

Unsupervised Semantic Segmentation by Contrasting Object Mask Proposals

Wouter Van Gansbeke^{1*} Simon Vandenhende^{1*} Stamatis Georgoulis² Luc Van Gool^{1,2}
¹KU Leuven/ESAT-PSI ²ETH Zurich/CVL, TRACE

Abstract

Being able to learn dense semantic representations of images without supervision is an important problem in computer vision. However, despite its significance, this problem remains rather unexplored, with a few exceptions that considered unsupervised semantic segmentation on small-scale datasets with a narrow visual domain. In this paper, we make a first attempt to tackle the problem on datasets that have been traditionally utilized for the supervised case. To achieve this, we introduce a two-step framework that adopts a predetermined mid-level prior in a contrastive optimization objective to learn pixel embeddings. This marks a large deviation from existing works that relied on proxy tasks or end-to-end clustering. Additionally, we argue about the importance of having a prior that contains information about objects, or their parts, and discuss several possibilities to obtain such a prior in an unsupervised manner.

Experimental evaluation shows that our method comes with key advantages over existing works. First, the learned pixel embeddings can be directly clustered in semantic groups using K-Means on PASCAL. Under the fully unsupervised setting, there is no precedent in solving the semantic segmentation task on such a challenging benchmark. Second, our representations can improve over strong baselines when transferred to new datasets, e.g. COCO and DAVIS. The code is available ¹.

1. Introduction

The problem of assigning dense semantic labels to images, formally known as *semantic segmentation*, is of great importance in computer vision as it finds many applications, including autonomous driving, augmented reality, human-computer interaction, etc. To achieve state-of-the-art performance in this task, fully convolutional networks [47] are typically trained on datasets [16, 22, 46] that contain a large number of fully-annotated images. However, obtaining accurate, pixel-wise semantic labels for every image in a dataset is a labor-intensive process that costs significant

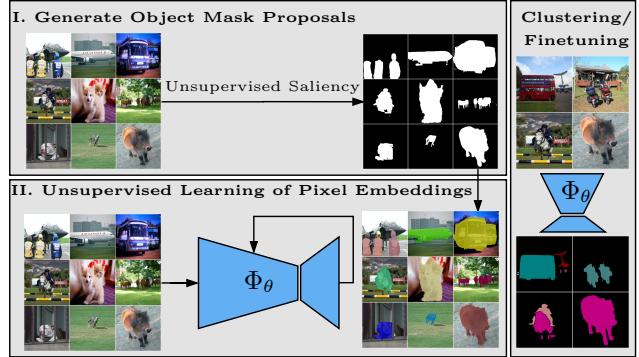


Figure 1. We learn pixel embeddings for semantic segmentation in an unsupervised way. First, we predict object mask proposals using unsupervised saliency. Second, we use the obtained masks as a prior in a self-supervised optimization objective. Finally, the pixel embeddings can be clustered or fine-tuned to a semantic segmentation of the image.

amounts of time and money [4]. To remedy this situation, weakly-supervised methods leveraged weaker forms of supervision, such as scribbles [45, 67, 68, 74, 82], bounding boxes [17, 39, 58, 82], clicks [5], and image-level tags [58, 68, 82], while semi-supervised methods [17, 28, 30, 58, 59] used only a fraction of the dataset as labeled examples, both of which require substantially less human annotation effort. Despite the continued progress, the vast majority of semantic segmentation works still rely on some form of annotations to train the neural network models.

In this paper, we look at the problem from a different perspective, namely self-supervised representation learning. More concretely, we aim to learn pixel-level representations or embeddings for semantic segmentation without using ground-truth. If we obtain a good pixel embedding [18] that is discriminative w.r.t. the semantic classes, we can directly cluster the pixels into semantic groups using K-Means. This tackles the semantic segmentation problem under the fully unsupervised setup. Alternatively, if a limited number of annotated examples are available, the representations can be further fine-tuned under a semi-supervised or transfer learning setup. In this paper, we primarily focus on the fully unsupervised setup, but include additional fine-

*Authors contributed equally

¹github.com/wvangansbeke/Unsupervised-Semantic-Segmentation.git

tuning experiments for the sake of completeness.

Unsupervised or self-supervised techniques [38] were recently being employed to learn rich and effective visual representations without external supervision. The obtained representations can subsequently be used for a variety of purposes, including task transfer learning [26], image clustering [2, 3, 73], semi-supervised classification [13], etc. Popular representation learning techniques used an instance discrimination task [80], that is treating every image as a separate class, to generate representations in an unsupervised way. Images and their augmentations are considered as positive examples of the class, while all other images are treated as negatives. In practical terms, the instance discrimination task is formulated as a non-parametric classification problem, and a contrastive loss [25, 56] is used to model the distribution of negative instance classes.

Purushwalkam and Gupta [63] showed that contrastive self-supervised methods learn to encode semantic information, since two views of the same image will always show a part of the same object, and no objects from other categories. However, under this setup, there is no guarantee that the representations also learn to differentiate between pixels belonging to different semantic classes. For example, when foreground-background pairs frequently co-occur, e.g. cattle grazing on farmland, pixels belonging to the two classes can share their representation. This renders existing works based on instance discrimination less appropriate w.r.t. our goal of learning semantic pixel embeddings. To address these limitations, we propose to learn pixel-level, rather than image-level representations, in a self-supervised way.

The proposed method consists of two steps. First, we leverage an unsupervised saliency estimator to mine object mask proposals from the dataset. This mid-level visual prior transfers well across different datasets. In the second step, we use a contrastive framework to learn pixel embeddings. The object mask proposals are employed as a prior - we pull embeddings from pixels belonging to the same object together, and contrast them against pixels from other objects. The generated representations are evaluated on the semantic segmentation task following standard protocols. The framework is illustrated in Figure 1.

Our contributions are: (1) We propose a two-step framework for unsupervised semantic segmentation, which marks a large deviation from recent works that relied on proxy tasks or end-to-end clustering. Additionally, we argue about the importance of having a mid-level visual prior which incorporates object-level information. This contrasts with earlier works that grouped pixels together based upon low-level vision tasks like boundary detection. (2) The proposed method is the first able to tackle the semantic segmentation task on a challenging dataset like PASCAL under the fully unsupervised setting. (3) Finally, we report promising results when transferring our representations to other datasets.

This shows that adopting a mid-level visual prior can be useful for self-supervised representation learning.

2. Related Work

As our method is mostly related to unsupervised semantic segmentation and representation learning, in what follows we discuss representative works from each topic.

Unsupervised semantic segmentation. There have only been a few attempts in the literature to tackle semantic image segmentation under the fully unsupervised setting. Some works [36, 57] followed an end-to-end approach - maximizing the discrete mutual information between augmented views to learn a clustering function. However, these methods were only applied to small-scale datasets, covering a narrow visual domain, e.g. separating sky from vegetation, using satellite imagery, etc. In contrast, our method applies to more challenging scenarios, and decouples feature learning from clustering.

A few works [31, 97] used segments obtained from boundaries to learn pixel embeddings in a self-supervised way. However, it is unclear whether the representations could be post-processed with an off-line clustering criterion to obtain discrete labels. In particular, the evaluation only considered semantic segment retrieval which requires an annotated train set. Furthermore, Hwang *et al.* [31] still relied on additional supervision sources like ImageNet pre-training and boundary annotations [1, 81].

Representation learning. These methods aim at learning visual representations by solving pre-designed *pretext tasks*, which do not require manual annotations. Examples of such pretext tasks include colorizing images [32, 42, 96], predicting context [19, 51], solving jigsaw puzzles [53, 55], generating images [65], clustering [2, 8, 84], predicting noise [6], spotting artifacts [35], using adversarial training [20, 21], predicting optical flow [49, 90], counting [54], inpainting [60], predicting transformation parameters [23, 94], using predictive coding [56], performing instance discrimination [9, 10, 12, 24, 26, 43, 50, 70, 71, 80, 87], and so on. The learned representations can subsequently be transferred to learn a separate down-stream task, e.g. object detection.

In a similar vein, some works tried to learn pixel-level representations for semantic segmentation by solving proxy tasks, e.g. colorization [32, 42, 89, 96], optical flow [49, 90], using co-occurrences [33], etc. Differently, in this paper, we avoid the use of a proxy task.

