larger batch size also improves results on all tasks. Interestingly, the addition of two learnable vision blocks on top of the vision encoder significantly improves retrieval results, showing that the task requires a better visual alignment to the text. Increasing the text embedding size from 768 to 1280 also induces a large gain on all tasks. Finally, we can observe the importance of the combined image- and text-based data curation, which is further discussed below.

| curation | | *class.* | *retrieval* | *seg.* |
|---|---|---|---|---|
| image | text | IN1K | COCO | ADE |
| | | 80.3 | 42.9 | 20.0 |
| | ✓ | 80.8 | 43.9 | 20.5 |
| ✓ | | 80.9 | 43.7 | 20.4 |
| ✓ | ✓ | **81.4** | **45.4** | **20.6** |

Table 4. **Impact of the image- and text-based curation**.

**Impact of dataset curation.** Table 4 decouples the impact of each data curation strategy on `dino.txt` results. Both text- and image-based data curation help to re-balance the long-tailed data distribution, and boost performance. Our proposed combination of them leads to the best performance on all of the three tasks. This result highlights

the important of curating data based on both text and visual modality for visual-language training.

## 4.4. Comparisons to state of the art

**Zero-shot classification and retrieval.** We compare `dino.txt` with state-of-the-art baselines in Table 5 on two image-level understanding tasks: zero-shot image classification and cross-modal retrieval. Our model is on par or better than alternative CLIP-like models on classification benchmarks, setting the state-of-the-art performance on IN-v2, IN-A and the challenging iNaturalist datasets. It can also be observed that the perfomance of `dino.txt` is lower than competitors such as SigLIP [93] on text-image retrieval tasks. This is likely due to the unsatisfactory quality of our trained text encoder, which in turn is a consequence of freezing the vision encoder, as discussed later in Section 4.5. However, we see next that `dino.txt` largely outperforms SigLIP in open-vocabulary segmentation task.



Figure 4. **High-resolution inference. Left**: input image. **Middle**: result of $k$-means clustering ($k$=32) on the features. **Right**: open-vocabulary predictions with the ADE20K class names.

**Open-vocabulary segmentation.** On segmentation, our approach greatly outperform alternatives as shown in Table 6, and performs on par or better than specialized models, without any engineering refinement: we simply apply the classification model on the local features. We note that, compared to other methods, the performance trends lower on VOC20 while being higher on other datasets, which we attribute to a domain gap, as VOC20 frequently contains only one close-up centered object. We ablate in the appendix the impact of using our proposed representation versus MaskCLIP. We also show the quality of our results with our high-resolution strategy in Figure 4.

**Training efficiency.** We show in Figure 1 how zero-shot classification performance on the IN1K validation set evolves as a function of training GFLOP for `dino.txt` and CLIP, trained on the same LVTD-2.3B dataset, described below. On 128 A100 GPUs, 19 hours of training are enough to reach 81.4% on IN1K. In comparison, CLIP requires 110 hours to obtain 79.0%. Furthermore, restricting CLIP training to match the GFLOP budget of our best performing model only achieves 73% accuracy, 8.4% below our model.

| Method | Res. | Dataset | classification | | | | | | retrieval | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | IN1K | IN-v2 | IN-A | ObjNet | iNat21 | PL205 | COCO | Flickr |
| OpenCLIP [39] | 224 | LAION-2B | 74.0 | 66.4 | 48.0 | – | – | – | 46.1 | 75.0 |
| CLIP [64] | 336 | WIT-400M | 76.6 | 70.9 | 77.5 | 72.3 | 5.7 | 59.2 | 37.1 | 67.3 |
| MetaCLIP [88] | 224 | MetaCLIP-2.5B | 79.2 | 72.6 | 72.3 | 75.3 | 10.2 | 61.8 | 45.5 | 76.9 |
| EVA-02-CLIP [73] | 336 | Merged-2B | 80.4 | 73.8 | 82.9 | **78.5** | 6.4 | 61.7 | 47.9 | 77.9 |
| DFN [27] | 224 | DFN-2B | 81.4 | 74.5 | 66.8 | 74.1 | 17.6 | **62.0** | 47.0 | 76.0 |
| SigLIP [93] | 256 | WebLi | 80.4 | 73.8 | 62.0 | 66.6$^\dagger$ | 14.8 | 58.6 | 51.1 | 79.9 |
| SigLIP [93] | 384 | WebLi | **82.1** | 75.9 | 76.4 | 74.0$^\dagger$ | 17.5 | 59.6 | **52.7** | **81.9** |
| dino.txt | 224 | LVTD-2.3B | 81.4 | 75.7 | 80.0 | 72.4 | 18.8 | 61.2 | 45.4 | 77.1 |
| dino.txt | 336 | LVTD-2.3B | 81.6 | **75.9** | **83.2** | 74.5 | **19.4** | 61.2 | 44.9 | 77.6 |

Table 5. **Zero-shot classification and retrieval** results with ViT-L models. $^\dagger$ indicates that the SigLIP results on ObjectNet (authors report 77.9 and 81.0) could not be reproduced despite obtaining matching results on ImageNet-1K.

