A CLUSTERING BASELINE FOR OBJECT-CENTRIC REPRESENTATIONS

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028 029

031

Paper under double-blind review

ABSTRACT

Object-centric learning aims to discover and represent visual entities as a small set of object embeddings and masks, which can be later used for downstream tasks. Recent methods for object-centric learning build upon vision foundation models trained with self supervision because of the rich semantic features they produce. However, they often involve additional training to optimize for object mask accuracy for a specific granularity of objects on a test dataset, while overlooking the evaluation of the quality of the object embeddings which is arguably more important. In this work, we demonstrate how to discover objects and parts with a simple multi-scale application of k-means to the features of an off-the-shelf backbone. Our method is fast and flexible, produces interpretable masks, preserves the quality of the backbone embeddings, does not require additional training, and can capture different part/whole structures. We evaluate the quality of the obtained representation on a variety of downstream tasks including scene classification and action recognition in videos, showing that it surpasses the performance of fine-tuned object-centric learning methods. Object masks produced by our method also effectively capture real-world objects and parts at various granularity, with comparable quality to specialized methods when evaluated on unsupervised segmentation benchmarks. These results suggest rethinking the current approach to object-centric learning, with a greater focus on the quality of the representation.

1 INTRODUCTION

Self-supervised learning (SSL) enables computer vision models to learn from vast amounts of 033 unlabeled data with minimal supervision. At its inception, SSL relied on handcrafted pretext tasks to 034 learn useful representations from data, such as predicting augmentations (Zhang et al., 2016; Gidaris et al., 2018), reconstructions (Kingma & Welling, 2014), or contrastive assignments (Chen et al., 2020; Grill et al., 2020). Over time, the field has shifted towards more streamlined architectures (Dosovitskiy et al., 2020) and training objectives (Caron et al., 2021; Radford et al., 2021; He et al., 037 2022; Assran et al., 2023; Oquab et al., 2023), which have shown impressive performance on a wide range of downstream tasks. However, existing SSL models focus primarily on learning global or local representations, expressed respectively as a single embedding – a vector – that represents the entire 040 image (Chen et al., 2020; Grill et al., 2020; Caron et al., 2021) or a dense feature map that represents 041 local image patches (Zhou et al., 2022; Oquab et al., 2023). These two levels of representation are 042 useful for a wide range of tasks, such as image classification, object detection, and segmentation. 043 However, they lack an understanding of the structure in the represented parts, objects, and groups 044 of entities that we typically use to describe the world (Whitehead, 1985; Spelke & Kinzler, 2007; Johnson-Laird, 2010). Bridging the gap between the two levels, *object-centric* representations aim to offer a more structured and interpretable representation of the world (Lake et al., 2017; Greff 046 et al., 2020; Hinton, 2021; LeCun, 2022). Informally, object-centric representations are a set of 047 features associated to distinct parts or objects in an image. Downstream tasks can benefit from such a 048 representation in terms of compute, interpretability, and memory usage. 049

Learning object-centric representations remains at present an open challenge. Supervised or weakly supervised approaches struggle with the compositional explosion of annotations. For example, SAM (Kirillov, 2023) is trained explicitly for object segmentation, but requires extensive annotations of
 every object in an image. Text-supervised models like CLIP (Radford et al., 2021) can be used for open-vocabulary segmentation, assuming that the training captions contained a detailed description

054 of every possible item. On the other hand, self-supervised approaches need to acquire a notion of 055 objectness from the data itself. Such general notion is by itself not trivial, as the definition of objects 056 can be ambiguous and task-specific (e.g. distinguishing between a car and a truck, or the wheel of a car 057 and the car itself), which in turn makes it harder to define clear benchmarks and evaluation metrics. To 058 address these challenges, early object-centric methods made use of specialized architectures (Hinton et al., 2018; Locatello et al., 2020) and training objectives (Burgess et al., 2019; Greff et al., 2019) with built-in inductive biases to better steer the model towards the desired behavior. More recent 060 methods build upon pre-trained vision backbones - more powerful and general than ad-hoc models -061 and learn additional adapters to extract object-centric representations (Vo et al., 2020; Hamilton et al., 062 2021; Seitzer et al., 2023). Even though object-centric learning aims at producing both meaningful 063 object-centric representation and object masks, these works have mostly been focusing on generating 064 high-quality object masks, and evaluate them for the unsupervised segmentation task. This work aims 065 to re-establish the two-sided nature of object-centric learning, with focus on the quality of both the 066 object masks and the representation. 067

