# CLIP-S$^4$: Language-Guided *Self-Supervised Semantic Segmentation*

Wenbin He      Suphanut Jamonnak      Liang Gou      Liu Ren

Bosch Research North America & Bosch Center for Artificial Intelligence (BCAI)

{wenbin.he2, suphanut.jamonnak, liang.gou, liu.ren}@us.bosch.com

## Abstract

*Existing semantic segmentation approaches are often limited by costly pixel-wise annotations and predefined classes. In this work, we present CLIP-S$^4$ that leverages self-supervised pixel representation learning and vision-language models to enable various semantic segmentation tasks (e.g., unsupervised, transfer learning, language-driven segmentation) without any human annotations and unknown class information. We first learn pixel embeddings with **pixel-segment contrastive learning** from different augmented views of images. To further improve the pixel embeddings and enable language-driven semantic segmentation, we design two types of consistency guided by vision-language models: 1) **embedding consistency**, aligning our pixel embeddings to the joint feature space of a pre-trained vision-language model, CLIP [34]; and 2) **semantic consistency**, forcing our model to make the same predictions as CLIP over a set of carefully designed target classes with both known and unknown prototypes. Thus, CLIP-S$^4$ enables a new task of class-free semantic segmentation where no unknown class information is needed during training. As a result, our approach shows consistent and substantial performance improvement over four popular benchmarks compared with the state-of-the-art unsupervised and language-driven semantic segmentation methods. More importantly, our method outperforms these methods on unknown class recognition by a large margin.*

## 1. Introduction

Semantic segmentation aims to partition an input image into semantically meaningful regions and assign each region a semantic class label. Recent advances in semantic segmentation [6, 27, 48] heavily rely on pixel-wise human annotations, which have two limitations. First, acquiring pixel-wise annotations is extremely labor intensive and costly, which can take up to 1.5 hours to label one image [31]. Second, human annotations are often limited to a set of predefined semantic classes, with which the learned models lack the ability to recognize unknown classes [25].
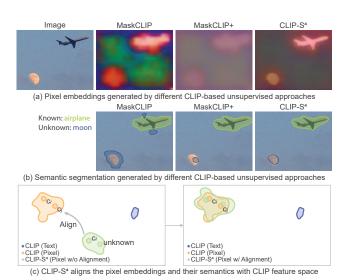


Figure 1. (a) Pixel embeddings from different CLIP-based unsupervised methods: Our method, CLIP-S$^4$, generates sharper and more coherent pixel embeddings than MaskCLIP [49] and MaskCLIP+'s [49]; (b) Language-driven semantic segmentation by different methods: CLIP-S$^4$ can recognize challenging unknown classes (e.g., moon); (c) The key idea behind CLIP-S$^4$: aligning the pixel embeddings and their semantics with CLIP feature space.

Various approaches have been proposed to tackle these limitations, among which we are inspired by two lines of recent research in particular. First, for unsupervised semantic segmentation (i.e., without human annotations), self-supervised pixel representation learning approaches [14, 18, 19, 23, 40] have shown promising results on popular unsupervised benchmarks. The main idea is to extend self-supervised contrastive learning [7, 16] from images to pixels by attracting each pixel's embedding to its positive pairs and repelling it from negative pairs. The prior of pairs can be contours [18, 19], hierarchical groups [23], salience maps [40], and pre-trained models [14]. Although these approaches can group pixels into semantically meaningful clusters, human annotations are still needed to assign class labels to the clusters for semantic segmentation [37].

Second, for unknown classes in semantic segmentation, large-scale vision-language models such as CLIP [34]

have shown great potential. This line of research, called *language-driven semantic segmentation*, aims to segment images with arbitrary classes defined by texts during testing time [25, 37, 46, 49]. Among these methods, most still need training time annotations, such as pixel annotations [25] and captions [46]. Only a few recent work, MaskCLIP & MaskCLIP+ [49] attempts to address this without using additional supervision: MaskCLIP directly extracts pixel embeddings correlated with texts from CLIP, but these pixel embeddings are coarse and noisy (Fig. 1a). To address this issue, MaskCLIP+ [49] trains a segmentation model on the pseudo-labels generated by MaskCLIP for a set of predefined classes. However, the pixel embeddings of MaskCLIP+ are distorted by the predefined classes (Fig. 1a), which limits its ability to recognize unknowns (Fig. 1b). Also, it needs unknown class information during training, which hinders its real-world applications.

We propose a language-guided self-supervised semantic segmentation approach, CLIP-S$^4$, which takes advantage of the strengths from both lines of research and addresses their limitations accordingly. The key idea is to learn consistent pixel embeddings with respect to visual and conceptual semantics using self-supervised learning and the guidance of a vision-language model, CLIP.

Specifically, we first train pixel embeddings with *pixel-segment contrastive learning* from different augmented image views [18, 19, 23] such that images can be partitioned into visually meaningful regions. To further improve pixel embedding quality and enable language-driven semantic segmentation, we introduce vision-language model guided consistency to regularize our model (Fig. 1c). The consistency is enforced from two aspects: *embedding consistency* and *semantic consistency*. First, embedding consistency aims to align the pixel embeddings generated by our model with the joint feature space of texts and images of CLIP by minimizing the distance between the pixel embeddings generated by our model and CLIP. Second, semantic consistency forces our model to make the same prediction as CLIP over a set of carefully designed target classes with both *known* and *unknown* prototypes. Note that unlike the previous methods [25, 49] that use a predefined set of *known classes*, CLIP-S$^4$ also learns the representation of *unknown classes* from images during training.

In the end, CLIP-S$^4$ also enables a new task, namely *class-free semantic segmentation*, as shown in Tab. 1. This new task does not need any human annotations and even assumes NO class names are given during training. This is a more challenging task than the recent work [49] that requires class names of both known and unknown.