| Method | Base | Train data. | ADE | City. | VOC | Stuff | C59 |
|---|---|---|---|---|---|---|---|
| *Trained models specialized for segmentation* | | | | | | | |
| GViT [89] | S/16 | CC12M+RC | 9.2 | 11.1 | 79.7 | 15.3 | 23.4 |
| CoCu [87] | S/16 | CC12+3M+Y | 12.3 | 22.1 | 51.4 | 15.2 | 23.6 |
| CLIPpy [65] | B/16 | H-134M | 13.5 | - | 52.2 | - | - |
| TCL [12] | B/16 | CC12+3M | 14.9 | 23.1 | 77.5 | 19.6 | 30.3 |
| CoDe [83] | B/16 | CC12+3M | 17.7 | 28.9 | 57.7 | 23.9 | 30.5 |
| *Generalist image-text encoders* | | | | | | | |
| MetaCLIP [88] | L/14 | MCLIP-2.5B | 2.5 | 1.7 | 23.3 | 3.9 | 6.5 |
| CLIP [64] | B/32 | WIT-400M | 5.0 | 8.6 | 34.7 | 9.0 | 14.2 |
| DFN [27] | L/14 | DFN-2B | 5.8 | 8.0 | 25.9 | 5.1 | 10.0 |
| OpenCLIP [39] | L/16 | LAION-2B | 5.9 | 9.8 | 30.0 | 8.3 | 13.1 |
| CLIP [64] | L/14 | WIT-400M | 6.0 | 11.5 | 24.8 | 7.3 | 10.9 |
| SigLIP [93] | L/14 | WebLi | 9.1 | 18.3 | 30.3 | 9.5 | 13.7 |
| OpenCLIP [39] | B/32 | LAION-2B | 9.9 | 18.1 | 42.9 | 12.7 | 19.0 |
| OpenCLIP [39] | B/16 | LAION-2B | 12.7 | 20.2 | 45.4 | 16.4 | 24.2 |
| dino.txt | L/14 | LVTD-2.3B | **20.6** | **32.1** | **62.1** | **20.9** | **30.9** |
| HR(dino.txt) | L/14 | LVTD-2.3B | *25.1* | *41.0* | *67.6* | *24.1* | *36.7* |

Table 6. **Open-vocabulary segmentation** performance in mIoU (%). All 'generalist image-text encoders' methods are evaluated using MaskCLIP defined in Section 4.1. We put in gray methods specialized for segmentation and bold the sections separately. For reference, we also produce results with our high-resolution inference procedure (noted 'HR' and in italic).

## 4.5. Further Analysis

**Failure modes.** We conduct an error analysis on the open-vocabulary segmentation task, on the ADE20K dataset.

*Object boundaries.* We replace the $k$-means operation in our inference procedure by the ground-truth masks to perform open-vocabulary predictions. This leads to a perfect-boundary topline of 38.9 mIoU on ADE20K, meaning that the remaining performance gap can be attributed to the misalignment between image and text. In this case, errors include predicting "shower" where the real label is "wall", in a bathroom photo, suggesting the patch features take the context into account to some extent.

*Object overlaps.* We observed multiple cases where overlapping objects are predicted but not annotated due to overlaps. For example, in Figure 4 we can observe that the seat of the motorbike is predicted separately, while in usual annotations the whole motorbike is labeled as a single entity. This decreases the benchmark score for this class of models, solely due to the dataset collection procedure.

*Class names.* We observe during evaluation that some class names do not always correspond best to their meaning in everyday language as found in captions for image-text paired data. For example, the building class of ADE20K almost systematically corresponds to a facade; the class name "person" is rarely used in captions, while "people" is more frequent. Similarly, the "vegetation" class name in CityScapes is problematic. To account for these discrepancies, we search for the optimal class names on ADE20K by averaging the token embeddings for the ground-truth masks, and using the closest word in the embedding space. We obtain a new list of class names, that we present in appendix. This procedure can add +2.1 mIoU on ADE20K which shows that the dataset class names, chosen arbitrarily at collection time, have a strong impact on the result.

This suggests that existing datasets are not well-suited to evaluate open-vocabulary semantic segmentation: first because classes naturally overlap (windows are often included in buildings), second because the class names are not aligned with their use in common natural language.

**Quality of the text encoder.** We analyze the quality of our trained text encoder by evaluating it on text classification, clustering, reranking, and pair classification tasks in the text embedding benchmark suite MTEB [57]. Our text encoder is outperformed by CLIP's text encoder by a margin of 4.2% on average on these tasks. Moreover, when comparing dino.txt with and without two learnable blocks on top of the vision encoder, we observe that removing the blocks further decreases the performance by 3.2% on average. These results show that freezing the vision encoder might hurt the text encoder and lead to lower performance on tasks such as image-text retrieval. They also suggest that we need to find a better trade-off between

exploiting the quality of frozen vision encoder and allowing it to adapt to new data domain, left for future work.

## 5. Conclusion

We have presented a training recipe, named `dino.txt`, which aligns from scratch a text encoder to a frozen self-supervised vision model, specifically DINOv2 [18, 60], unlocking open-vocabulary abilities. The approach includes a self-supervised data curation technique with no human annotation and allows for fast training, leading to strong zero-shot classification performance, on par with the state-of-the-art. The resulting text encoder is also aligned to patch-level features, therefore providing precise dense open-vocabulary segmentation capabilities thanks to the quality of the frozen vision encoder. We also argue that classic semantic segmentation benchmarks require rethinking for open-vocabulary, as they do not allow for overlapping concepts nor finer granularity in prediction than the annotations.