To this end, we consider a simple yet effective clustering-based method for obtaining high-quality 068 object-centric representations. Such representations are crucial for downstream tasks that require 069 an in-depth understanding of the objects in a scene, their position and relations. Section 3 describes 070 our method: Given patch features extracted from state-of-the-art vision models such as DINOv2 071 (Oquab et al., 2023), we apply multiple rounds of k-means (Lloyd, 1982) to obtain a small set of 072 cluster centroids that represent objects at different levels of granularity, each associated with an 073 object mask that localizes it in the image. In Section 4, we provide a comprehensive evaluation of 074 our approach, focusing both on unsupervised segmentation tasks and on assessing the quality of 075 the learned representations, which is often overlooked in the object-centric literature. We show that our simple approach can rival existing object-centric methods without complex pre-training losses, 076 lengthy fine-tuning runs, or dataset-specific hyperparameters. Our clustering-based representations 077 yields excellent performance in downstream tasks such as image and video classification, which 078 require compositional understanding of complex scenes and frame sequences, at a much lower 079 computational cost than processing the dense representation. Following previous work (Locatello 080 et al., 2020; Hénaff et al., 2022), we also evaluate how well our object masks correspond to real-world 081 objects on common unsupervised segmentation benchmarks. We demonstrate that even though 082 k-means assignment masks are not as clean as those produced by complex, specialized methods, 083 they are more than capable of capturing objects and parts at different levels of granularity, yielding a 084 flexible and interpretable representation of the image.

085 086 087

880

2 RELATED WORK

Self-supervised learning. The goal of self-supervised learning (SSL) is to learn useful representations 089 from unlabeled collections of images. At a high level, this is achieved by training a model to solve 090 a pretext task with no external annotation through which the model is encouraged to learn useful 091 features that generalize well to other tasks. Early examples include denoising (Vincent et al., 2008), 092 colorization (Zhang et al., 2016), or inverting geometric augmentations (Noroozi & Favaro, 2016; 093 Gidaris et al., 2018). Modern SSL methods can be categorized into two families depending on 094 the pretext task and the domain in which their losses are computed. Generative models are based on the auto-encoder paradigm (Vincent et al., 2008), in which the model tries to reconstruct an 096 image from its altered versions (Kingma & Welling, 2014; He et al., 2022; El-Nouby et al., 2024). 097 Latent-space methods avoid pixel reconstruction and formulate their loss directly in the model's 098 embedding space. Notable examples include contrastive learning (van den Oord et al., 2018; Chen et al., 2020; Grill et al., 2020), self-distillation at the global or patch level (Caron et al., 2021; Zhou 099 et al., 2022), or predictive methods (Assran et al., 2023). SSL features exhibit an ability to capture 100 structure and semantics of the input (Caron et al., 2021), which can be further developed using explicit 101 object-centric losses (Hénaff et al., 2022), additional tokens (Darcet et al., 2023), or used as-is for 102 object discovery, as discussed below. This work builds on DINOv2 (Oquab et al., 2023) and seeks to 103 efficiently extract a set-structured representation of objects and parts from the learned patch features. 104

Object discovery. Given an unlabelled set of images, object discovery aims at grouping images that
 contain similar objects together while also localizing these objects in forms of bounding boxes or
 masks. Early works focus on discovering object categories, using techniques such as probabilistic
 modeling (Weber et al., 2000; Sivic et al., 2005; Russell et al., 2006) or non-negative matrix fac-



Figure 1: Dense features produced by a strong SSL backbone are clustered into distinct objects and parts using the classic k-means algorithm. Top: our approach: we cluster DINOv2 features using k-means for $k \in \{2, 4, 8\}$. The segment boundaries roughly align with object boundaries, capturing part-whole hierarchies as exemplified by the laptop, screen and keyboard segments. Bottom left: PCA visualization of DINOv2 features. Bottom right: cutouts of segments and the closest ImageNet sample retrieved using the centroid as query.