3. Method

In this paper, we aim to learn a pixel embedding function for semantic segmentation from an unlabeled dataset of images. Since the goal of semantic segmentation is to assign a

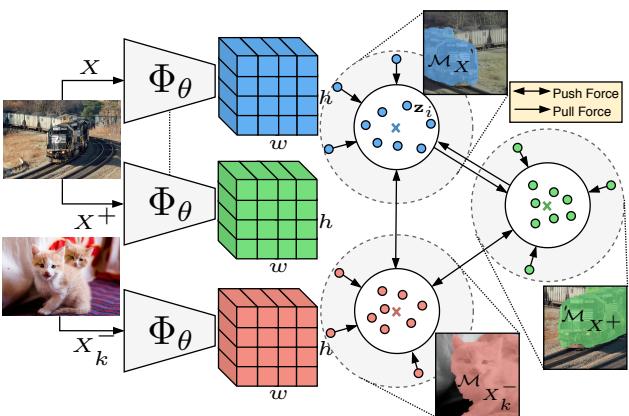


Figure 2. **MaskContrast** learns pixel embeddings for unsupervised semantic segmentation in the following way. We use a saliency estimator to generate positive pairs of object-centric crops (X, X^+) and negative pairs X_k^- . The model Φ_θ is trained to maximize the agreement between embeddings of pixels belonging to the objects in X, X^+ , while minimizing the agreement with pixels from objects in X_k^- .

class label to every pixel of an image, a good pixel embedding should be discriminative w.r.t. the semantic classes. If the latter holds true, the embedding function can be directly used to cluster the pixels into semantic groups, or be further fine-tuned under a semi-supervised setup.

To tackle the aforementioned problem, we follow a divide-and-conquer strategy. We argue that it is more difficult to directly cluster the pixels into semantic groups following an end-to-end pipeline, while it is easier to first look for image regions where pixels are likely to belong together. Although this information does not directly result in a semantic segmentation of the scene, it gives us a useful starting point to learn the pixel embeddings. In particular, we can leverage the obtained regions as a prior by grouping their pixels together. Since the prior is determined before the feature learning step, we reduce the dependence on the network initialization. This is an intentional divergence from existing end-to-end learning pipelines [36, 57], which are prone to latch onto low-level image cues - like color, contrast, etc. - as shown in [73].

The proposed method named *MaskContrast* consists of two steps. In a first step, we determine a prior by identifying objects in the images for which pixels can be grouped together. Mid-level visual groups, like objects, transfer well across datasets, since they do not depend on any pre-defined ground-truth classes. In the second step, we employ the obtained prior in a contrastive loss [25, 56] to generate pixel embeddings. More specifically, we pull pixels belonging to the same object together, and contrast them against pixels from other objects, as shown in Figure 2. This forces the model to map pixels from visually similar objects closer to-

gether, while pushing pixels from dissimilar objects further apart. In this way, the model discovers a pixel embedding space that can serve as a dense semantic representation of the scene.

The method section is further organized as follows. Section 3.1 motivates the use of object mask proposals as a prior for semantic segmentation. Section 3.2 analyzes the use of an unsupervised saliency estimator to mine the object masks from unlabeled datasets. Section 3.3 integrates the prior in a contrastive loss to learn pixel embeddings.

3.1. A Mid-Level Visual Prior for Grouping Pixels

As a starting point for unsupervised semantic segmentation, we try to define an appropriate prior. Several works have emerged in the literature that tried to group pixels by solving a proxy task. Examples include colorizing images [32, 42, 96], predicting optical flow [49, 90], using co-occurrences [33], etc. Unfortunately, there is no guarantee that the generated representations will align with the semantic classes, as the latter are co-variant to the proxy task’s output. For example, a colorization network will be sensitive to color changes, even though these do not necessarily alter the semantics of the scene. This behavior is unwanted for the objective of semantic segmentation.

To overcome these limitations, we follow an alternative route that avoids the use of a proxy task. In particular, we mine object mask proposals which cover patches that are likely to contain an object. A prior can then be defined from the masks based upon *shared pixel ownership*, i.e. if a pair of pixels belongs to the same mask, we assume that they should be grouped together, and maximize the agreement between their pixel embeddings. We hypothesize that this is a more reliable pixel grouping strategy compared to the use of proxy tasks. In particular, our approach builds a high-level image segmentation by first identifying mid-level visual groups, instead of directly producing a complete segmentation by solving a proxy task. A motivation for this bottom-up approach is also provided in [66].

At the same time, the proposed prior can be seen as an object-centric approach to unsupervised semantic segmentation, which brings several advantages to the table. First, using mid-level visual cues, like object information, regularizes the feature representations. In particular, the model can not simply rely on low-level information like color to group the pixels together, but needs to learn more semantically meaningful image characteristics. This differs from competing works [31, 97] that used superpixels or image boundaries as a prior. Second, object cues can be highly informative of the semantic segmentation task. Evidence for the latter has been provided in the literature for weakly-supervised methods that utilized annotations containing object information. As an example, several works [17, 39, 58, 82] reported strong results on the seg-

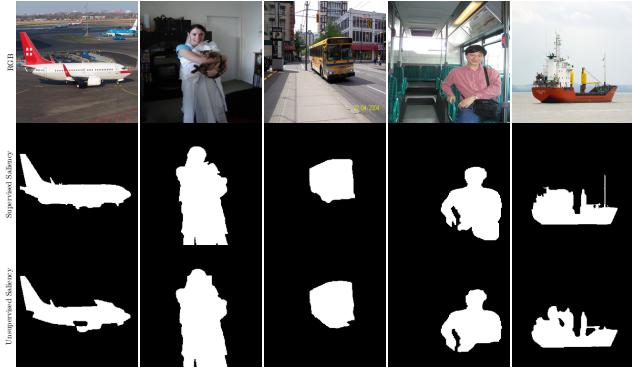


Figure 3. **Mask Proposals.** We train a supervised (middle) and unsupervised (bottom) saliency estimator on the DUTS and MSRA datasets respectively. We make predictions on PASCAL.

mentation task by employing object bounding boxes.

Next, we show how an unsupervised saliency estimator can be used to generate the object mask proposals.

3.2. Mining Object Mask Proposals

We need to retrieve a set of object mask proposals for the images in our dataset. The literature [1, 52, 62, 72] offers a multitude of ways to do this. We prefer to use a simple strategy to verify whether unsupervised semantic segmentation benefits from adopting a mid-level visual prior. Moreover, we would like to use a method that does not rely on external supervision, or can be trained with a limited amount of annotations. In the latter case, the object mask proposal mechanism should generalize well to new scenes.

Based upon our requirements, we propose the use of saliency estimation [7, 79] to generate object masks proposals. Most importantly, various unsupervised methods can be used for this purpose. Several of these works [52, 91, 93] used predictions obtained with hand-crafted priors [37, 44, 98, 100] as pseudo-labels to train a deep neural network. Others [85, 86] relied on videos to learn a salient object detector. Furthermore, on a variety of datasets [15, 75, 83] unsupervised saliency methods have shown to perform on par with their supervised counterparts [29, 48, 64, 77, 92, 95]. Finally, the model predictions transfer well to novel unseen datasets as shown by [52].

For completeness, in Section 4 we explore both unsupervised [52] and supervised [64] saliency estimation methods to predict the object masks, and showcase the potential of our method. Figure 3 shows some examples.

3.3. MaskContrast: Learning Pixel Embeddings by Contrasting Salient Objects

Consider a dataset of images \mathcal{X} with associated non-overlapping object mask proposals $\{\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_N\}$ obtained using a saliency estimator. Our goal is to learn a *pixel embedding function* $\Phi_\theta : \mathcal{X} \rightarrow \mathcal{Z}$ parameterized

by a neural network with weights θ , that maps each pixel i in an image to a point \mathbf{z}_i on a D -dimensional normalized hyper-sphere. We chose a normalized embedding space, so that the output of Φ_θ is bounded. Note that, the use of such scale-invariant embeddings decouples the loss from other design choices that could implicitly limit the range of distances, e.g. weight decay, as shown in [41].

We construct an optimization objective to learn the embedding function Φ_θ as follows. First, we describe how to learn semantically meaningful image feature using a contrastive learning objective. Second, we modify the criterion to learn pixel embeddings.