## References

[1] Common Crawl. https://commoncrawl.org.

[2] Amro Abbas, Kushal Tirumala, Dániel Simig, Surya Ganguli, and Ari S Morcos. SemDeDup: Data-efficient learning at web-scale through semantic deduplication. *arXiv preprint arXiv:2303.09540*, 2023.

[3] Ahmed Abdelreheem, Ivan Skorokhodov, Maks Ovsjanikov, and Peter Wonka. SATR: Zero-Shot Semantic Segmentation of 3D Shapes. In *ICCV*, 2023.

[4] Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat, Yann LeCun, and Nicolas Ballas. Self-Supervised Learning from Images with a Joint-Embedding Predictive Architecture. *arXiv preprint arXiv:2301.08243*, 2023.

[5] Hangbo Bao, Li Dong, and Furu Wei. BEiT: BERT Pre-Training of Image Transformers. *arXiv preprint arXiv:2106.08254*, 2021.

[6] Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh Tenenbaum, and Boris Katz. ObjectNet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. *NeurIPS*, 32, 2019.

[7] Walid Bousselham, Felix Petersen, Vittorio Ferrari, and Hilde Kuehne. Grounding everything: Emerging localization properties in vision-language transformers. In *CVPR*, pages 3828–3837, 2024.

[8] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. COCO-Stuff: Thing and stuff classes in context. In *CVPR*, 2018.

[9] Mathilde Caron, Piotr Bojanowski, Julien Mairal, and Armand Joulin. Unsupervised Pre-Training of Image Features on Non-Curated Data. In *ICCV*, 2019.

[10] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised Learning of Visual Features by Contrasting Cluster Assignments. In *NeurIPS*, 2020.

[11] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021.

[12] Junbum Cha, Jonghwan Mun, and Byungseok Roh. Learning to Generate Text-grounded Mask for Open-world Semantic Segmentation from Only Image-Text Pairs. In *CVPR*, 2023.

[13] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing Web-Scale Image-Text Pre-Training To Recognize Long-Tail Visual Concepts. In *CVPR*, 2021.

[14] Xinlei Chen, Saining Xie, and Kaiming He. An Empirical Study of Training Self-Supervised Vision Transformers. *arXiv preprint arXiv:2104.02057*, 2021.

[15] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. InternVL: Scaling up Vision Foundation Models and Aligning for Generic Visual-Linguistic Tasks. In *CVPR*, pages 24185–24198, 2024.

[16] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *CVPR*, 2023.

[17] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The Cityscapes Dataset for Semantic Urban Scene Understanding. In *CVPR*, pages 3213–3223, 2016.

[18] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision Transformers Need Registers. *ICLR*, 2024.

[19] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR*, 2009.

[20] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805*, 2018.

[21] Aniket Didolkar, Andrii Zadaianchuk, Anirudh Goyal, Mike Mozer, Yoshua Bengio, Georg Martius, and Maximilian Seitzer. Zero-shot object-centric representation learning. In *arXiv preprint arXiv:2408.09162*, 2024.

[22] Jian Ding, Nan Xue, Gui-Song Xia, and Dengxin Dai. Decoupling Zero-Shot Semantic Segmentation. In *CVPR*, 2022.

[23] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *arXiv preprint arXiv:2010.11929*, 2020.

[24] Pinar Duygulu, Kobus Barnard, Joao FG de Freitas, and David A Forsyth. Object Recognition as Machine Translation: Learning a Lexicon for a Fixed Image Vocabulary. In *ECCV*, 2002.

[25] Johan Edstedt, Qiyu Sun, Georg Bökman, Mårten Wadenbäck, and Michael Felsberg. RoMa: Robust Dense Feature Matching. In *CVPR*, 2024.

[26] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results, 2012.

[27] Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander T Toshev, and Vaishaal Shankar. Data Filtering Networks. In *ICLR*, 2024.

[28] Ali Farhadi, Mohsen Hejrati, Mohammad Amin Sadeghi, Peter Young, Cyrus Rashtchian, Julia Hockenmaier, and David Forsyth. Every Picture Tells a Story: Generating Sentences from Images. In *ECCV*, 2010.

[29] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, Eyal Orgad, Rahim Entezari, Giannis Daras, Sarah Pratt, Vivek Ramanujan, Yonatan Bitton, Kalyani Marathe, Stephen Mussmann, Richard Vencu, Mehdi Cherti, Ranjay Krishna, Pang Wei Koh, Olga Saukh, Alexander Ratner, Shuran Song, Hannaneh Hajishirzi, Ali Farhadi, Romain Beaumont, Sewoong Oh, Alex Dimakis, Jenia Jitsev, Yair Carmon, Vaishaal Shankar, and Ludwig Schmidt. DataComp: In search of the next generation of multimodal datasets. In *NeurIPS*, 2024.

[30] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Scaling Open-Vocabulary Image Segmentation with Image-Level Labels. In *ECCV*, 2022.

[31] Priya Goyal, Dhruv Mahajan, Abhinav Gupta, and Ishan Misra. Scaling and Benchmarking Self-Supervised Visual Representation Learning. In *ICCV*, 2019.