126

127

128

129

132 torization (Tang & Lewis, 2008), or consider only simple scenarios with a few object classes (Zhu et al., 2012; Zhang et al., 2015). Cho et al. (2015) studies object discovery in the more challenging 133 in-the-wild image domain, proposing an algorithm that alternates between finding similar images and 134 localizing objects. This approach is further developed into different optimization formulations (Vo 135 et al., 2019; 2020; 2021; Choudhury et al., 2021) that could discover multiple objects per image and 136 handle million-scale datasets. With the advent of powerful SSL backbones which already excel at 137 at grouping similar images (Caron et al., 2021, DINO), the focus shifts entirely to the localization 138 task. LOST (Siméoni et al., 2021) is the first to leverage DINO for finding objects with a graph-based 139 method. Following works improve upon LOST with spectral clustering (Wang et al., 2023b; Melas-140 Kyriazi et al., 2022; Rambhatla et al., 2023; Wang et al., 2023a) or mask/feature correspondence 141 distillation (Hamilton et al., 2021; Van Gansbeke et al., 2022; Zadaianchuk et al., 2022). From a 142 methodological perspective, this work is also based on processing dense features learned by SSL 143 backbones, though the focus is on representing objects rather than localizing them exactly. As most 144 object discovery methods do not output object embeddings, they are not directly comparable to our 145 work, which is instead more aligned with the object-centric learning literature.

146 **Object-centric learning.** This line of work evolved in parallel to object discovery and shares similar 147 goals, though the focus is on the learned representations rather than object localization. Many 148 object-centric methods rely on compositional generative models (Greff et al., 2019; Burgess et al., 149 2019) whose latent space is structured around a small set of vectors, each representing a distinct 150 object or part of the scene. In their seminal work, Locatello et al. (2020) introduce the term *slots* to describe the attention-based set-structured bottleneck of an auto-encoder that learns to group spatial 151 features into semantically-coherent regions. Subsequent work improve the bottleneck mechanism 152 (Chang et al., 2022) and the decoder (Singh et al., 2022a; Wu et al., 2023; Jiang et al., 2023), with 153 extensions to videos (Kipf et al., 2021; Singh et al., 2022b; Wu et al., 2022; Sun et al., 2023) and 3D 154 data (Stelzner et al., 2021; Sajjadi et al., 2022). In an effort to scale from simple synthetic data to 155 complex natural images, more recent works replace the pixel reconstruction objective with optical 156 flow prediction (Kipf et al., 2021) or reconstruction of semantic features from pretrained SSL models 157 (Seitzer et al., 2023; Kakogeorgiou et al., 2024). In this work, we borrow ideas and terminology 158 from the object-centric literature, but we depart from the training-based approach, which requires 159 specialized architectures and dataset-specific hyperparameters. We instead probe whether pre-trained 160 SSL models already expose object-centric representations that can be extracted with a simpler, more flexible and more interpretable clustering-based algorithm. 161

¹⁶² 3 METHOD

164 3.1 GOAL AND PROBLEM SETUP

Given a set of unlabeled images, object-centric learning aims to extract a representation that captures distinct objects or parts in a scene, possibly learned directly from the data without supervision. In line with open-ended nature of SSL methods, we use the term "object" loosely, including all parts and subparts that may be relevant depending on the context. For generality, we also allow for overlaps between the objects, to account for occlusion and part-whole object hierarchies, *e.g.* one may want to capture not only a shirt as a whole, but also its pockets and the buttons on it.