Learning Image-Level Representations. Existing contrastive self-supervised methods (e.g. [12, 26, 80]) learn visual representations through an instance discrimination task defined at the image-level. Positive views (X, X^+) of the same image are acquired for which it is guaranteed that both images contain a part of the same object. Similarly, examples of negative pairs $\{(X, X_0^-), (X, X_1^-), \dots, (X, X_K^-)\}$ can be found that never contain the same object. In practice, we impose additional invariances by applying augmentations. The positives and negatives can now be used in a contrastive framework to learn image representations that encode semantic information about the objects.

We realize this concept by training an *image embedding function* Ψ_η to maximize the agreement between positive pairs (X, X^+) , while minimizing the agreement between negative pairs $\{(X, X_0^-), (X, X_1^-), \dots, (X, X_K^-)\}$. If we measure the similarity between pairs using a dot product, the contrastive loss [25, 56] is defined as

$$\mathcal{L} = -\log \frac{\exp(\Psi_\eta(X)^T \cdot \Psi_\eta(X^+)/\tau)}{\sum_{k=0}^K \exp(\Psi_\eta(X)^T \cdot \Psi_\eta(X_k^-)/\tau)}, \quad (1)$$

where the temperature τ relaxes the dot product. As shown by [63], the model learns to encode object information because the positive examples always preserve a part of the same object. Moreover, since the representational capacity of the network is intentionally limited, visually similar objects will tend to be mapped closer together by Ψ_η . The combination of these two properties results in image representations that can be directly clustered into semantic groups (see also [73] for a more detailed explanation).

The above observations showed how to train a model that encodes semantic object information. Next, we modify the contrastive loss from Equation 1 to learn representations at the pixel level.

Learning Pixel-Level Representations. We adopt the following notation. Let i be a pixel with \mathbf{z}_i its pixel embedding. Let $m(i)$ be the index of the object mask that pixel i belongs to, i.e. $i \in \mathcal{M}_{m(i)}$. Finally, let the mean pixel

embedding $\mathbf{z}_{\mathcal{M}_n}$ of an object mask \mathcal{M}_n be defined as

$$\mathbf{z}_{\mathcal{M}_n} = \frac{1}{|\mathcal{M}_n|} \sum_{i \in \mathcal{M}_n} \mathbf{z}_i. \quad (2)$$

The optimization objective is derived from a pull- and push-force in the pixel embedding space.

Pull-force. In Section 3.1, we motivated the use of a prior based upon *shared pixel ownership* to pull pixels together in the embedding space. More concretely, if two pixels i, j belong to the same object, i.e. $m(i) = m(j)$, we maximize the agreement between their pixel embeddings $\mathbf{z}_i, \mathbf{z}_j$. In practice, the agreement is maximized between pixels and the mean embedding of their object mask in order to obtain a criterion that scales linearly with the number of pixels, rather than quadratically.

Push-force. Additionally, we require a push-force to avoid mode collapse in the embedding space. Moreover, the push-force should drive pixels from visually similar objects to lie close together in the embedding space, while pixels from dissimilar objects to be mapped further apart. As motivated in the previous paragraph, this can be achieved by adopting a contrastive loss that takes augmented views of objects as positive pairs, and views of other objects as negatives. In this case, the push-force is found between different objects. We represent the objects by their mean pixel embedding.

Optimization objective. We modify the contrastive loss from Equation 1 to include the proposed pull- and push-forces. Positive pairs of object-centric crops ($\Psi_\eta(X), \Psi_\eta(X^+)$) are replaced with positive pairs of pixel embeddings: $(\mathbf{z}_i, \mathbf{z}_{\mathcal{M}_{X^+}})$ for $i \in \mathcal{M}_X$. In a similar way, the negative pairs $(\Psi_\eta(X), \Psi_\eta(X^-_k))$ are replaced with $(\mathbf{z}_i, \mathbf{z}_{\mathcal{M}_{X^-_k}})$. We obtain the following optimization criterion for a pixel $i \in \mathcal{M}_X$

$$\mathcal{L}_i = -\log \frac{\exp(\mathbf{z}_i \cdot \mathbf{z}_{\mathcal{M}_{X^+}} / \tau)}{\sum_{k=0}^K \exp(\mathbf{z}_i \cdot \mathbf{z}_{\mathcal{M}_{X^-_k}} / \tau)}. \quad (3)$$

The pixel embedding function Φ_θ maximizes the agreement between pixels and an augmented view of the object they belong to, while minimizing the agreement with other objects. We apply the pixel-wise loss \mathcal{L}_i to all foreground pixels. The background pixels are not contrasted, since there could be multiple background objects on which we have no conclusive information. In this case, however, the network does not need to discriminate between pixels that fall inside or outside the object masks. As a consequence, the pixel embeddings can collapse to a single vector across an image. To prevent this, we regularize the feature space by including a separate linear head that predicts the saliency masks. We refer to the pseudo-code in the supplementary materials for an overview of MaskContrast (Algorithm 1).

Interestingly, the proposed objective can also be viewed in an alternative way. Wang and Isola [78] showed that a contrastive loss optimizes two properties: (1) alignment of features from positive pairs and (2) uniformity of the feature distribution on a normalized hyper-sphere. From this viewpoint, our optimization objective can also be interpreted as optimizing the alignment of pixel embeddings based upon shared pixel ownership, while spreading pixel embeddings uniformly across the hyper-sphere \mathcal{Z} .

4. Experiments

4.1. Experimental Setup

Datasets. We conduct the bulk of our experimental analysis on the PASCAL [22] dataset following prior work [31, 97]. The `train_aug` and `val` splits are used during training and evaluation respectively. We perform additional experiments on the COCO [46] and DAVIS-2016 [61] datasets to verify if the pixel embeddings transfer to novel scenes. We use the annotations from Kirillov *et al.* [40] for the semantic segmentation task on COCO and evaluate on the PASCAL classes. On DAVIS-2016, the representations are used to compute correspondences for propagating object masks in videos. Only the first frame is annotated and we evaluate the propagated masks on the rest of the frames. We adopt the evaluation protocol from [34], and report the region similarity \mathcal{J} and contour-based accuracy \mathcal{F} scores.

Training setup. We use a DeepLab-v3 [11] model with dilated [88] ResNet-50 backbone [27]. The backbone is initialized from MoCo v2 [14] pre-trained on ImageNet, unless defined otherwise. We train the model for 60 epochs using batches of size 64. The model weights are updated through SGD with momentum 0.9 and weight decay $1e^{-4}$. The initial learning rate is set to 0.004 and decayed with a poly learning rate scheme. We use the same set of augmentations as SimCLR [12] to generate positive pairs (X, X^+) , while making sure that each image contains at least a part of the salient object (area $> 10\%$). The features of negatives $\{\mathbf{z}_{\mathcal{M}_{X^-_0}}, \dots, \mathbf{z}_{\mathcal{M}_{X^-_K}}\}$ are saved in a memory bank, with K set to 128. The negatives are encoded with a momentum-updated version of the network following [26]. We use dimension $D = 32$ and temperature $\tau = 0.5$.

Saliency estimation. We test both unsupervised and supervised saliency estimators to mine the object mask proposals. We adopt the BAS-Net [64] architecture. The *supervised saliency model* is trained on DUTS [76]. Differently, the *unsupervised saliency model* is trained on MSRA [15] using the approach from DeepUSPS [52]. MSRA considers less complex scenes from which the unsupervised training benefits. However, directly transferring the predictions

Method	LC (MIoU)
Supervised Saliency Model	6.5
MoCo v2 [14] (Unsupervised)	45.0
ImageNet (IN) Classifier (Supervised)	53.1
MaskContrast (MoCo v2 Init. - Unsup. Sal. Model)	58.4
MaskContrast (MoCo v2 Init. - Sup. Sal. Model)	62.2
MaskContrast (IN Classifier Init. - Unsup. Sal. Model)	61.0
MaskContrast (IN Classifier Init. - Sup. Sal. Model)	63.9

Table 1. **Baseline comparison** under the linear evaluation protocol on PASCAL.

to our target datasets, e.g. PASCAL, results in low-quality mask proposals when using the unsupervised model. We employ a simple bootstrapping procedure to improve the predictions on the target datasets. In particular, we obtain our final saliency estimator from training BAS-Net on pseudo-labels generated with the unsupervised DeepUSPS model on MSRA.