[32] Priya Goyal, Mathilde Caron, Benjamin Lefaudeux, Min Xu, Pengchao Wang, Vivek Pai, Mannat Singh, Vitaliy Liptchinsky, Ishan Misra, Armand Joulin, and Piotr Bojanowski. Self-supervised Pretraining of Visual Features in the Wild. *arXiv preprint arXiv:2103.01988*, 2021.

[33] Priya Goyal, Quentin Duval, Isaac Seessel, Mathilde Caron, Mannat Singh, Ishan Misra, Levent Sagun, Armand Joulin, and Piotr Bojanowski. Vision Models Are More Robust And Fair When Pretrained On Uncurated Images Without Supervision. *arXiv preprint arXiv:2202.08360*, 2022.

[34] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised learning. In *NeurIPS*, 2020.

[35] Abhinav Gupta and Larry S Davis. Beyond Nouns: Exploiting Prepositions and Comparative Adjectives for Learning Visual Classifiers. In *ECCV*, 2008.

[36] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum Contrast for Unsupervised Visual Representation Learning. In *CVPR*, 2020.

[37] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked Autoencoders Are Scalable Vision Learners. *arXiv preprint arXiv:2111.06377*, 2021.

[38] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural Adversarial Examples. In *CVPR*, pages 15262–15271, 2021.

[39] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. OpenCLIP, 2021.

[40] iNaturalist 2021 competition dataset. `https://github.com/visipedia/inat_comp/tree/master/2021`, 2021.

[41] Krishna Murthy Jatavallabhula, Alihusein Kuwajerwala, Qiao Gu, Mohd Omama, Tao Chen, Shuang Li, Ganesh Iyer, Soroush Saryazdi, Nikhil Keetha, Ayush Tewari, Joshua B. Tenenbaum, Celso Miguel de Melo, Madhava Krishna, Liam Paull, Florian Shkurti, and Antonio Torralba. ConceptFusion: Open-set Multimodal 3D Mapping. In *RSS*, 2023.

[42] Armand Joulin, Laurens Van Der Maaten, Allan Jabri, and Nicolas Vasilache. Learning Visual Features from Large Weakly Supervised Data. In *ECCV*, 2016.

[43] Dahyun Kang and Minsu Cho. In Defense of Lazy Visual Grounding for Open-Vocabulary Semantic Segmentation. In *ECCV*, 2024.

[44] Laurynas Karazija, Iro Laina, Andrea Vedaldi, and Christian Rupprecht. Diffusion Models for Open-Vocabulary Segmentation. *arXiv preprint arXiv:2306.09316*, 2023.

[45] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. LERF: Language Embedded Radiance Fields. In *ICCV*, 2023.

[46] Girish Kulkarni, Visruth Premraj, Vicente Ordonez, Sagnik Dhar, Siming Li, Yejin Choi, Alexander C Berg, and Tamara L Berg. BabyTalk: Understanding and Generating Simple Image Descriptions. *IEEE TPAMI*, 35(12):2891–2903, 2013.

[47] Mengcheng Lan, Chaofeng Chen, Yiping Ke, Xinjiang Wang, Litong Feng, and Wayne Zhang. ProxyCLIP: Proxy Attention Improves CLIP for Open-Vocabulary Segmentation. In *ECCV*, 2024.

[48] Samuel Lavoie, Polina Kirichenko, Mark Ibrahim, Mahmoud Assran, Andrew Gordon Wilson, Aaron Courville, and Nicolas Ballas. Modeling Caption Diversity in Contrastive Vision-Language Pretraining. *arXiv preprint arXiv:2405.00740*, 2024.

[49] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and Rene Ranftl. Language-driven Semantic Segmentation. In *ICLR*, 2022.

[50] Dylan Li and Gyungin Shin. ProMerge: Prompt and Merge for Unsupervised Instance Segmentation. In *ECCV*, 2024.

[51] 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. In *CVPR*, 2023.

[52] 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 *ECCV*, 2022.

[53] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X. Yu. Large-Scale Long-Tailed Recognition in an Open World. In *CVPR*, 2019.

[54] Huaishao Luo, Junwei Bao, Youzheng Wu, Xiaodong He, and Tianrui Li. SegCLIP: Patch aggregation with learnable centers for open-vocabulary semantic segmentation. In *ICML*, 2023.

[55] George A Miller. WordNet: a lexical database for English. *Communications of the ACM*, 38(11):39–41, 1995.

[56] 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 *CVPR*, 2014.

[57] Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. MTEB: Massive Text Embedding Benchmark. *arXiv preprint arXiv:2210.07316*, 2022.

[58] Jishnu Mukhoti, Tsung-Yu Lin, Omid Poursaeed, Rui Wang, Ashish Shah, Philip HS Torr, and Ser-Nam Lim. Open Vocabulary Semantic Segmentation with Patch Aligned Contrastive Learning. In *CVPR*, 2023.

[59] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

[60] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision. *TMLR*, 2024.

[61] Evin Pınar Örnek, Yann Labbé, Bugra Tekin, Lingni Ma, Cem Keskin, Christian Forster, and Tomas Hodan. FoundPose: Unseen Object Pose Estimation with Foundation Features. In *ECCV*, pages 163–182, 2025.