Formally, we aim to learn a function $f : \mathbb{R}^{H \times W \times 3} \to (S_N, \mathcal{M}_N)$ that maps an image $x \in \mathbb{R}^{H \times W \times 3}$ to a set of **object vectors** $S_N = \{s_1, \ldots, s_N \mid s_n \in \mathbb{R}^D\}$ and their corresponding **assignment masks** $\mathcal{M}_N = \{m_1, \ldots, m_N \mid m_n \in \{0, 1\}^{H \times W}\}$. The role of object vectors, or *embeddings*, and masks is complementary: The former represents the "content" of the objects while the latter allows to localize them in the image. This set-structured representation complements the standard output of self-supervised models: the *global* representation $g \in \mathbb{R}^D$, which captures the image content as a whole, and the *dense* representation $p \in \mathbb{R}^{H \times W \times D}$, which preserves the spatial dimensions of the input by associating each coordinate to a feature vector.

The number of relevant objects in a scene depends on the data domain, the tasks at hand, as well as the users' subjective notion. For simplicity, we assume that N is fixed a priori to capture all relevant objects and possibly more. In practice, this high-recall behavior fits the scenario where object tokens are fed to a downstream model, *e.g.* a classifier or a vision-language model. Processing a few more tokens than strictly necessary remains advantageous compared to the alternative of using all patch tokens. Also, hardware accelerators are better optimized to handle constant-size inputs.

186 187

3.2 BACKGROUND: SLOT-BASED METHODS

We briefly characterize *state-of-the-art* methods for object-centric learning such as DINOSAUR (Seitzer et al., 2023) and its follow-up SPOT (Kakogeorgiou et al., 2024). These models are implemented as an auto-encoder with a slot-attention bottleneck (Locatello et al., 2020) built on top of a pre-trained SSL backbone (Caron et al., 2021). At the bottleneck, a set of query vectors compete over spatial locations with a specialized attention mechanism to compress backbone features p into "slots" S_N . After training, the slot vectors are used as object embeddings, and corresponding masks are obtained through an arg max over the attention maps, associating each pixel to a single object.

195 Slot-based auto-encoders have proven effective in capturing objects in images, though not without 196 limitations. Limited flexibility: While some formulations allow for a variable (Locatello et al., 197 2020) or adaptive (Fan et al., 2024) number of slots, this feature is not adopted by recent works that instead train each model with a fixed number of queries (Seitzer et al., 2023; Kakogeorgiou 198 et al., 2024). No-overlaps: Predicting a single class for a given pixel can not express part-whole 199 hierarchies, e.g. "wheel" and "car", which are key to understanding real world scenes (Lin et al., 2014; 200 Zhou et al., 2017). As a consequence, a trained object-centric model is implicitly biased toward a 201 certain granularity and object density in images, which may not generalize to other tasks and datasets. 202 **Entangled embeddings:** Due to the reconstruction objective, slot vectors must contain both semantic 203 and position information, which degrades downstream performance compared to the SSL backbone 204 they are fine-tuned from (see Section 4.1). In the following section, we depart from the slot-based 205 formulation and propose a simpler method that is able to capture a more general representation of 206 objects and parts without training or fine-tuning.

207 208

209

3.3 CLUSTERING DENSE SSL FEATURES

Our approach stems from two observations. First, the dense feature maps of modern SSL backbones contain a rich semantic representation of the input image (Oquab et al., 2023; Darcet et al., 2023). Referencing the PCA in Figure 1, the information about objects and parts is already present, we only need to extract it. Second, the slot-based methods act as a soft and parametric approximation of *k*-means clustering (Locatello et al., 2020; Chang et al., 2022). The substantial training effort required by slot-based methods is only needed to align the learned clustering to the objects that are annotated in the test datasets. Therefore, we investigate k-means clustering (Lloyd, 1982) as a simple, interpretable, flexible, and learning-free method for extracting object-centric representations.

Given the patch features produced by an SSL backbone, k-means groups them into K clusters identified by $C_K = \{c_1, \dots, c_K\}$ centroids. Without additional processing, we directly use the centroids as object vectors and the cluster assignments and object masks. To capture objects and parts at various scales, we run k-means multiple times with different values of K, as exemplified in the top row of Figure 1. In our experiments, we build a set of object vectors that includes the global representation g and all the k-means centroids with K chosen on a geometric progression:

224 225

$$\mathcal{S}_N = \{\boldsymbol{g}\} \cup \mathcal{C}_1 \cup \mathcal{C}_{2^1} \cup \ldots \cup \mathcal{C}_{2^{\log_2(N)-1}},\tag{1}$$

where N is a power of 2, and g and C_1 are associated to dummy masks that cover the whole image. The choice of N and K values can of course be refined if domain-specific knowledge is available.