In summary, the contributions of this paper are threefold:

- We propose a self-supervised semantic segmentation approach that combines pixel-segment contrastive learning with the guidance of pre-trained vision lan-

| | Known | | Unknown | | |
| | Annot. | Cls Name | Annot. | Cls Name | Add. Info. |
|---|---|---|---|---|---|
| Un/Self-supervised ( [19] etc.) | ✗ | ✗ | ✗ | ✗ | Fine-Tuning |
| Supervised ( [27] etc.) | ✓ | ✓ | N/A | N/A | N/A |
| Zero-shot ( [3] etc.) | ✓ | ✓ | ✗ | ✓ | Word2Vec, etc. |
| Language- *MaskCLIP+* [49] | ✗ | ✓ | ✗ | ✓ | CLIP |
| Driven *CLIP-S$^4$* | ✗ | ✓ | ✗ | ✗ | CLIP |

Table 1. Comparison of information required for training over different tasks. CLIP-S$^4$ enables a new task called *class-free semantic segmentation*. Compared with MaskCLIP+ [49], the new task assumes unknown class names are NOT given during training.

guage models. Our method can generate high-quality pixel embeddings without any human annotations and be applied to a variety of semantic segmentation tasks.

- We open up new research potentials for language-driven semantic segmentation without any human annotations by introducing and addressing a new task of *class-free semantic segmentation* (Tab. 1). Unlike previous work that assumes all the class names are known during training, our method can discover unknown classes from unlabelled image data without even knowing unknown class names.

- Consistent and substantial gains are observed with our approach over the state-of-the-art unsupervised and language-driven semantic segmentation methods on four popular datasets. More importantly, our method significantly outperforms the state-of-the-art on the segmentation of unknown classes.

## 2. Related Work

**Unsupervised Semantic Segmentation.** There are two groups of recent unsupervised semantic segmentation methods. One group of methods learns to generate consistent pixel representations or predictions between different augmentation of images with the guidance of mutual information [22, 30], clusters [8], contours [18, 19], hierarchical groups [23, 47], and saliency masks [40]. The other group of methods extracts dense features from pre-trained models based on saliency maps [36], augmentations [39], spectral decomposition [28], and feature correspondences [14]. While these methods can generate pixel embeddings with semantically meaningful clusters, annotations are needed to assign class labels to the clusters (e.g., $k$-nearest neighbor search [19] and Hungarian algorithm [40]). Our work combines pixel-segment self-supervision with pre-trained vision-language models to enable semantic segmentation without any human annotations.

**Language-Driven Semantic Segmentation.** Recently, vision-language models (e.g., CLIP [34]) trained on large-
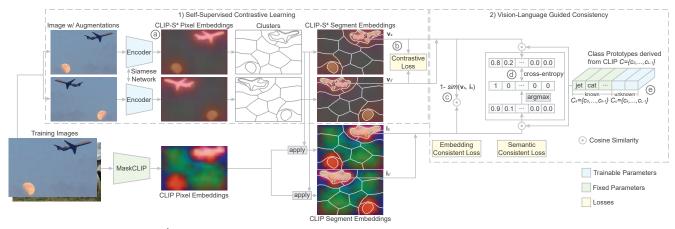
Figure 2. Framework of CLIP-S$^4$. (a) It starts with an encoder that maps images into pixel embeddings for semantic segmentation. Then it follows with two components: 1) *self-supervised contrastive learning* and 2) *vision-language model guided consistency*. Specifically, shown in (b), the self-supervised contrastive learning forces pixel embeddings to be consistent within visually coherent regions and among different augmented views of the same image. For vision-language model guided consistency, this framework introduces (c) *embedding consistency* that aligns the pixel embeddings generated by our model with CLIP embeddings, and (d) *semantic consistency* that forces our model to make the same predictions as CLIP for a set of *target classes* with both known and unknown prototypes. For known classes, the prototypes are pre-computed and fixed during training. For unknown classes, the prototypes are learned during training via clustering ((e)).

scale image-text datasets have shown great potential on various downstream tasks such as image synthesis [21, 41], out-of-distribution detection [11], and object detection [13]. To extend vision-language models for semantic segmentation, one active research, *language-driven semantic segmentation*, aims to segment images with arbitrary unknown classes defined by texts during testing time [25, 37, 46, 49]. Some methods [25, 44] use pixel-wise annotations to train language-guided semantic segmentation models. Other methods [10, 46] perform large-scale pre-training on image-text pairs specifically for semantic segmentation.

By contrast, we directly use vision-language models that are pre-trained for classification tasks. Along this line of research, a few approaches have been proposed [35, 37, 49]. The most relevant approach to our method is MaskCLIP [49], which extends the embeddings generated by pre-trained vision-language models from image to pixel level, but these embeddings are often coarse and noisy. To address this issue, MaskCLIP+ [49] fine-tunes the pixel embeddings by the pseudo-labels of a specific set of classes on top of MaskCLIP. However, it needs unknown class names during training, which may not be possible in real-world cases. Compared with [49], our method can recognize unknown classes without knowing any unknown class information during training time, and also learns fine-grained and sharper pixel embeddings with self-supervision.

# 3. Method

Our method (Fig. 2) segments images by learning a pixel embedding function with self-supervised contrastive learning and the guidance of a pre-trained vision-language model, CLIP. We use self-supervised contrastive learning

to force the consistency of pixel embeddings within visually coherent regions (e.g., superpixels) and among different augmented views of the same image (Sec. 3.1). We also introduce two vision-language model guided consistency (i.e., *embedding consistency* and *semantic consistency*) to further regularize the model (Sec. 3.2). The two components are complementary to each other. On the one hand, contrastive learning mitigates the noise introduced by CLIP. On the other hand, with the knowledge extracted from CLIP, the quality of the pixel embeddings can be improved. More importantly, this approach enables us to perform language-driven semantic segmentation with our carefully designed *target class prototypes* of both knowns and unknowns. In the following, we discuss the two components in detail.