[62] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *NeurIPS*, pages 8024–8035. 2019.

[63] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k Entities: Collecting Region-to-Phrase Correspondences for Richer Image-to-Sentence Models. In *ICCV*, pages 2641–2649, 2015.

[64] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning Transferable Visual Models From Natural Language Supervision. In *ICML*, pages 8748–8763, 2021.

[65] Kanchana Ranasinghe, Brandon McKinzie, Sachin Ravi, Yinfei Yang, Alexander Toshev, and Jonathon Shlens. Perceptual Grouping in Contrastive Vision-Language Models. In *ICCV*, 2023.

[66] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou, and Jiwen Lu. DenseCLIP: Language-Guided Dense Prediction with Context-Aware Prompting. In *CVPR*, 2022.

[67] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do ImageNet classifiers generalize to ImageNet? In *ICML*, pages 5389–5400. PMLR, 2019.

[68] 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. In *ICLR*, 2023.

[69] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. LAION-400M: Open Dataset of CLIP-Filtered 400 Million Image-Text Pairs. *arXiv preprint arXiv:2111.02114*, 2021.

[70] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5B: An open large-scale dataset for training next generation image-text models. In *NeurIPS*, pages 25278–25294, 2022.

[71] Gyungin Shin, Weidi Xie, and Samuel Albanie. ReCo: Retrieve and Co-segment for Zero-shot Transfer. *NeurIPS*, 2022.

[72] Gyungin Shin, Weidi Xie, and Samuel Albanie. NamedMask: Distilling Segmenters from Complementary Foundation Models. In *CVPR*, pages 4961–4970, 2023.

[73] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. EVA-CLIP: Improved Training Techniques for CLIP at Scale. *arXiv preprint arXiv:2303.15389*, 2023.

[74] Quan Sun, Jinsheng Wang, Qiying Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, and Xinlong Wang. EVA-CLIP-18B: Scaling CLIP to 18 Billion Parameters. *arXiv preprint arXiv:2402.04252*, 2024.

[75] Jamie Tolan, Hung-I Yang, Ben Nosarzewski, Guillaume Couairon, Huy V. Vo, John Brandt, Justine Spore, Sayantan Majumdar, Daniel Haziza, Janaki Vamaraju, Theo Moutakanni, Piotr Bojanowski, Tracy Johns, Brian White, Tobias Tiecke, and Camille Couprie. Sub-meter resolution canopy height maps using self-supervised learning and a vision transformer trained on Aerial and GEDI Lidar. *arXiv preprint arXiv:2304.07213*, 2023.

[76] Tsung-Yi, Genevieve Patterson, Matteo R. Ronchi, Yin Cui, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays Georgia, Pietro Perona, Deva Ramanan, Larry Zitnick, and Piotr Dollár. COCO 2017: Common Objects in Context 2017, 2017.

[77] Narek Tumanyan, Assaf Singer, Shai Bagon, and Tali Dekel. DINO-Tracker: Taming DINO for Self-Supervised Point Tracking in a Single Video. In *ECCV*, 2024.

[78] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. In *NeurIPS*, 2017.

[79] Huy V. Vo, Vasil Khalidov, Timothée Darcet, Théo Moutakanni, Nikita Smetanin, Marc Szafraniec, Hugo Touvron, Camille Couprie, Maxime Oquab, Armand Joulin, Herve Jegou, Patrick Labatut, and Piotr Bojanowski. Automatic Data Curation for Self-Supervised Learning: A Clustering-Based Approach. *TMLR*, 2024.

[80] Feng Wang, Jieru Mei, and Alan Yuille. SCLIP: Rethinking Self-Attention for Dense Vision-Language Inference. In *ECCV*, pages 315–332, 2024.

[81] Haoxiang Wang, Pavan Kumar Anasosalu Vasu, Fartash Faghri, Raviteja Vemulapalli, Mehrdad Farajtabar, Sachin Mehta, Mohammad Rastegari, Oncel Tuzel, and Hadi Pouransari. SAM-CLIP: Merging Vision Foundation Models towards Semantic and Spatial Understanding. In *CVPR*, 2024.

[82] Jianzong Wu, Xiangtai Li, Shilin Xu, Haobo Yuan, Henghui Ding, Yibo Yang, Xia Li, Jiangning Zhang, Yunhai Tong, Xudong Jiang, Bernard Ghanem, and Dacheng Tao. Towards Open Vocabulary Learning: A Survey. *IEEE TPAMI*, 2024.

[83] Ji-Jia Wu, Andy Chia-Hao Chang, Chieh-Yu Chuang, Chun-Pei Chen, Yu-Lun Liu, Min-Hung Chen, Hou-Ning Hu, Yung-Yu Chuang, and Yen-Yu Lin. Image-Text Co-Decomposition for Text-Supervised Semantic Segmentation. *arXiv preprint arXiv:2404.04231*, 2024.

[84] Size Wu, Wenwei Zhang, Lumin Xu, Sheng Jin, Xiangtai Li, Wentao Liu, and Chen Change Loy. CLIPSelf: Vision Transformer Distills Itself for Open-Vocabulary Dense Prediction. *arXiv preprint arXiv:2310.01403*, 2023.