228 Compared to learning-based approaches, this method based on k-means clustering offers several 229 advantages. Flexibility: The algorithm requires no additional training or fine-tuning, and can be applied to any off-the-shelf model, allowing for quick experimentation with different backbones 230 and number of objects. Granularity: Running multiple k-means accounts for the complexity of 231 real-world images where each pixel can be associated to multiple objects and parts. Semantic 232 **embeddings:** Each k-means centroid is obtained by averaging the patch features assigned to it, 233 therefore preserving the quality of the SSL representation. Since object vectors "live" in the same 234 space as the dense features, we can easily visualize their semantic content by retrieving related images 235 as in Figure 1. 236

Resolution and hierarchical k-means. For simplicity, Section 3.1 does not distinguish between 237 the input resolution and the resolution of the dense features. In practice, all modern SSL backbones 238 introduce a significant downsampling, e.g. a DINO ViT-B/16 (Caron et al., 2021) produces a 14×14 239 feature map from its native resolution of 224×224 . Applying k-means clustering at this resolution 240 has no noticeable effect on the object vectors, but results in "blocky" object masks. Unless specified 241 otherwise, we use the native resolution when evaluating representation quality (Section 4.1), and 242 process the images at $4 \times$ the backbone resolution for evaluating segmentation (Section 4.2). In this 243 case, the ViT-B/16 feature map would now be 56×56 , so we apply k-means in a hierarchical manner: 244 an initial clustering step with K = 256 followed by another step at the desired K. Hierarchical 245 k-means also has the desirable effect of reducing the bias towards equally-sized clusters, finding valid 246 object masks even when there is an important imbalance between object sizes (Vo et al., 2024).

247 248

249

4 EXPERIMENTS

250 We evaluate object-centric representations along two axes: the quality of the object vectors as a 251 compact scene representation for downstream tasks (Section 4.1), and the ability to associate each 252 vector to an actual object in the image (Section 4.2). The former is directly related to the goal of 253 object-centric learning, i.e. turning a visual input into a small set of tokens that represent a composition 254 of entities in a scene. The latter ensures that the learned representations are indeed aligned with a 255 human notion of "objects", and are not just a set of vectors that happens to be useful for downstream tasks. We point out that the object-centric literature often focuses on the latter, pushing to improve 256 unsupervised segmentation metrics, while overlooking the usefulness of the learned representations. 257 As segmentation is not the only goal, we give equal importance to both aspects in our evaluations. 258

259 260

4.1 EVALUATING OBJECT-CENTRIC REPRESENTATIONS

In the SSL literature, the standard practice to evaluate the quality of global $g \in \mathbb{R}^D$ and dense representations $p \in \mathbb{R}^{H \times W \times D}$, is to probe them with linear classifiers on a variety of tasks. This approach assumes linear separability in the feature domain and introduces a minimal number of trainable parameters, so to focus on the representation itself. We evaluate the object-centric representations in a similar spirit, but adapt the classifier to take as input an unordered set of vectors S, which is described in Equation 1 for our approach and corresponds to slot embeddings for SPOT (Kakogeorgiou et al., 2024). Specifically, the classifier is implemented as a permutation-invariant attention pooling model:

$$y = \boldsymbol{W}^{\text{out}} \left(\sum_{\boldsymbol{s} \in \mathcal{S}} \frac{\exp(\boldsymbol{q}^T \boldsymbol{W}^{\text{key}} \boldsymbol{s})}{\sum_{\boldsymbol{s}' \in \mathcal{S}} \exp(\boldsymbol{q}^T \boldsymbol{W}^{\text{key}} \boldsymbol{s}')} \boldsymbol{W}^{\text{value}} \boldsymbol{s} \right),$$
(2)

where $q \in \mathbb{R}^D$, $W^{\{\text{key,value}\}} \in \mathbb{R}^{D \times D}$, and $W^{\text{out}} \in \mathbb{R}^{C \times D}$ are learnable parameters trained by gradient descent, and $y \in \mathbb{Y} \subset \mathbb{R}^C$ is a single- or multi-label prediction depending on the task.