Implementation. We provide the implementation details of every method in the supplementary materials. The code and pre-computed saliency masks will be made available.

Scope. We adopt standard evaluation protocols [36, 97] for unsupervised semantic segmentation to benchmark our method. More specifically, we use linear probes (Sec. 4.3), direct clustering (Sec. 4.4) and a segment retrieval approach (Sec. 4.5) to quantify if the pixel embeddings are disentangled according to the semantic classes. This experimental setup differs from the typical setting used in self-supervised representation learning, where the evaluation focuses on fine-tuning the feature representations to various down-stream tasks. For completeness, we include additional fine-tuning experiments in Sections 4.6 - 4.7.

4.2. Ablation Studies

We examine the influence of the different components of our framework under the linear evaluation protocol following existing work [97]. The network weights are kept fixed and we train a 1×1 convolutional layer on top to predict the class assignments. Since the discriminative power of a linear classifier is low, the pixel embeddings need to be informative of the semantic class to solve the task in this way.

Baseline comparison. Table 1 compares several baselines. Applying a linear classifier on top of the saliency features results in the lowest performance (6.5%). This is to be expected since the saliency estimator only discriminates between two groups of pixels, i.e. the salient object vs. background. Differently, our method discovers a semantically structured embedding space, where pixels from visually similar objects lie close together, while pixels from dissimilar objects end up far apart. This allows a linear classifier to correctly group the pixels ($> 58.4\%$). Importantly,

the results improve over the models from which the backbone weights were initialized (45.0% to 58.4% for MoCo and 53.1% to 61.0% for supervised pre-training). We conclude that the performance of our method can not be attributed to the use of a specific initialization. Also, it is beneficial to learn representations at pixel-, rather than at image-level, for the segmentation task. Finally, we observe further performance gains when including additional supervision, e.g. supervised pre-training on ImageNet (58.4% to 61.0%), or a supervised saliency estimator (58.4% to 62.2% and 61.0% to 63.9%).

Mask proposals. Table 2a compares three mask proposal strategies. Better numbers are reported when using salient object masks. We found that the regions extracted with the hierarchical segmentation algorithm were often too small to be representative of an object or part. In this way, the model does not learn useful information for the segmentation task. This confirms the hypothesis from Section 3.1, i.e. a good prior expresses object information.

Training mechanisms. Table 2b ablates some of the included training mechanisms. First, using augmented views to sample positive pairs improves the results, as we learn additional invariances. Second, including a memory bank results in further performance gains, because we can better estimate the distribution of negatives. Third, it is helpful to encode the negatives with a momentum-updated version of the network Φ_θ , as this enforces consistency in the memory bank (see also [26]). In summary, all three mechanisms positively contribute to the results.

Hyperparameter study. Table 2c studies the influence of the used temperature τ and number of negatives K . We conclude that the proposed algorithm is not very hyperparameter sensitive based upon the reported standard deviations.

4.3. Linear Classifier

Table 3a compares our method against competing works under the linear evaluation protocol on PASCAL.

MaskContrast vs. proxy tasks. The method substantially outperforms works based on proxy tasks. It is unlikely that a proxy task aligns the embeddings with the semantic groups in the dataset. In contrast, combining our proposed prior, i.e. shared pixel ownership, with a contrastive loss results in more semantically meaningful pixel embeddings.

MaskContrast vs. clustering. We outperform IIC [36] which used a clustering objective. As discussed earlier, the clusters strongly depend on the network initialization, which negatively impacts the learned features as the network can latch onto low-level information, like color, texture, contrast, etc. Differently, we suppress these problems by decoupling the prior from the network initialization.

Mask Proposals	LC (MIoU)	Augmented Views	Memory	Momentum Encoder	LC (MIoU)	Hyperparameter	Range	LC (MIoU)					
Hierarchical Seg. [1, 81]	30.5	✗	✗	✗	52.4	Temperature τ	[0.1-1]	56.2 \pm 1.4					
Unsupervised Sal. Model	58.4	✓	✗	✗	54.0	Negatives K	[64-1024]	57.0 \pm 0.6					
Supervised Sal. Model	62.2	✓	✓	✗	55.0	(c) Hyperparameter study. We report the mean and standard deviation.							
(a) Comparison of three mask proposal mechanisms.								(b) Analysis of the used training mechanisms.					

Table 2. **Ablation studies** of our method under the linear evaluation protocol on PASCAL. Tables 2b- 2c report results with masks from the unsupervised saliency estimator. We use MoCo v2 initial weights.

MaskContrast vs. contrastive learning. The method reports higher accuracy compared to existing contrastive self-supervised approaches. This group of works defined the contrastive loss at the global image- or patch-level. Naturally, our pixel embeddings are more predictive of the semantic segmentation task as we defined a contrastive learning objective at the pixel-level.

MaskContrast vs. boundary based. Finally, we outperform methods that relied on boundary detectors to group pixels together. We argue that the employed saliency masks incorporate higher level visual information compared to the regions obtained from boundary detectors. The results are indicative of the benefits that can be obtained from using such information, which supports our earlier claims from Section 3.1.

4.4. Clustering

We verify whether the feature representations can be directly clustered in semantically meaningful groups using an off-line clustering criterion like K-Means. The number of clusters equals the number of ground-truth classes. The Hungarian matching algorithm is used to match the predicted clusters with the ground-truth classes and the results are averaged across five runs. Table 3b shows the results. Our learned pixel embeddings can be successfully clustered using K-Means on PASCAL. In contrast, the features representations obtained in prior works do not exhibit this behavior. We include additional results in the suppl. materials when applying overclustering.

4.5. Semantic Segment Retrieval

Next, we adopt a retrieval approach to examine our representations on PASCAL. First, we compute a feature vector for every salient object by averaging the pixel embeddings within the predicted mask. Next, we retrieve the nearest neighbors of the val set objects from the train_aug set. Table 4 shows a quantitative comparison with the state-of-the-art for the following 7 classes: bus, airplane, car, person, cat, cow and bottle. As before, we outperform prior works by significant margins. To facilitate future comparison, we also include results when evaluating on all 21 PASCAL classes. Figure 4 shows some qualitative results.

Method	LC	K-Means
<i>Proxy task based:</i>		
Co-Occurrence [33]	13.5	-
CMP [90]	16.5	-
Colorization [96]	25.5	-
<i>Clustering based:</i>		
IIC [36]	28.0	9.8
<i>Contrastive learning based:</i>		
Inst. Discr. [80]	26.8	-
MoCo v2 [26]	45.0	-
InfoMin [71]	45.2	-
SWAV [9]	50.7	-
<i>Boundary based:</i>		
SegSort [31]	36.2	-
Hierarch. Group. [97]	48.8	-
ImageNet (IN) Classifier (Supervised)	53.1	-
MaskContrast (MoCo Init. + Unsup. Sal.)	58.4	35.0
MaskContrast (MoCo Init. + Sup. Sal.)	62.2	38.9
MaskContrast (IN Sup. Init. + Unsup. Sal.)	61.0	41.6
MaskContrast (IN Sup. Init. + Sup. Sal.)	63.9	44.2

(a) Linear classifier.

(b) K-Means.

Table 3. **State-of-the-art comparison** on the PASCAL val set (MIoU). We indicate results on par with random guessing with '-'.

Method	MIoU (7 classes)	MIoU (21 classes)
SegSort [31]	10.2	-
Hierarch. Group. [97]	24.6	-
MoCo v2 [14]	48.0	39.0
MaskContrast (Unsup. Sal.)	53.4	43.3
MaskContrast (Sup. Sal.)	62.3	49.6

Table 4. **State-of-the-art comparison** for semantic segment retrieval on the PASCAL val set. We use MoCo v2 initial weights.

4.6. Transfer Learning

We study the transferability of our pixel embeddings. Table 5 shows the results when pretraining on ImageNet and evaluating the generated pixel embeddings on a different target dataset. Interestingly, our representations transfer well across various datasets. Training a linear classifier to solve the segmentation task on PASCAL improves over the MoCo v2 baseline (55.4% for MaskContrast vs. 45.0% for MoCo when using an unsupervised saliency model). A similar effect can be observed on COCO (45.0% for MaskCon-



Figure 4. **Nearest neighbors** for queries (1st col.) on PASCAL.