[85] Monika Wysoczańska, Michaël Ramamonjisoa, Tomasz Trzciński, and Oriane Siméoni. CLIP-DIY: CLIP Dense Inference Yields Open-Vocabulary Semantic Segmentation For-Free. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2024.

[86] Monika Wysoczańska, Oriane Siméoni, Michaël Ramamonjisoa, Andrei Bursuc, Tomasz Trzciński, and Patrick Pérez. CLIP-DINOiser: Teaching CLIP a few DINO tricks for Open-Vocabulary Semantic Segmentation. *ECCV*, 2024.

[87] Yun Xing, Jian Kang, Aoran Xiao, Jiahao Nie, Ling Shao, and Shijian Lu. Rewrite Caption Semantics: Bridging Semantic Gaps for Language-Supervised Semantic Segmentation. In *NeurIPS*, pages 68798–68809, 2023.

[88] Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Demystifying CLIP data. In *ICLR*, 2024.

[89] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. GroupViT: Semantic Segmentation Emerges from Text Supervision. *CVPR*, 2022.

[90] Jilan Xu, Junlin Hou, Yuejie Zhang, Rui Feng, Yi Wang, Yu Qiao, and Weidi Xie. Learning Open-vocabulary Semantic Segmentation Models From Natural Language Supervision. In *CVPR*, 2023.

[91] Andrii Zadaianchuk, Maximilian Seitzer, and Georg Martius. Object-Centric Learning for Real-World Videos by Predicting Temporal Feature Similarities. *NeurIPS*, 36, 2024.

[92] Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas Steiner, Daniel Keysers, Alexander Kolesnikov, and Lucas Beyer. LiT: Zero-Shot Transfer with Locked image text Tuning. In *CVPR*, pages 18123–18133, 2022.

[93] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid Loss for Language Image Pre-Training. In *ICCV*, pages 11975–11986, 2023.

[94] Yiwu Zhong, Jianwei Yang, Pengchuan Zhang, Chunyuan Li, Noel Codella, Liunian Harold Li, Luowei Zhou, Xiyang Dai, Lu Yuan, Yin Li, and Jianfeng Gao. RegionCLIP: Region-based Language-Image Pretraining. In *CVPR*, 2022.

[95] Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. Learning Deep Features for Scene Recognition using Places Database. In *NeurIPS*, 2014.

[96] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic Understanding of Scenes through the ADE20K Dataset. *IJCV*, 2019.

[97] Chong Zhou, Chen Change Loy, and Bo Dai. Extract Free Dense Labels from CLIP. In *ECCV*, 2022.

## A. Additional details

**Computational cost.** In this work, we show that we can obtain good global and local vision-language alignment with minimal additional cost thanks to powerful pre-trained SSL models. This appears to be a more efficient paradigm than training CLIP from scratch. The computational costs for training our models and different CLIP models are reported in Table 7. For completeness' sake, we also include the pretraining cost of the ViT-g DINOv2 vision encoder as well as the cost of distilling this model into a ViT-L. In practice, such additional costs should however be considered amortized over the multiple downstream adaptations of the DINOv2 backbone.

| Method | Samples seen | Batch size | GPUs | total GPU.h | GPU arch. |
|---|---|---|---|---|---|
| CLIP | 12.8B | 32768 | 256 | 73728 | V100 |
| OpenCLIP | 12.8B | 38400 | 400 | 50800 | A100 40 GB |
| MetaCLIP | 12.8B | 32768 | 128 | 92160 | V100 |
| EVA-02-CLIP | 2B | 61000 | 128 | – | A100 40 GB |
| DINOv2 ViT-g pretraining | – | – | 256 | 22000 | A100 80 GB |
| DINOv2 ViT-L distillation | – | – | – | 8000 | A100 80 GB |
| dino.txt | 3.2B | 65536 | 128 | 2432 | A100 80 GB |
| dino.txt @336 | 3.2B | 65536 | 256 | 4096 | A100 80 GB |

Table 7. **Computational cost of different models in GPU hours.**

**ADE20K class names for the error analysis discussion.** In Section 4.5, we discuss the failure modes of our zero-shot semantic segmentation method. In particular, we show that class names can be optimized to boost results, instead of using the default ones from each dataset. This is not surprising, the 150 class names of ADE20K were originally chosen to identify each category and were not intended as holistic descriptors for zero-shot segmentation via a vision-language model. In our experiments, we have observed that some class names are too broad, *e.g.*, *building*, or ambiguous, *e.g.*, *throne*, and consequently result in incorrect predictions. In Table 10, we include the optimized class names for ADE20K that improve open-vocabulary segmentation by 2.1 mIoU points, as reported in the discussion about failure modes in Section 4.5. Please note that for all experiments in the main text, we use the original class names to facilitate comparison with previous work.

## B. Additional ablation studies

| Inference embedding | segmentation ADE | City. |
|---|---|---|
| [value] (MCLIP) | 7.0 | 11.7 |
| [CLS patch] | 19.9 | 26.2 |
| [value patch] | 20.0 | 29.0 |
| [patch] | **20.6** | **32.1** |

Table 8. **Ablation of the embedding in dense zero-shot segmentation inference.** We show segmentation results with different embeddings to represent a patch, on the datasets ADE20K and Cityscapes. 'MCLIP' corresponds to MaskCLIP [97] strategy, which we also name here value.