## 3.1. Pixel-Segment Contrastive Learning

We train a pixel embedding function to generate consistent pixel embeddings within visually coherent regions through pixel-segment contrastive learning [18, 23]. Specifically, the embedding function transforms each pixel $p$ of an image to a unit-length embedding vector $\mathbf{z}_p$ of dimension $d$ via a deep neural network. The image is then partitioned into $|\mathcal{S}|$ segments by clustering the pixel embeddings. The embedding $\mathbf{v}_s$ of each segment $s$ is calculated as the average of the pixel embeddings $\mathbf{v}_s = \sum_{p \in s} \mathbf{z}_p / |s|$, which is also normalized into a unit-length vector $\mathbf{v}_s = \mathbf{v}_s / \|\mathbf{v}_s\|$. For each pixel $p$, the segments are grouped into two sets including a positive set $\mathcal{S}^+$ and a negative set $\mathcal{S}^-$. The positive set $\mathcal{S}^+$ of a pixel contains segments within the same visually coherent region of the pixel. Following the prior work [18, 23], the visually coherent region can be derived from super-pixels [1] or contours [2]. We also use data augmentation (e.g., random cropping and color jitter-

3

ing) to generate consistent pixel embeddings between different augmented views of the same image. Hence, segments within the same region of the pixel in any augmented views are considered as the positive set $\mathcal{S}^+$. Other segments in the image and segments from other images in the same batch are included in the negative set $\mathcal{S}^-$. The pixel embedding $\mathbf{z}_p$ is then attracted to the segments in positive set $\mathcal{S}^+$ and repelled from the segments in negative set $\mathcal{S}^-$ with *contrastive loss*:

$$\mathcal{L}_t(p) = -log\frac{\sum_{s\in\mathcal{S}^+} exp(sim(\mathbf{z}_p, \mathbf{v}_s)\kappa)}{\sum_{s\in\mathcal{S}^+\cup\mathcal{S}^-} exp(sim(\mathbf{z}_p, \mathbf{v}_s)\kappa)}, \quad (1)$$

where $\kappa$ is the concentration constant and $sim(\mathbf{z}_p, \mathbf{v}_s)$ is the cosine similarity between the pixel embedding $\mathbf{z}_p$ and the segment embedding $\mathbf{v}_s$.

## 3.2. Vision-Language Model Guided Consistency

To enable language-driven semantic segmentation and improve the quality of pixel embeddings, we use a pre-trained vision-language model such as CLIP [34] to guide the training of the pixel embedding function. The key idea is to align the output space of our pixel embedding function consistent with the feature space of CLIP (Fig. 1c). Specifically, two types of consistency are considered during training including *embedding consistency* and *semantic consistency*, which are detailed as follows.

**Embedding Consistency.** Our goal is to align the pixel embeddings generated from our self-supervised method (the green contour in Fig. 1c) with CLIP's pixel embeddings (the orange contour in Fig. 1c). This is done by minimizing the distance between the two pixel embedding spaces.

We first obtain the pixel embeddings of an input image from CLIP by modifying the attention-based pooling layer of the CLIP image encoder following [49]. Specifically, we 1) remove the query and key projection layers and 2) reformulate the value projection layer and the last linear layer as two consecutive fully connected layers. In the following, we use $clip$-$i(\cdot)$ as the modified CLIP image encoder and $clip$-$t(\cdot)$ as CLIP text encoder.

Then we obtain the pixel embeddings of CLIP for different augmented views of the image. Note that we use the original image to generate the CLIP pixel embeddings and perform augmentation afterwards to make sure that the CLIP pixel embeddings are consistent among different augmented views. In the end, we minimize the distance of embeddings between **segments** instead of pixels from our self-supervised and CLIP embedding spaces. This is because the pixel embeddings of CLIP are noisy (Fig. 2), which can be mitigated by aggregating over segments. Hence, we use the pixel embeddings generated by our model to derive segments (clusters) and then apply them to the CLIP's pixel embeddings. In the end, for each segment $s$, the *embedding*
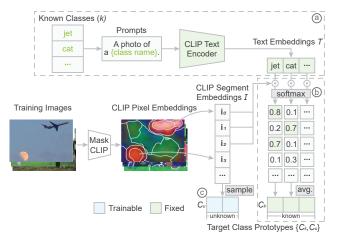


Figure 3. Computation of **target class prototypes** with both *knowns* and *unknowns*, $C = \{C_k, C_u\}$. For a set of known classes (e.g., bird, cat), we first obtain their CLIP text embeddings, $T$, via a set of prompt templates [13,49] (shown in ⓐ); then, we calculate the normalized (via softmax) similarity between text embeddings, $T$, and all segments' CLIP embeddings, $I$, from training images, and average the top-$m$ similar segments' CLIP embedding as the embedding prototype for this class, as shown in ⓑ; For each unknown class, we randomly select the CLIP embedding of a segment as the initial prototype (shown in ⓒ).

*consistent loss* is defined as:

$$\mathcal{L}_e(s) = 1 - sim(\mathbf{v}_s, \mathbf{i}_s), \quad (2)$$

where $\mathbf{v}_s$ and $\mathbf{i}_s$ are the segment embeddings derived from our embedding function and CLIP, respectively. Here, $\mathbf{i}_s$ is the average of the CLIP pixel embedding from segment $s$, namely, $\mathbf{i}_s = \sum_{p\in s} clip\text{-}i(s)/|s|$.

**Semantic Consistency** In addition to embedding consistency, we introduce semantic consistency by forcing our model to make the same predictions of semantic classes as CLIP. The rationale is that we can generate better pixel embeddings if they can form distinctive clusters corresponding to different semantic classes, as the goal of semantic segmentation is to perform pixel-wise classification.

Semantic consistency is achieved via a similar idea of pseudo-labeling [38]. Again, we force the semantic consistency at the segment level (not the pixel level) to reduce the noise in pseudo-labels. Specifically, for each segment $s$, we first use CLIP to generate its pseudo-label $y_s$ over a set of target classes, which include both knowns and unknowns (we will introduce how to design these target classes later). The pseudo-label is generated based on the highest similarity between the segment embedding $\mathbf{i}_s$ with a set of prototypes, $C = \{\mathbf{c}_l\}_0^{L-1}$, of the target classes in the pixel embedding space of CLIP, namely, $y_s = \mathbf{argmax}_{l\in L}(sim(\mathbf{i}_s, \mathbf{c}_l))$.