Model	PASCAL (MIoU) \uparrow	COCO (MIoU) \uparrow	DAVIS '16 $\mathcal{J}_m \uparrow$	$\mathcal{F}_m \uparrow$
MoCo v2	45.0	35.2	77.1	77.2
MaskContrast (Unsup. Sal.)	55.4	45.0	78.0	77.8
MaskContrast (Sup. Sal.)	57.2	47.2	82.0	80.9

Table 5. **Transfer learning setup.** All models were pre-trained on ImageNet. We use MoCo v2 initial weights. Results on PASCAL and COCO are reported for a linear classifier. On DAVIS, we freeze the representations and adopt the protocol from [34].

trast vs. 35.2% for MoCo). Finally, our representations also transfer well to the semantic object segmentation task on DAVIS-2016. This dataset covers a rich set of natural image augmentations like viewpoint changes, occlusions, etc., for which our pixel embeddings have learned invariances.

The gains observed across all three benchmarks show that the learned representations are not limited to a specific dataset. We conclude that the use of a mid-level visual prior can be useful for self-supervised representation learning.

4.7. Semi-Supervised Learning

The proposed method can alternatively be used as a pre-training strategy for semantic segmentation. That is, the model is fine-tuned in a semi-supervised way on PASCAL. We use 1%, 2%, 5%, 12.5% and 100% of the `train_aug` split as labeled examples. We initialize our model from supervised pre-training on ImageNet. This weight initialization is commonly used in semantic segmentation. Furthermore, directly fine-tuning a model initialized in the same way serves as a strong baseline. Table 6 shows the results.

The representations generated with our method yield higher performance after fine-tuning, compared to supervised pre-training on ImageNet. This holds true when using both an unsupervised and supervised saliency estimator to predict the object mask proposals. Predictably, the gains become smaller when more labeled examples are available (see also [99]). We conclude that unsupervised learning of pixel embeddings can complement an existing pre-training strategy based on an optimization criterion defined at the

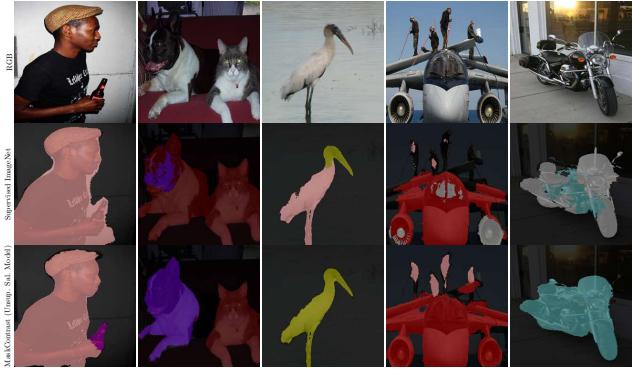


Figure 5. **Qualitative comparison** after fine-tuning on PASCAL using 1 % of labeled data. We use supervised pre-training on ImageNet (middle) or our method (bottom) to initialize the weights before fine-tuning.

Label Fraction	1%	2%	5%	12.5%	100%
ImageNet Classifier Init.	43.4	55.2	62.7	68.4	78.0
+ MaskContrast (Unsup. Sal.)	50.5	57.2	64.5	69.0	78.4
+ MaskContrast (Sup. Sal.)	51.5	59.6	65.3	69.4	78.6

Table 6. **Semi-supervised fine-tuning** on PASCAL (MIoU).

global image- or patch-level. We hope that this observation can spur further research efforts in this direction.

5. Discussion and Limitations

This work presented a general two-step framework based upon a mid-level visual prior for tackling unsupervised semantic segmentation. The proposed setup prevents the model from latching onto low-level image features, a problem present in prior works that relied on end-to-end clustering, proxy tasks or low-level visual cues. Instead, MaskContrast learns pixel embeddings which incorporate more semantically meaningful information (see Figure 4). As a result, we were able to tackle the semantic segmentation task under a fully unsupervised setup on a diverse dataset like PASCAL. Further, experimental evaluation showed that our pixel embeddings have several other interesting properties: the ability for semantic segment retrieval, transfer learning and semi-supervised fine-tuning.

Still, there are some limitations that we did not yet address. The object mask proposals were obtained using a salient object estimator. A disadvantage of using saliency estimation is that only a limited number of object masks can be obtained per image. We argue that alternative ways to mine the object mask proposals can be explored for tackling even more challenging datasets where many objects can exist per image. In particular, we could see additional sensory data [69] or other techniques [62] being used that are better suited for this type of images. The optimization criterion from Equation 3 could then be extended accordingly. Given the viability of our framework, we believe these are interesting research directions.

Algorithm 1 Pseudocode of MaskContrast.

```

# f_q, f_k: encoder-decoder networks for query and key
# queue: dictionary of K prototype keys (CxK)
# m: momentum
# t: temperature
# H, W = height, width of an image x
# P : number of salient pixels in a batch

f_k.params = f_q.params # initialize
for (x, s) in loader:
    # load a batch with N samples and N saliency masks
    # constrain aug s.t. object area > threshold
    x_q, s_q = aug(x, s) # augmented version
    x_k, s_k = aug(x, s) # another augmented version

    q, aux = f_q.forward(x_q) # q: NxCxHxW, aux: NxHxW
    k, _ = f_k.forward(x_k) # k: NxCxHxW

    # salient objects are non zero
    valid_ids = s_q.nonzero() # valid_ids: Px1
    # remap each object to a unique id in {0..N-1}
    s_r = remap(s_q) # s_r: Px1

    # key prototypes: NxC
    p_k = bmm(k.view(N,C,H,W), s_k.view(N,H,W,1))
    p_k = normalize(p_k, dim=1) # L2-normalize
    p_k = p_k.detach() # no gradient to prototypes

    # select embeddings of salient objects: PxC
    q = index_select(q.view(H,WxC), index=valid_ids)

    # positive logits: PxN
    l_pos = mm(q.view(P,C), p_k.view(C,N))

    # negative logits: PxK
    l_neg = mm(q.view(P,C), queue.view(C,K))

    # logits: Px(N+K)
    logits = cat([l_pos, l_neg], dim=1)

    # contrastive loss: positives are the s_r-th
    MaskContrast_loss = CrossEntropyLoss(logits/t, s_r)

    # auxiliary BCE loss to prevent collapse
    aux_loss = BCE(aux, s_q)
    total_loss = MaskContrast_loss + aux_loss

    # SGD update: query network
    total_loss.backward()
    update(f_q.params)

    # momentum update: key network
    f_k.params = m*f_k.params+(1-m)*f_q.params

    # update dictionary
    enqueue(queue, p_k) # enqueue current prototypes
    dequeue(queue) # dequeue earliest prototypes

```

bmm: batch matrix multiplication; mm: matrix multiplication; cat: concatenation; BCE: binary cross-entropy loss; remap: custom function.

Acknowledgment. The authors thankfully acknowledge support by Toyota via the TRACE project and MACCHINA (KU Leuven, C14/18/065). This work is also sponsored by the Flemish Government under the Flemish AI program.

A. Supplementary Materials

The supplementary materials include the pseudocode of our algorithm and the implementation details of the experimental evaluation. Additionally, we provide overclustering results in Section E and qualitative results for semi-supervised fine-tuning in Section G.

B. Pseudo-code

Algorithm 1 shows the pseudo-code of MaskContrast. The saliency masks are obtained by running the public code of existing saliency estimators [52, 64]. The default hyper-parameter settings were used. A PyTorch implementation of our method and pre-computed saliency masks will be made publicly available.

C. Pre-training

This section describes the pre-training setup for the models included in the experiments section of the main paper. In the majority of cases, we were able to use the pre-trained weights made available by the authors of the respective works.

Co-Occurrence. We adopt the training setup from the original work [33]. The features before the output layer of the network are used for the purpose of training a linear classifier and applying K-Means clustering.

Colorization. The pre-trained colorizer from Zhang *et al.* [96] is used. It is argued that the intermediate representations in the network will extract semantic information in order to solve the colorization task. As a consequence, it is non-trivial from what layer we should tap the features to tackle the semantic segmentation task. To resolve this, we tried using features from various intermediate layers, and report the best results when training a linear classifier or applying K-Means.