**Impact of the embedding in segmentation.** Table 8 presents open-vocabulary segmentation results on the challenging datasets ADE20K and Cityscapes. We follow the evaluation protocol of TCL [12]. Following only MaskCLIP patch representation ([value]) leads to the worst results. Using solely the model's output patch descriptor ([patch]) and their corresponding part in the text embedding leads to the best results. This is the setup used in the main paper. We also observe that concatenating the [CLS] token to the patch representation hurts the performance *vs.* [patch] only, particularly in Cityscapes: we found this to be due to the dominance of the salient visual concept in the [CLS].

**Impact of the image embedding size at training.** We show in Table 9 that the benefit of using the concatenated representation **g** (noted here [CLS avg]) when training dino.txt does not come from higher dimensionality of the image embedding. To this end, we have conducted an additional experiment in which we project the [CLS] token from the dimension of 1024 to 2048 before passing it to the vision blocks. Little impact is observed from this dimensionality change. This additionally shows that the gain (from 30.9 to 34.7) in the retrieval task is largely due to the concatenation of the [CLS] token with [avg].

| Training embedding | proj | *class.* IN1K | *retr.* COCO |
|---|---|---|---|
| [CLS] | | 78.8 | 30.2 |
| [CLS] | $1024 \rightarrow 2048$ | 78.8 | 30.9 |
| [CLS avg] | | **79.2** | **34.7** |

Table 9. **Analysis of the image embedding size at training time.** Projecting the [CLS] embedding to dimension 2048 (second row) yields minimal gain on the benchmarks.

## C. Additional qualitative results

**Open-vocabulary semantic segmentation.** Figures 5-6 demonstrate that the segmentation results of dino.txt with images and texts in the wild. For each image, we select a small number of descriptive text prompts and run the zero-shot semantic segmentation pipeline described in Section 4.4. Our model is able to segment complex scenes with multiple semantic objects and specific text inputs, *e.g.*, "pesto bruschetta" and "nautical rope".

| Original | Optimized | | Original | Optimized |
|---|---|---|---|---|
| wall | wall | | swivel chair | swivel chair |
| building, edifice | **facade, frontage, frontal** | | boat | boat |
| sky | sky | | bar | bar |
| floor, flooring | **floor** | | arcade machine | arcade machine |
| tree | tree | | hovel, hut, hutch, shack, shanty | **hovel** |
| ceiling | ceiling | | bus, autobus, coach, charabanc, double-decker, jitney, motorbus, motorcoach, omnibus, passenger vehicle | **bus** |
| road, route | **road** | | towel | towel |
| bed | bed | | light, light source | **skylight, fanlight** |
| windowpane, window | **windowpane** | | truck, motortruck | **truck** |
| grass | grass | | tower | tower |
| cabinet | cabinet | | chandelier, pendant, pendent | **chandelier** |
| sidewalk, pavement | sidewalk, pavement | | awning, sunshade, sunblind | **awning** |
| person, individual, someone, somebody, mortal, soul | **people** | | streetlight, street lamp | **streetlight** |
| earth, ground | **ground, earth** | | booth, cubicle, stall, kiosk | **newsstand** |
| door, double door | **interior door** | | television receiver, television, television set, tv, tv set, idiot box, boob tube, telly, goggle box | **television receiver** |
| table | table | | airplane, aeroplane, plane | **airplane** |
| mountain, mount | **mountain** | | dirt track | dirt track |
| plant, flora, plant life | **plant** | | apparel, wearing apparel, dress, clothes | **clothes closet, clothespress** |
| curtain, drape, drapery, mantle, pall | **curtain** | | pole | pole |
| chair | chair | | land, ground, soil | **land** |
| car, auto, automobile, machine, motorcar | **car** | | banister, banister, balustrade, balusters, handrail | banister, banister, balustrade, balusters, handrail |
| water | water | | escalator, moving staircase, moving stairway | **escalator** |
| painting, picture | **painting** | | ottoman, pouf, pouffe, puff, hassock | **footstool, footrest, ottoman, tuffet** |
| sofa, couch, lounge | sofa, couch, lounge | | bottle | bottle |
| shelf | shelf | | buffet, counter, sideboard | **china cabinet, china closet** |
| house | house | | poster, posting, placard, notice, bill, card | **poster** |
| sea | sea | | stage | stage |
| mirror | mirror | | van | van |
| rug, carpet, carpeting | **rug** | | ship | ship |
| field | field | | fountain | fountain |
| armchair | armchair | | conveyer belt, conveyor belt, conveyer, conveyor, transporter | **conveyer belt** |
| seat | seat | | canopy | **baldachin** |
| fence, fencing | **fence** | | washer, automatic washer, washing machine | **washer** |
| desk | desk | | plaything, toy | **plaything** |
| rock, stone | **rock** | | swimming pool, swimming bath, natatorium | **swimming pool** |
| wardrobe, closet, press | **wardrobe** | | stool | stool |
| lamp | lamp | | barrel, cask | **barrel** |
| bathtub, bathing tub, bath, tub | **bathtub** | | basket, handbasket | **basket** |
| railing, rail | **railing** | | waterfall, falls | **waterfall** |
| cushion | **pillow** | | tent, collapsible shelter | **tent** |
| base, pedestal, stand | **stall, stand, sales booth** | | bag | bag |
| box | box | | minibike, motorbike | **motorcycle, bike** |
| column, pillar | **column** | | cradle | **baby bed, baby's bed** |
| signboard, sign | **signboard** | | oven | oven |
| chest of drawers, chest, bureau, dresser | **chest of drawers** | | ball | ball |
| counter | **reception desk** | | food, solid food | **food** |
| sand | sand | | step, stair | **pedestal, plinth, footstall** |
| sink | sink | | tank, storage tank | **tank** |
| skyscraper | skyscraper | | trade name, brand name, brand, marque | **trade name** |
| fireplace, hearth, open fireplace | fireplace, hearth, open fireplace | | microwave, microwave oven | **microwave** |
| refrigerator, icebox | **refrigerator** | | pot, flowerpot | **pot** |
| grandstand, covered stand | **grandstand** | | animal, animate being, beast, brute, creature, fauna | **animal** |
| path | path | | bicycle, bike, wheel, cycle | **bicycle** |
| stairs, steps | **stairs** | | lake | lake |
| runway | runway | | dishwasher, dish washer, dishwashing machine | **dishwasher** |
| case, display case, showcase, vitrine | case, display case, showcase, vitrine | | screen, silver screen, projection screen | **screen** |
| pool table, billiard table, snooker table | **pool table** | | blanket, cover | **blanket** |
| pillow | **pillow sham** | | sculpture | sculpture |
| screen door, screen | **shower** | | hood, exhaust hood | **range hood** |
| stairway, staircase | **stairway** | | sconce | sconce |
| river | river | | vase | vase |
| bridge, span | **bridge** | | traffic light, traffic signal, stoplight | **traffic light** |
| bookcase | bookcase | | tray | tray |
| blind, screen | **blind** | | ashcan, trash can, garbage can, wastebin, ash bin, ash-bin, ashbin, dustbin, trash barrel, trash bin | ashcan, trash can, garbage can, wastebin, ash bin, ash-bin, ashbin, dustbin, trash barrel, trash bin |
| coffee table, cocktail table | **coffee table** | | fan | fan |
| toilet, can, commode, crapper, pot, potty, stool, throne | **toilet** | | pier, wharf, wharfage, dock | **pier** |
| flower | flower | | crt screen | crt screen |
| book | book | | plate | **plate, collection plate** |
| hill | **hillside** | | monitor, monitoring device | **computer screen, computer display** |
| bench | bench | | bulletin board, notice board | **bulletin board** |
| countertop | countertop | | shower | shower |
| stove, kitchen stove, range, kitchen range, cooking stove | stove, kitchen stove, range, kitchen range, cooking stove | | radiator | radiator |
| palm, palm tree | **cabbage palm, cabbage tree, Livistona australis** | | glass, drinking glass | **glass** |
| kitchen island | kitchen island | | clock | clock |
| computer, computing machine, computing device, data processor, electronic computer, information processing system | **desktop computer** | | flag | flag |