Then we define the *semantic consistent loss* as the cross entropy between our model's prediction $\varphi(\mathbf{v}_s)$ over the tar-

4

get classes and the pseudo-label $y_s$:

$$\mathcal{L}_s(s) = \mathbf{H}(y_s, \varphi(\mathbf{v}_s)), \tag{3}$$

where $\varphi(\mathbf{v}_s) = \mathbf{softmax}(sim(\mathbf{v}_s, C))$.

**Target Class Prototypes** The design of *target classes* and associated *class prototypes*, $C = \{\mathbf{c}_l\}_0^{L-1}$, is crucial to achieve the semantic consistency. Here, a class prototype, $\mathbf{c}_l$, is an embedding vector that can represent a class in an embedding space. For example, it can be the mean vector of embeddings of all segments of a class "car". Currently, most existing methods [25, 49] assume that the target classes are already predefined, which is not feasible in real-world use cases without any human annotations. Thus, those methods cannot handle unknown classes hidden in the data. To address this issue, we introduce two sets of class prototypes of *known*, $C_k = \{\mathbf{c}_0, \ldots, \mathbf{c}_{k-1}\}$, and *unknown classes*, $C_u = \{\mathbf{c}_k, \ldots, \mathbf{c}_{k+u}\}$, where the known classes are predefined by leveraging CLIP and the unknown classes are learned from image data during training. Thus, we have $C = \{\mathbf{c}_l\}_0^{L-1} = \{\mathbf{c}_0, \ldots, \mathbf{c}_{k-1}, \mathbf{c}_k, \ldots, \mathbf{c}_{k+u}\}, L = k + u$.

For known classes, a natural choice is to use the text embeddings generated by CLIP as their class prototype embeddings [25, 49]. However, even though the text embeddings are trained to align with image/pixel embeddings [34], there is still a huge gap between the text and image/pixel embeddings in the joint space of CLIP (Fig. 1c). Therefore, it is challenging to learn meaningful unknown classes from image features when using text embeddings as class prototypes. Hence, in this work, we use the prototype of CLIP pixel embeddings to represent each known class.

To this end, for a set of known classes (e.g., bird, cat), $K = \{0, \ldots, k-1\}$, we first obtain their CLIP text embeddings, $T = \{\mathbf{t}_k\} = \{clip\text{-}t(k)\}$, via a set of prompt templates following [13, 49], as shown in Fig. 3a. We also get a set of CLIP segment embeddings, $I = \{\mathbf{i}_{\hat{s}}\}$, for all training images by a) feeding training images into the modified image encoder of CLIP to get pixel embeddings; b) clustering the pixel embeddings as segments, $\hat{\mathcal{S}}$; c) averaging the pixel embeddings in each segment, $\hat{s}$. Hence, we have embeddings for each segment: $\mathbf{i}_{\hat{s}} = \sum_{p \in \hat{s}} clip\text{-}i(p)/|\hat{s}|$. Then, we calculate the similarity between text embeddings of known classes, $T$, and all CLIP segment embeddings $I$, and normalize the similarities over all classes by softmax. Finally, we average the top-$m$ similar segments' embedding as the embedding prototype for each class, $C_k = \{\mathbf{c}_k\} = avg_m(top\text{-}m_{\hat{s}}(\mathbf{softmax}_k(sim(I, T))))$.

The prototype embeddings of the unknown classes, $C_u$, are randomly initialized by sampling the CLIP embeddings of all segments, namely $C_u = random(clip\text{-}i(\hat{\mathcal{S}}), u)$, with a size of the unknown class of $u$ (Fig. 3c). During training, the embedding $\mathbf{c}_u$ of each unknown class prototype is updated by minimizing its distance to all segments that are classified as this unknown class (similar to updating the cen-

troids in $k$-means clustering):

$$\mathcal{L}_u = \sum_{s \in S_u} (1 - sim(\mathbf{c}_u, clip\text{-}i(s)))/|S_u|, \tag{4}$$

where $S_u$ are the segments classified as the unknown classes. In this way, our model can also learn the pixel representation of unknown classes.

### 3.3. Training and Inference

In summary, we train the pixel embedding function by combining the pixel-segment contrastive loss, embedding consistent loss, and semantic consistent loss:

$$\mathcal{L} = \mathcal{L}_t + \mathcal{L}_e + \mathcal{L}_s. \tag{5}$$

During training, we also update the embeddings for the unknown classes with $\mathcal{L}_u$.

For inference, we use the trained model to generate pixel embeddings for each input image and use the pixel embeddings for different downstream tasks, including language-driven and unsupervised semantic segmentation. For language-driven semantic segmentation, we first obtain the text embeddings of arbitrary inference classes by feeding the prompt-engineered texts into the text encoder of CLIP. Then we assign each pixel with the class label whose text embedding is the closest to CLIP-S[4] pixel embedding. For unsupervised semantic segmentation, we follow previous work [19, 40] that uses $k$ nearest neighbor search or linear classifier to perform semantic segmentation.

## 4. Experiments

We evaluate our model on three tasks: 1) language-driven semantic segmentation for both known and unknown classes; 2) unsupervised semantic segmentation with $k$-means clustering/linear classification; 3) transfer learning of generated pixel embeddings for instance mask tracking. We also conduct ablation studies to understand the components of our model.

### 4.1. Datasets

**Pascal VOC 2012** [12] contains 20 object classes and a background class. It has 1,464 and 1,449 images for training and validation, respectively. Following common practice [27, 48], we augment the training data with additional annotations [15], resulting in 10,582 training images.

**Pascal Context** [29] extends Pascal VOC 2010 [12] with additional annotations on 4,998 training and 5,105 validation images. Following the prior work [49], we use the most common 59 classes for evaluation.

**COCO-Stuff** [5] labels MS COCO [26] with 171 object/stuff classes. It contains 118,287 and 5,000 images for training and validation, respectively.

5

**DAVIS 2017** [33] contains video sequences for instance mask tracking. Following the prior work [18, 47], we train pixel embeddings on Pascal VOC 2012 and evaluate the validation sequences without fine-tuning.