CMP. We follow the strategy from the colorization task for training a linear classifier or applying K-Means. The pre-trained model from Zhan *et al.* [90] is used.

IIC. We follow the implementation strategy from [36].

Contrastive-Learning Methods. We used the weights from a ResNet-50 model pre-trained on ImageNet. The weights were made available by the authors of the respective works, i.e. the instance discrimination task [80], SWAV [9], MoCo v2 [26] and InfoMin [71]. In some cases, multiple variants of the model were released, e.g. when using different augmentation strategies during training. We chose the best available model each time.

The contrastive learning models were only pre-trained on ImageNet, as we could not see any substantial improvements from further pre-training them on the target dataset, i.e. PASCAL. To obtain dense predictions, we apply dilated convolutions in the last residual block. We use the features from the backbone for training a linear classifier or applying K-Means.

MaskContrast. We use a dilated ResNet-50 model with

DeepLab-v3 head as outlined in the main paper. The final 1×1 convolutional layer is split into two linear heads. The first head predicts the pixel embeddings, while the second head predicts the saliency mask. During linear evaluation, we replace the final layer by a randomly initialized 1×1 convolutional layer. Other details were already provided in the paper.

D. Linear Classifier

This section describes the training setup used for the linear evaluation protocol. We train a 1×1 convolutional layer for 60 epochs using batches of size 16. The complete train set is used during training. We optimize the weights through stochastic gradient descent with momentum 0.9, weight decay 0.0001 and initial learning rate 0.1. The learning rate is reduced to 0.01 after 40 epochs. We found that increasing the train time, or modifying the learning rate did not improve the results.

E. Clustering

This section specifies how to obtain discrete class assignments by clustering the representations using K-Means. We follow the evaluation strategy from [36] to calculate the mean IoU metric. In particular, we first match the predicted clusters with the ground-truth classes using a Hungarian algorithm. We subsequently calculate the mean IoU from the re-assigned clusters and the ground-truth labels. We report the average from five runs.

Contrastive based methods. As described in Section C, we apply K-Means clustering to the backbone features. The cluster assignments are upsampled to match the original image resolution, before applying the Hungarian algorithm.

IIC. No specific post-processing is required. We simply match the predicted clusters with the ground-truth classes following the original work [36].

Proxy-task based methods (Co-Ocurrence, Colorization, CMP). We select the features for applying K-Means as described in Section C. The predictions are up-sampled to match the original image resolution before applying the Hungarian algorithm.

MaskContrast. We compute the mean embeddings of the foreground objects and apply K-Means using the L2-normalized feature vectors. All pixels belonging to the object are assigned the same label as the mean-pixel embedding after clustering. The predictions from the saliency estimation head are used to identify the background class. We match the predictions with the ground-truth classes using the Hungarian algorithm.

Clusters	Init. Sup. Sal.	MoCo v2		Sup.	
		X	✓	X	✓
21		35.0	38.9	41.6	44.2
50		41.4	48.8	46.2	51.4
100		43.3	49.5	47.3	52.5
200		45.0	51.1	48.5	53.6
500		48.1	54.2	51.3	57.0

Table S1. Overclustering on PASCAL with MaskContrast (MIoU). We use MoCo or supervised ImageNet initial weights, and supervised (✓) or unsupervised (X) saliency.

Overclustering results. K-Means does not employ any prior world knowledge, i.e. the ground-truth or target clusters are unknown. Therefore, it is unlikely that the predicted clusters will match the target ones on a complex and imbalanced dataset like PASCAL. To better understand the semantic structure discovered by the embedding space, we apply overclustering. In this case, a many-to-one mapping exists between the predicted and target clusters. Table S1 shows the results. The accuracy improves as we increase the number of predicted clusters. We hypothesize that local neighborhoods in the embedding space contain pixels of the same or visually similar objects, which benefits the performance when overclustering.

F. Semi-Supervised Learning

This section describes the semi-supervised learning setup. In each case, we report the average result for three randomly sampled splits.

ImageNet Pre-Trained Baseline. We load the pre-trained ImageNet weights into a ResNet-50 backbone with dilated convolutions. We use batch size of 8 and stochastic gradient descent with momentum 0.9 and learning rate 0.004 in all data regimes. The learning rate was selected after performing a grid search. Additionally, we explored the use of different parameter groups with specific learning rate, e.g. the decoder used 10 times higher learning rate compared to the encoder. However, this did not result in any further improvements. We include a weight decay term 0.0001. A poly learning rate scheduler is used.

MaskContrast. We use a batch size of 8 and learning rate of 0.004 when fine-tuning with 5%, 12.5% and 100% of the labels. Differently, when using 1% and 2% of the labels, the learning rate is set to 0.001 for all layers in the network, except for the final convolutional layer which uses learning rate 0.1. The latter is well-motivated, since the complete network, including both encoder and decoder, were already pre-trained for the semantic segmentation task. The batch norm stats are frozen. We use stochastic gradient descent with momentum 0.9 and a weight decay term 0.0001. The learning rate is decayed using a poly learning rate scheduler.

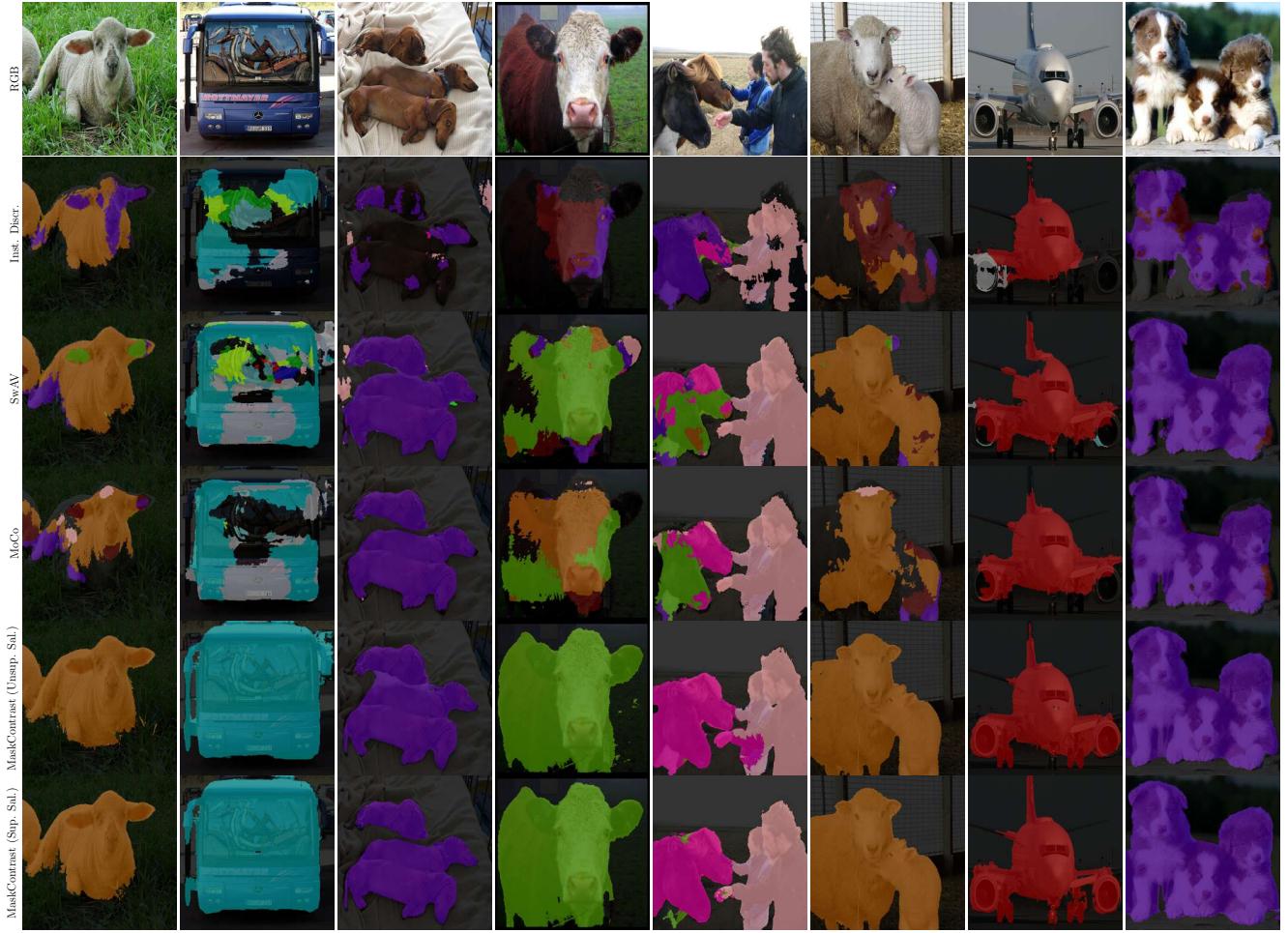


Figure S1. Qualitative comparison of the results after training a linear classifier on PASCAL. We use the MoCo weights to initialize our backbone.