Table 10. **ADE20K dataset:** original class names *vs.* optimized class names for zero-shot semantic segmentation. Modified class names are highlighted in bold. The new class names have been picked through manual analysis to increase specificity, for example *building* to *facade*, or to remove potential confusion, for example *throne* for *toilet*.

| Color | Name |
|---|---|
| | Wine glass |
| | Wine bottle |
| | Stone cutting board |
| | Cherry tomatoes |
| | White ceramic bowl |
| | French cheese |
| | Salami slices |
| | Wooden table |
| | Sliced baguette |
| | Green grapes |
| | Pesto bruschetta |
| | Red pepper spread on bread |

| Color | Name |
|---|---|
| | Window |
| | White cabinet |
| | Black television screen |
| | Wooden sofa table |
| | Gray couch |
| | Candle |
| | Potted plant |
| | Books |
| | Indoor wall |
| | Parquet floor |

| Color | Name |
|---|---|
| | Red pickup truck |
| | Stone wall |
| | Lush tree |
| | Bush |
| | Paved road |
| | Chair |
| | Facade |
| | Blue sky |

Figure 5. **Open-vocabulary semantic segmentation, part 1/2.** The input resolution is 896×896 pixels.

15

| Color | Name |
|---|---|
| (blue) | Tall giraffe |
| (orange) | Blue automobile |
| (green) | Tanned man in shirt and pants |
| (red) | Open sky |
| (purple) | Trees, bushes |
| (brown) | Dirt road, sandy ground |
| (pink) | Wood railing, fence |

| Color | Name |
|---|---|
| (blue) | Wood rowing canoe |
| (orange) | Inflatable motor boat |
| (green) | Peaceful lake |
| (red) | Wooden pier |
| (purple) | Bush |
| (brown) | Blue sky |
| (pink) | Tree |
| (gray) | Nautical rope |

| Color | Name |
|---|---|
| (blue) | Pedestrian |
| (orange) | Tram |
| (green) | Car |
| (red) | Electric wires |
| (purple) | Facade |
| (brown) | Window |
| (pink) | Open sky |
| (gray) | Road, pavement |

Figure 6. **Open-vocabulary semantic segmentation, part 2/2.** The input resolution is 896×896 pixels.