It is worth mentioning that **no ground truth** labels of any datasets are used during training. Instead, we perform self-supervised learning on pseudo segments generated by contour detectors and owt-ucm [2]. Two contour detectors are used including HED [45] for Pascal VOC 2012 and Pascal Context and PMI [20] for COCO-Stuff.

## 4.2. Implementation Details

For self-supervised contrastive learning, images are augmented with the same set of data augmentations as Sim-CLR [7], including random resizing, cropping, flipping, color jittering, and Gaussian blurring. The concentration constant $\kappa$ is set to 10, and the number of segments is set to 36 for each augmented view.

For vision-language guidance, we use pre-trained CLIP models [34] with modified image encoders following [49]. We use prompt-engineered texts with 85 prompt templates to generate text embeddings following [13, 49]. We use the average embedding of the top 32 segments of high probabilities as the prototype of each known class. We set the number of unknown classes to $u = 64$.

Following the prior work [18, 19, 23], we use PSP-Net [48] with a dilated ResNet-50 [17] backbone as the network architecture. The backbone is pre-trained on the ImageNet [9] dataset. We train our model on Pascal VOC 2012 and Pascal Context for 20k iterations and on COCO-Stuff for 80k iterations. We set the batch size to 8 with additional memory banks that cache the segment embeddings of the previous 2 batches. We set the initial learning rate to 0.001 and decay it with a polynomial learning rate policy. We use the CLIP model trained with ViT-B/16 backbone unless otherwise stated.

## 4.3. Language-Driven Semantic Segmentation

For language-driven semantic segmentation, no human annotations are used for either training or inference. At the inference time, each pixel is assigned an arbitrarily given class label whose CLIP text embedding is the closest to this pixel's CLIP-S$^4$ embedding.

We first compare the performance of our method with the state-of-the-art language-driven semantic segmentation approaches [37, 46, 49] on the Pascal Context and COCO-Stuff datasets. The performance is evaluated with the mean Intersection over Union (mIoU). For MaskCLIP and MaskCLIP+ [49], we obtain the results using the same hyper-parameter setting as our approach with CLIP models of two different backbones of ResNet50 and ViT-B/16. Meanwhile, GroupViT [46], ReCo [37], and ReCo+ [37] use completely different training mechanisms compared

| Method | CLIP Model | Pascal Context | COCO-Stuff |
|---|---|---|---|
| | | mIoU | mIoU |
| GroupViT [46] | - | 22.4 | - |
| ReCo [37] | ResNet50x16 + | 26.6 | - |
| ReCo+ [37]† | ViT-L/14@336px | - | 18.4 |
| MaskCLIP [49] | ResNet50 | 18.6 | 10.6 |
| | ViT-B/16 | 25.2 | 15.2 |
| MaskCLIP+ [49]† | ResNet50 | 23.4 | 13.9 |
| | ViT-B/16 | 32.2 | 20.7 |
| CLIP-S$^4$† | ResNet50 | **28.5** (+5.1) | **16.7** (+2.8) |
| | ViT-B/16 | **33.6** (+1.4) | **22.1** (+1.4) |

Table 2. **Language-guided semantic segmentation benchmarks (mIoU).** CLIP-S$^4$ consistently outperforms the state-of-the-art methods on both Pascal Context and COCO-Stuff datasets with CLIP models of different backbones. † indicates the models are fine-tuned on target datasets.

with our method. GroupViT is trained on image-caption pairs, and ReCo/ReCo+ combines image retrieval and co-segmentation. For comparison, we take the best results from [37, 46] for GroupViT, ReCo, and ReCo+. Tabel 2 shows the benchmarking results of the aforementioned methods. Our method consistently outperforms the state-of-the-art on both datasets with CLIP models of different backbones.

To evaluate models' performance for *class-free semantic segmentation* with both known and unknown classes, we split the 59 classes of Pascal Context into 4 folds, where each fold includes around 15 classes. For each experiment, classes from one fold are considered as unknown and **excluded during training**. The mIoUs of known and unknown classes, as well as their harmonic mean (hIoU) are reported in Tab. 3. The performance of our method is averaged across 5 runs with randomly initialized prototypes of unknown classes. Our method achieves significant gains over MaskCLIP+ on unknown classes, which are comparable to MaskCLIP. Also, our method outperforms both MaskCLIP and MaskCLIP+ on known classes, which leads to better overall performance.

Qualitatively, the visualization in Fig. 4a offers us some insights into why our approach can achieve better results: our model yields *consistent* embeddings aligned with the pre-trained CLIP model. Fig. 4a visualizes the projection of segment embeddings generated by different methods on Pascal Context and COCO-Stuff. We observe that segment embeddings generated by MaskCLIP+ are distorted by the given text embeddings. Meanwhile, CLIP-S$^4$ generates segment embeddings that are well aligned with the segment embeddings derived from the pre-trained CLIP model. Hence, segment embeddings generated by CLIP-S$^4$ can better capture both known and unknown classes. Fig. 4b shows

| Method | fold0 | | | fold1 | | | fold2 | | | fold3 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $mIoU_u$ | $mIoU_k$ | hIoU | $mIoU_u$ | $mIoU_k$ | hIoU | $mIoU_u$ | $mIoU_k$ | hIoU | $mIoU_u$ | $mIoU_k$ | hIoU |
| MaskCLIP | 29.7 | 23.7 | 26.3 | **23.7** | 25.7 | 24.6 | **23.9** | 25.7 | 24.7 | 23.4 | 25.8 | 24.5 |
| MaskCLIP+ | 3.6 | 28.5 | 6.3 | 3.0 | 29.2 | 5.4 | 4.8 | 29.2 | 8.2 | 4.5 | 29.9 | 7.8 |
| CLIP-S$^4$ | **32.0±0.8** | **29.4±0.3** | **30.6±0.5** | 22.3±0.9 | **32.8±0.4** | **26.5±0.6** | 22.4±0.5 | **32.1±0.5** | **26.4±0.4** | **28.6±0.8** | **31.5±0.2** | **30.0±0.5** |
| *vs. MaskCLIP+* | **+28.4** | **+0.9** | **+24.3** | **+19.3** | **+3.6** | **+21.1** | **+17.6** | **+2.9** | **+18.2** | **+24.1** | **+1.6** | **+22.2** |

Table 3. **Language-guided semantic segmentation benchmarks (mIoU) for unknown classes.** The classes of Pascal Context are split into 4 folds with around 15 classes each fold. For each experiment, classes of one fold are considered as unknown. The performance of CLIP-S$^4$ is averaged over 5 runs with randomly initialized unknown class embeddings. CLIP-S$^4$ significantly outperforms MaskCLIP+ on unknown classes. Meanwhile, CLIP-S$^4$ archives consistent gains on known classes over MaskCLIP and MaskCLIP+, and hence leads to better overall performance.