G. Qualitative Results

Figure S1 shows a qualitative comparison when training a linear classifier on top of the pre-trained representations. We compare the representations learned by our method using an unsupervised (5th row) or supervised (6th row) saliency estimator, against the ones from instance discrimination (2nd row) [80], SwAV (3rd row) [9] and MoCo v2 (4th row) [14]. The qualitative results support the claim that our pixel embeddings learn semantically meaningful information.

References

- [1] Pablo Arbelaez, Michael Maire, Charless Fowlkes, and Jitendra Malik. Contour detection and hierarchical image segmentation. *T-PAMI*, 2010. 2, 4, 7
- [2] Yuki M Asano, Mandela Patrick, Christian Rupprecht, and Andrea Vedaldi. Labelling unlabelled videos from scratch with multi-modal self-supervision. In *NeurIPS*, 2020. 2
- [3] Yuki M. Asano, Christian Rupprecht, and Andrea Vedaldi. Self-labelling via simultaneous clustering and representation learning. In *ICLR*, 2020. 2
- [4] Amy Bearman, Olga Russakovsky, Vittorio Ferrari, and Li Fei-Fei. What's the point: Semantic segmentation with point supervision. In *ECCV*, 2016. 1
- [5] Amy Bearman, Olga Russakovsky, Vittorio Ferrari, and Li Fei-Fei. What's the point: Semantic segmentation with point supervision. In *ECCV*, 2016. 1
- [6] Piotr Bojanowski and Armand Joulin. Unsupervised learning by predicting noise. In *ICML*, 2017. 2
- [7] Ali Borji, Ming-Ming Cheng, Huaizu Jiang, and Jia Li. Salient object detection: A benchmark. *IEEE transactions on image processing*, 24(12):5706–5722, 2015. 4
- [8] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *ECCV*, 2018. 2
- [9] 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. 2, 7, 9, 11

- [10] Krishna Chaitanya, Ertunc Erdil, Neerav Karani, and Endre Konukoglu. Contrastive learning of global and local features for medical image segmentation with limited annotations. In *NIPS*, 2020. 2
- [11] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017. 5
- [12] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 2, 4, 5
- [13] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big self-supervised models are strong semi-supervised learners. In *NeurIPS*, 2020. 2
- [14] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 5, 6, 7, 11
- [15] Ming-Ming Cheng, Niloy J. Mitra, Xiaolei Huang, Philip H. S. Torr, and Shi-Min Hu. Global contrast based salient region detection. *T-PAMI*, 2015. 4, 5
- [16] 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*, 2016. 1
- [17] Jifeng Dai, Kaiming He, and Jian Sun. Boxsup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In *ICCV*, 2015. 1, 3
- [18] Bert De Brabandere, Davy Neven, and Luc Van Gool. Semantic instance segmentation with a discriminative loss function. *arXiv preprint arXiv:1708.02551*, 2017. 1
- [19] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *ICCV*, 2015. 2
- [20] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. In *ICLR*, 2017. 2
- [21] Jeff Donahue and Karen Simonyan. Large scale adversarial representation learning. In *NIPS*, 2019. 2
- [22] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 2010. 1, 5
- [23] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *ICLR*, 2018. 2
- [24] 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, et al. Bootstrap your own latent: A new approach to self-supervised learning. In *NeurIPS*, 2020. 2
- [25] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *AISTATS*, 2010. 2, 3, 4
- [26] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020. 2, 4, 5, 6, 7, 9
- [27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5
- [28] Seunghoon Hong, Hyeonwoo Noh, and Bohyung Han. Decoupled deep neural network for semi-supervised semantic segmentation. In *NeurIPS*, 2015. 1
- [29] Qibin Hou, Ming-Ming Cheng, Xiaowei Hu, Ali Borji, Zhuowen Tu, and Philip HS Torr. Deeply supervised salient object detection with short connections. In *CVPR*, 2017. 4
- [30] Wei-Chih Hung, Yi-Hsuan Tsai, Yan-Ting Liou, Yen-Yu Lin, and Ming-Hsuan Yang. Adversarial learning for semi-supervised semantic segmentation. *arXiv preprint arXiv:1802.07934*, 2018. 1
- [31] 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*, 2019. 2, 3, 5, 7
- [32] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. Let there be color! joint end-to-end learning of global and local image priors for automatic image colorization with simultaneous classification. In *SIGGRAPH*, 2016. 2, 3
- [33] Phillip Isola, Daniel Zoran, Dilip Krishnan, and Edward H Adelson. Learning visual groups from co-occurrences in space and time. *arXiv preprint arXiv:1511.06811*, 2015. 2, 3, 7, 9
- [34] Allan Jabri, Andrew Owens, and Alexei A. Efros. Space-time correspondence as a contrastive random walk. In *NeurIPS*, 2020. 5, 8
- [35] Simon Jenni and Paolo Favaro. Self-supervised feature learning by learning to spot artifacts. In *CVPR*, 2018. 2
- [36] Xu Ji, João F Henriques, and Andrea Vedaldi. Invariant information clustering for unsupervised image classification and segmentation. In *ICCV*, 2019. 2, 3, 6, 7, 9, 10
- [37] Bowen Jiang, Lihe Zhang, Huchuan Lu, Chuan Yang, and Ming-Hsuan Yang. Saliency detection via absorbing markov chain. In *ICCV*, 2013. 4
- [38] Longlong Jing and Yingli Tian. Self-supervised visual feature learning with deep neural networks: A survey. *T-PAMI*, 2020. 2
- [39] Anna Khoreva, Rodrigo Benenson, Jan Hendrik Hosang, Matthias Hein, and Bernt Schiele. Simple does it: Weakly supervised instance and semantic segmentation. In *CVPR*, 2017. 1, 3
- [40] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *CVPR*, 2019. 5
- [41] Shu Kong and Charless C Fowlkes. Recurrent pixel embedding for instance grouping. In *CVPR*, 2018. 4
- [42] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Colorization as a proxy task for visual understanding. In *CVPR*, 2017. 2, 3
- [43] Junnan Li, Pan Zhou, Caiming Xiong, Richard Socher, and Steven CH Hoi. Prototypical contrastive learning of unsupervised representations. *arXiv preprint arXiv:2005.04966*, 2020. 2
- [44] Xiaohui Li, Huchuan Lu, Lihe Zhang, Xiang Ruan, and Ming-Hsuan Yang. Saliency detection via dense and sparse reconstruction. In *ICCV*, 2013. 4