(a) Embeddings projection for PASCAL Context (left) and COCO-Stuff (right)



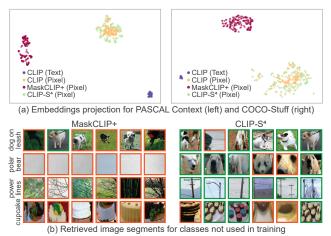(b) Retrieved image segments for classes not used in training

Figure 4. (a) Projection of pixel embeddings generated by CLIP, MaskCLIP+, and CLIP-S$^4$ trained on Pascal Context with CLIP's ViT-B/16 model (left) and COCO-Stuff with CLIP's ResNet50 model (right). MaskCLIP+ distorts pixel embeddings with respect to the given text embeddings, while CLIP-S$^4$ aligns pixel embeddings with the pre-trained CLIP model. (b) Image segments retrieved for classes that are not used during training. Correct and incorrect retrievals are outlined in green and orange, respectively. Compared with MaskCLIP+, CLIP-S$^4$ retrieves images that are more closely related to the classes.

image segments retrieved from the COCO-Stuff validation set using MaskCLIP+ and CLIP-S$^4$ for a set of classes that are not used in training. For each class, we obtain its text embedding and compare it with segment embeddings to obtain top retrievals from both methods. Due to better alignment with the pre-trained CLIP model, CLIP-S$^4$ retrieves images that are more closely related to the unknown classes compared with MaskCLIP+.

## 4.4. Unsupervised Semantic Segmentation

To study whether CLIP-S$^4$ can generate pixel embeddings that form distinctive clusters, we evaluate CLIP-S$^4$ on the unsupervised semantic segmentation task for Pascal VOC 2012. To derive semantic segmentation from pixel embeddings, we use and test two approaches, including $k$ nearest neighbor ($k$-NN) search [19] and linear classification [40]. For $k$-NN search, we assign each segment a class label by the majority vote of its nearest neighbors from the

| Method | mIoU $k$-NN | mIoU Linear Classifier |
|---|---|---|
| IIC [22] | - | 28.0 |
| SegSort [19] | 47.3 | 55.4 |
| Hierarch. Group. [47] | - | 48.8 |
| MaskContrast (Sup.) [40] | 53.9 | 63.9 |
| ConceptContrast [18] | 58.8 | 60.4 |
| HSG [23] | 61.7 | - |
| MaskCLIP [49] | 67.3 | 69.5 |
| MaskCLIP+ [49] | 65.1 | 70.0 |
| CLIP-S$^4$ | **72.0 (+4.7)** | **73.0 (+3.0)** |

Table 4. **Unsupervised semantic segmentation benchmarks (mIoU) on Pascal VOC 2012.** CLIP-S$^4$ consistently outperforms the state-of-the-art methods on both $k$-NN search and linear classification.

training set following [19]. For linear classification, we train a linear classifier on the learned pixel embeddings following [40].

We compare our method with the state-of-the-art unsupervised and language-guided semantic segmentation approaches. We train the state-of-the-art models using the same hyper-parameter setting as our approach except for IIC [22] and Hierarchical Grouping [47] as they use different training mechanisms. For comparison, we take the best results for IIC and Hierarchical Grouping. The benchmark results are shown in Tab. 4. With the vision-language guidance, our method achieves significant gains compared with the previous non-CLIP-based approaches (i.e., +9% for both $k$-NN search and linear classification). Meanwhile, our method also outperforms the language-guided semantic segmentation approaches by a large margin.

## 4.5. Instance Mask Tracking

We evaluate the transferability of pixel embeddings learned from the Pascal VOC 2012 dataset. We use the pixel embeddings to track instance masks in the DAVIS 2017 validation set, where the instance masks at the first frame are given for each video. Following the prior work [48], we use the similarity between pixel embeddings cross frames to propagate the instance masks to the rest of the video frames. We evaluate the performance using the region similarity $\mathcal{J}$

| Method | $\mathcal{J}$(Mean)↑ | $\mathcal{F}$(Mean)↑ |
|---|---|---|
| MaskTrack-B [32] | 35.3 | 36.4 |
| OSVOS-B [4] | 18.5 | 30.0 |
| Video Colorization [42] | 34.6 | 32.7 |
| CycleTime [43] | 41.9 | 39.4 |
| mgPFF [24] | 42.2 | 46.9 |
| Hierarch. Group. [47] | 47.1 | 48.9 |
| MaskContrast (Sup.) [40] | 34.3 | 36.7 |
| ConceptContrast [18] | 50.4 | 53.9 |
| MaskCLIP [49] | 48.1 | 49.2 |
| MaskCLIP+ [49] | 42.6 | 44.2 |
| CLIP-S$^4$ | **52.3 (+1.9)** | **56.8 (+2.9)** |

Table 5. **Quantitative evaluation of instance mask tracking on the DAVIS-2017 validation set.** The performance is measured by the region similarity $\mathcal{J}$ (IoU) and the contour-based accuracy $\mathcal{F}$ defined by [33]. Our method outperforms existing supervised, unsupervised, and language-guided approaches on both metrics.

| $\mathcal{L}_t$ | $\mathcal{L}_e$ | $\mathcal{L}_s$ | pAcc | mIoU | $avgsim$ |
|---|---|---|---|---|---|
| ✓ | - | - | 1.6 | 0.5 | -0.01 |
| ✓ | ✓ | - | 48.1 | 24.3 | 0.79 |
| ✓ | - | ✓ | 52.3 | 32.9 | 0.33 |
| - | ✓ | ✓ | 48.6 | 31.3 | - |
| ✓ | ✓ | ✓ | **53.7** | **33.6** | 0.66 |

Table 6. **Ablation study on the contribution of each loss of CLIP-S$^4$.** Experimented on language-guided semantic segmentation of Pascal Context. $avgsim$ represents the average cosine similarity between segment embeddings generated by CLIP-S$^4$ and CLIP. By combining the embedding and semantic consistent losses $\mathcal{L}_e$ and $\mathcal{L}_s$, CLIP-S$^4$ archives better semantic segmentation performance while maintaining the alignment with CLIP's embeddings.

(IoU) and the contour-based accuracy $\mathcal{F}$ defined by [33].

We compare our method with existing supervised [4, 32], unsupervised [18, 24, 40, 42, 43, 47], and language-guided [49] approaches (Tab. 5). Though not trained on any video sequences, our method outperforms the existing approaches by more than 1.9% and 2.9% in terms of the region similarity $\mathcal{J}$ and contour accuracy $\mathcal{F}$, respectively. Note that the pixel embeddings generated by MaskCLIP+ [49] are distorted by the classes from Pascal VOC 2012, which hinder their transferability.

### 4.6. Ablation Study

We study the contribution of different losses of our method using the Pascal Context dataset and the language-guided semantic segmentation task. The performance is evaluated with pixel accuracy (pAcc) and mIoU. We also calculate the average cosine similarity ($avgsim$) between our segment embeddings and CLIP's segment embeddings to quantify the alignment. Tab. 6 shows the study results. We observe that by introducing embedding consistent loss $\mathcal{L}_e$ the learned segment embeddings are well aligned with

| #unknowns($u$) | mIoU |
|---|---|
| 16 | 71.6 |
| 32 | 71.9 |
| 64 | 72.0 |
| 128 | 72.2 |
| 256 | 71.9 |

| top-$m$ segments | $avgsim$ |
|---|---|
| 1 | 0.921 |
| 4 | 0.977 |
| 16 | 0.996 |
| 64 | 0.997 |
| 256 | 0.993 |

Table 7. **Ablation study** on the influence of the number of unknown class prototypes.

Table 8. **Ablation study** on different numbers of top-$m$ segments for class prototypes.

CLIP's embeddings with an average cosine similarity of 0.79. However, the learned segment embeddings do not perform well on the language-guided semantic segmentation task (24.3 vs. 33.6), because the segment embeddings are not optimized to classify target classes. Meanwhile, by using semantic consistent loss $\mathcal{L}_s$ without embedding consistent loss, the learned segment embeddings have the discriminative power to classify different classes but are not aligned with CLIP's embeddings as the average cosine similarity is 0.33. As a result, the segment embeddings are limited to the target classes used during training. Hence, we combine $\mathcal{L}_e$ and $\mathcal{L}_s$ to balance the discriminative power over target classes and the alignment with CLIP. Meanwhile, we observe that with pixel-segment contrastive learning, the model can archive better performance.

Also, we study the influence of the number of unknown class prototypes on the Pascal VOC dataset for the unsupervised semantic segmentation task. The results Tab. 7 show that the semantic segmentation performance is robust to the tested number of unknown class prototypes as the mIoU varies only 0.6%. Furthermore, we investigate how the size of top-$m$ segments impacts the embeddings of class prototypes. We compare the embeddings of class prototypes generated with different numbers of top-$m$ segments on the Pascal VOC dataset. We use the embeddings of class prototypes generated with $m = 32$ segments as the reference, and compute the cosine similarity between the reference prototypes embeddings and ones with different top-$m$ segments. For each case, the cosine similarity is averaged over all class prototypes. We observe that the embeddings of class prototypes are relatively stable if moderate top-$m$ segments (e.g., $m = 32$ in this work) are used (Tab. 8).

## 5. Conclusion

We propose CLIP-S$^4$, a novel pixel representation learning approach for semantic segmentation. Our method combines self-supervised contrastive learning and guidance of CLIP to learn consistent pixel embeddings with respect to visual and conceptual semantics. Our experiments on popular semantic segmentation benchmarks demonstrate consistent gains over the state-of-the-art unsupervised semantic segmentation and language-driven semantic segmentation methods, especially for unknown classes.

# References

[1] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk. SLIC superpixels compared to state-of-the-art superpixel methods. *IEEE TPAMI*, 34(11):2274–2282, 2012. 3

[2] Pablo Arbeláez, Michael Maire, Charless Fowlkes, and Jitendra Malik. Contour detection and hierarchical image segmentation. *IEEE TPAMI*, 33(5):898–916, 2011. 3, 6

[3] Maxime Bucher, Tuan-Hung VU, Matthieu Cord, and Patrick Pérez. Zero-shot semantic segmentation. In *NeurIPS*, pages 468–479, 2019. 2

[4] S. Caelles, K.-K. Maninis, J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool. One-shot video object segmentation. In *CVPR*, pages 5320–5329, 2017. 8

[5] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. COCO-Stuff: Thing and stuff classes in context. In *CVPR*, pages 1209–1218, 2018. 5

[6] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. *IEEE TPAMI*, 40(4):834–848, 2018. 1

[7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607, 2020. 1, 6

[8] Jang Hyun Cho, Utkarsh Mall, Kavita Bala, and Bharath Hariharan. PiCIE: Unsupervised semantic segmentation using invariance and equivariance in clustering. In *CVPR*, pages 16794–16804, 2021. 2

[9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009. 6

[10] Xiaoyi Dong, Yinglin Zheng, Jianmin Bao, Ting Zhang, Dongdong Chen, Hao Yang, Ming Zeng, Weiming Zhang, Lu Yuan, Dong Chen, Fang Wen, and Nenghai Yu. MaskCLIP: Masked self-distillation advances contrastive language-image pretraining. *arXiv preprint arXiv:2208.12262*, 2022. 3