- [45] Di Lin, Jifeng Dai, Jiaya Jia, Kaiming He, and Jian Sun. Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. In *CVPR*, 2016. 1
- [46] 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*, 2014. 1, 5
- [47] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 1
- [48] Zhiming Luo, Akshaya Mishra, Andrew Achkar, Justin Eichel, Shaozi Li, and Pierre-Marc Jodoin. Non-local deep features for salient object detection. In *CVPR*, 2017. 4
- [49] Aravindh Mahendran, James Thewlis, and Andrea Vedaldi. Cross pixel optical-flow similarity for self-supervised learning. In *ACCV*, 2018. 2, 3
- [50] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *CVPR*, 2020. 2
- [51] T Nathan Mundhenk, Daniel Ho, and Barry Y Chen. Improvements to context based self-supervised learning. In *CVPR*, 2018. 2
- [52] Tam Nguyen, Maximilian Dax, Chaithanya Kumar Mumadi, Nhung Ngo, Thi Hoai Phuong Nguyen, Zhongyu Lou, and Thomas Brox. Deepups: Deep robust unsupervised saliency prediction via self-supervision. In *NeurIPS*, 2019. 4, 5, 9
- [53] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*, 2016. 2
- [54] Mehdi Noroozi, Hamed Pirsiavash, and Paolo Favaro. Representation learning by learning to count. In *ICCV*, 2017. 2
- [55] Mehdi Noroozi, Ananth Vinjimoor, Paolo Favaro, and Hamed Pirsiavash. Boosting self-supervised learning via knowledge transfer. In *CVPR*, 2018. 2
- [56] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 2, 3, 4
- [57] Yassine Ouali, Céline Hudelot, and Myriam Tami. Autoregressive unsupervised image segmentation. In *ECCV*, 2020. 2, 3
- [58] George Papandreou, Liang-Chieh Chen, Kevin P Murphy, and Alan L Yuille. Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation. In *ICCV*, 2015. 1, 3
- [59] Deepak Pathak, Philipp Krahenbuhl, and Trevor Darrell. Constrained convolutional neural networks for weakly supervised segmentation. In *ICCV*, 2015. 1
- [60] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *CVPR*, 2016. 2
- [61] 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*, 2016. 5
- [62] Pedro OO Pinheiro, Ronan Collobert, and Piotr Dollár. Learning to segment object candidates. In *NeurIPS*, 2015. 4, 8
- [63] Senthil Purushwalkam and Abhinav Gupta. Demystifying contrastive self-supervised learning: Invariances, augmentations and dataset biases. In *NeurIPS*, 2020. 2, 4
- [64] Xuebin Qin, Zichen Zhang, Chenyang Huang, Chao Gao, Masood Dehghan, and Martin Jagersand. Basnet: Boundary-aware salient object detection. In *CVPR*, 2019. 4, 5, 9
- [65] Zhongzheng Ren and Yong Jae Lee. Cross-domain self-supervised multi-task feature learning using synthetic imagery. In *CVPR*, 2018. 2
- [66] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *T-PAMI*, 2000. 3
- [67] Meng Tang, Abdelaziz Djelouah, Federico Perazzi, Yuri Boykov, and Christopher Schroers. Normalized cut loss for weakly-supervised cnn segmentation. In *CVPR*, 2018. 1
- [68] Meng Tang, Federico Perazzi, Abdelaziz Djelouah, Ismail Ben Ayed, Christopher Schroers, and Yuri Boykov. On regularized losses for weakly-supervised cnn segmentation. In *ECCV*, 2018. 1
- [69] Hao Tian, Yuntao Chen, Jifeng Dai, Zhaoxiang Zhang, and Xizhou Zhu. Unsupervised object detection with lidar clues. *arXiv preprint arXiv:2011.12953*, 2020. 8
- [70] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019. 2
- [71] Yonglong Tian, C. Sun, Ben Poole, Dilip Krishnan, C. Schmid, and Phillip Isola. What makes for good views for contrastive learning. *arxiv preprint arXiv:2005.10243*, 2020. 2, 7, 9
- [72] Koen EA Van de Sande, Jasper RR Uijlings, Theo Gevers, and Arnold WM Smeulders. Segmentation as selective search for object recognition. In *ICCV*, 2011. 4
- [73] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc Van Gool. Scan: Learning to classify images without labels. In *ECCV*, 2020. 2, 3, 4
- [74] Paul Vernaza and Mamnoon Chandraker. Learning random-walk label propagation for weakly-supervised semantic segmentation. In *CVPR*, 2017. 1
- [75] Jingdong Wang, Huaizu Jiang, Zejian Yuan, Ming-Ming Cheng, Xiaowei Hu, and Nanning Zheng. Salient object detection: A discriminative regional feature integration approach. *IJCV*, 2017. 4
- [76] Lijun Wang, Huchuan Lu, Yifan Wang, Mengyang Feng, Dong Wang, Baocai Yin, and Xiang Ruan. Learning to detect salient objects with image-level supervision. In *CVPR*, 2017. 5
- [77] Tiantian Wang, Ali Borji, Lihe Zhang, Pingping Zhang, and Huchuan Lu. A stagewise refinement model for detecting salient objects in images. In *ICCV*, 2017. 4
- [78] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *ICML*, 2020. 5

- [79] Wenguan Wang, Qiuxia Lai, Huazhu Fu, Jianbing Shen, Haibin Ling, and Ruigang Yang. Salient object detection in the deep learning era: An in-depth survey. *arXiv preprint arXiv:1904.09146*, 2019. 4
- [80] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *CVPR*, 2018. 2, 4, 7, 9, 11
- [81] Saining Xie and Zhuowen Tu. Holistically-nested edge detection. In *ICCV*, 2015. 2, 7
- [82] Jia Xu, Alexander G Schwing, and Raquel Urtasun. Learning to segment under various forms of weak supervision. In *CVPR*, 2015. 1, 3
- [83] Q. Yan, L. Xu, J. Shi, and J. Jia. Hierarchical saliency detection. *CVPR*, 2013. 4
- [84] Xueting Yan, Ishan Misra, Abhinav Gupta, Deepti Ghadiyaram, and Dhruv Mahajan. Clusterfit: Improving generalization of visual representations. In *CVPR*, 2020. 2
- [85] Yanchao Yang, Brian Lai, and Stefano Soatto. Time-supervised primary object segmentation. *arXiv preprint arXiv:2008.07012*, 2020. 4
- [86] Yanchao Yang, Antonio Loquercio, Davide Scaramuzza, and Stefano Soatto. Unsupervised moving object detection via contextual information separation. In *CVPR*, 2019. 4
- [87] Mang Ye, Xu Zhang, Pong C Yuen, and Shih-Fu Chang. Unsupervised embedding learning via invariant and spreading instance feature. In *CVPR*, 2019. 2
- [88] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015. 5
- [89] Xiaohang Zhan, Ziwei Liu, Ping Luo, Xiaou Tang, and Chen Change Loy. Mix-and-match tuning for self-supervised semantic segmentation. In *AAAI*, 2018. 2
- [90] Xiaohang Zhan, Xingang Pan, Ziwei Liu, Dahua Lin, and Chen Change Loy. Self-supervised learning via conditional motion propagation. In *CVPR*, 2019. 2, 3, 7, 9
- [91] Dingwen Zhang, Junwei Han, and Yu Zhang. Supervision by fusion: Towards unsupervised learning of deep salient object detector. In *CVPR*, 2017. 4
- [92] Jing Zhang, Deng-Ping Fan, Yuchao Dai, Saeed Anwar, Fatemeh Sadat Saleh, Tong Zhang, and Nick Barnes. Ucnet: uncertainty inspired rgb-d saliency detection via conditional variational autoencoders. In *CVPR*, 2020. 4
- [93] Jing Zhang, Tong Zhang, Yuchao Dai, Mehrtash Harandi, and Richard Hartley. Deep unsupervised saliency detection: A multiple noisy labeling perspective. In *CVPR*, 2018. 4
- [94] Liheng Zhang, Guo-Jun Qi, Liqiang Wang, and Jiebo Luo. Aet vs. aed: Unsupervised representation learning by auto-encoding transformations rather than data. In *CVPR*, 2019. 2
- [95] Pingping Zhang, Dong Wang, Huchuan Lu, Hongyu Wang, and Xiang Ruan. Amulet: Aggregating multi-level convolutional features for salient object detection. In *CVPR*, 2017. 4
- [96] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *ECCV*, 2016. 2, 3, 7, 9
- [97] Xiao Zhang and Michael Maire. Self-supervised visual representation learning from hierarchical grouping. *NeurIPS*, 2020. 2, 3, 5, 6, 7
- [98] Wangjiang Zhu, Shuang Liang, Yichen Wei, and Jian Sun. Saliency optimization from robust background detection. In *CVPR*, 2014. 4
- [99] Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin D Cubuk, and Quoc V Le. Rethinking pre-training and self-training. In *NeurIPS*, 2020. 8
- [100] Wenbin Zou and Nikos Komodakis. Harf: Hierarchy-associated rich features for salient object detection. In *ICCV*, 2015. 4