[11] Sepideh Esmaeilpour, Bing Liu, Eric Robertson, and Lei Shu. Zero-shot out-of-distribution detection based on the pre-trained model CLIP. In *AAAI*, 2022. 3

[12] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 88(2):303–338, 2010. 5

[13] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. In *ICLR*, 2022. 3, 4, 5, 6

[14] Mark Hamilton, Zhoutong Zhang, Bharath Hariharan, Noah Snavely, and William T. Freeman. Unsupervised semantic segmentation by distilling feature correspondences. In *ICLR*, 2022. 1, 2

[15] Bharath Hariharan, Pablo Arbeláez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *ICCV*, pages 991–998, 2011. 5

[16] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, pages 9726–9735, 2020. 1

[17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 6

[18] Wenbin He, William Surmeier, Arvind Kumar Shekar, Liang Gou, and Liu Ren. Self-supervised semantic segmentation grounded in visual concepts. In *IJCAI*, pages 949–955, 2022. 1, 2, 3, 6, 7, 8

[19] Jyh-Jing Hwang, Stella X. Yu, Jianbo Shi, Maxwell D. Collins, Tien-Ju Yang, Xiao Zhang, and Liang-Chieh Chen. SegSort: Segmentation by discriminative sorting of segments. In *ICCV*, pages 7333–7343, 2019. 1, 2, 5, 6, 7

[20] Phillip Isola, Daniel Zoran, Dilip Krishnan, and Edward H. Adelson. Crisp boundary detection using pointwise mutual information. In *ECCV*, pages 799–814, 2014. 6

[21] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting NeRF on a diet: Semantically consistent few-shot view synthesis. In *ICCV*, pages 5865–5874, 2021. 3

[22] Xu Ji, Andrea Vedaldi, and Joao Henriques. Invariant information clustering for unsupervised image classification and segmentation. In *ICCV*, pages 9864–9873, 2019. 2, 7

[23] Tsung-Wei Ke, Jyh-Jing Hwang, Yunhui Guo, Xudong Wang, and Stella X. Yu. Unsupervised hierarchical semantic segmentation with multiview cosegmentation and clustering transformers. In *CVPR*, pages 2561–2571, 2022. 1, 2, 3, 6, 7

[24] Shu Kong and Charless Fowlkes. Multigrid predictive filter flow for unsupervised learning on videos. *arXiv preprint arXiv:1904.01693*, 2019. 8

[25] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. In *ICLR*, 2022. 1, 2, 3, 5

[26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, pages 740–755, 2014. 5

[27] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, pages 3431–3440, 2015. 1, 2, 5

[28] Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. Deep spectral methods: A surprisingly strong baseline for unsupervised semantic segmentation and localization. In *CVPR*, pages 8354–8365, 2022. 2

[29] 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*, pages 891–898, 2014. 5

[30] Yassine Ouali, Céline Hudelot, and Myriam Tami. Autoregressive unsupervised image segmentation. In *ECCV*, pages 142–158, 2020. 2

[31] Dim P. Papadopoulos, Ethan Weber, and Antonio Torralba. Scaling up instance annotation via label propagation. In *ICCV*, pages 15344–15353, 2021. 1

[32] Federico Perazzi, Anna Khoreva, Rodrigo Benenson, Bernt Schiele, and Alexander Sorkine-Hornung. Learning video object segmentation from static images. In *CVPR*, pages 3491–3500, 2017. 8

[33] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *CVPR*, pages 724–732, 2016. 6, 8

[34] 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. 1, 2, 4, 5, 6

[35] 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*, pages 18061–18070, 2022. 3

[36] Ramprasaath R. Selvaraju, Karan Desai, Justin Johnson, and Nikhil Naik. CASTing your model: Learning to localize improves self-supervised representations. In *CVPR*, pages 11053–11062, 2021. 2

[37] Gyungin Shin, Weidi Wie, and Samuel Albanie. ReCo: Retrieve and co-segment for zero-shot transfer. In *NeurIPS*, 2022. 1, 2, 3, 6

[38] Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Han Zhang, and Colin Raffel. FixMatch: Simplifying semi-supervised learning with consistency and confidence. In *NeurIPS*, pages 596–608, 2020. 4

[39] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, and Luc V Gool. Revisiting contrastive methods for unsupervised learning of visual representations. In *NeurIPS*, pages 16238–16250, 2021. 2

[40] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, and Luc Van Gool. Unsupervised semantic segmentation by contrasting object mask proposals. In *ICCV*, pages 10032–10042, 2021. 1, 2, 5, 7, 8

[41] Yael Vinker, Ehsan Pajouheshgar, Jessica Y. Bo, Roman Christian Bachmann, Amit Haim Bermano, Daniel Cohen-Or, Amir Zamir, and Ariel Shamir. CLIPasso: Semantically-aware object sketching. *ACM TOG*, 41(4):86:1–86:11, 2022. 3

[42] Carl Vondrick, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy. Tracking emerges by colorizing videos. In *ECCV*, pages 402–419, 2018. 8

[43] Xiaolong Wang, Allan Jabri, and Alexei A. Efros. Learning correspondence from the cycle-consistency of time. In *CVPR*, pages 2561–2571, 2019. 8

[44] Zhaoqing Wang, Yu Lu, Qiang Li, Xunqiang Tao, Yandong Guo, Mingming Gong, and Tongliang Liu. CRIS: Clip-driven referring image segmentation. In *CVPR*, pages 11676–11685, 2022. 3

[45] Saining Xie and Zhuowen Tu. Holistically-nested edge detection. In *ICCV*, pages 1395–1403, 2015. 6

[46] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. GroupViT: Semantic segmentation emerges from text supervision. In *CVPR*, pages 18134–18144, 2022. 2, 3, 6

[47] Xiao Zhang and Michael Maire. Self-supervised visual representation learning from hierarchical grouping. In *NeurIPS*, pages 16579–16590, 2020. 2, 6, 7, 8

[48] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *CVPR*, pages 6230–6239, 2017. 1, 5, 6, 7

[49] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from CLIP. In *ECCV*, 2022. 1, 2, 3, 4, 5, 6, 7, 8