273 4.1.1 IMAGE CLASSIFICATION274

Datasets and tasks. As a first benchmark, we use ImageNet (Deng et al., 2009), the *de-facto* standard 275 for image classification in self-supervised learning. ImageNet pictures typically contain a single 276 predominant object and do not require compositional understanding, yet it is a useful sanity check to 277 ensure that the object vectors preserve the quality of the SSL backbone. Moving to more complex 278 images, we evaluate on Places205 (Zhou et al., 2014) and SUN397 (Xiao et al., 2010), two scene 279 classification tasks that require more fine-grained understanding of objects and relationships. For these 280 three datasets, we report top-1 accuracy. Furthermore, we construct two multi-label classification 281 tasks from the CLEVR (Johnson et al., 2017) and COCO 2017 (Lin et al., 2014) datasets, where the 282 goal is to predict a multi-hot vector of all object categories present in an image, e.g. "red large metal cube" or "wine glass". These tasks probe whether the object vectors contain sufficient information 283 284 about all annotated objects and their attributes. For these two tasks, we report mean average precision (mAP) over 96 categories for CLEVR and 80 for COCO. 285

286 Models. We apply our method to DINO ViT-B/16 (Caron et al., 2021) and DINOv2 ViT-B/14 (Darcet 287 et al., 2023), using multiple rounds of k-means to extract S_8 or S_{16} as described in Section 3.3. 288 For comparison, we also evaluate a linear classifier on the CLS token of these models. At the 289 time of writing, the state-of-the-art model for object-centric learning is SPOT (Kakogeorgiou et al., 290 2024). We evaluate the open-source checkpoint based on a DINO ViT-B/16 backbone and fine-tuned on COCO using 7 slots, as well as an 8-slot and a 16-slot variant trained using the official code. 291 Additionally, we obtain multi-scale SPOT embeddings by ensembling multiple models with 2, 4, 8 292 slots, the backbone CLS token as well as the average output patch. We also attempted to train SPOT 293 with DINOv2, but it failed to converge. The second best, DINOSAUR (Seitzer et al., 2023), does not 294 provide open-source checkpoints for evaluation. 295

Results. We observe in Table 1 that our approach largely outperforms SPOT, especially on the more 296 object-centric tasks. Notably, all classifiers trained on slot embeddings perform worse than a linear 297 classifier trained on the CLS token of the original backbone. We attribute this to the auto-encoding 298 formulation of slot-based methods, which compresses both semantic and positional information into 299 small-dimensional vectors, leading to a degradation of the backbone embeddings (Section 3.2). On 300 the other hand, our clustering-based approach requires no additional training and improves over 301 the CLS token performance, as the attention pooling classifiers can leverage additional structured 302 information of each scene. In Figure 2, we assess the scaling trends for our method with respect to 303 the number of object vectors: S_4 , S_8 , S_{16} , or S_{32} ; and the size of the underlying DINOv2 backbone: 304 base, large and giant. We compare against linear classifiers trained on the CLS token alone, and with 305 attention-pooling classifiers trained on all patches (256 tokens) or on a square grid of average-pooled 306 patches, e.g. 3×3 or 4×4 . For simpler single-label classification tasks, performance saturates quickly 307 with the number of objects and does not spectacularly improve over the CLS token, suggesting that a 308

Table 1: Object-centric classification: top-1 accuracy for single-label tasks, *i.e.* ImageNet, SUN 397, and Places 205, and mean average precision (mAP) for multi-label tasks, *i.e.* CLEVR and COCO 2017. The object-centric methods are grouped based on their SSL backbones, namely DINO and DINOv2, of which we report the performance of a linear classifier trained on the CLS token as a reference. The 7-slot SPOT model is the open-source release, all other SPOT models are trained by us using the official code. Our method requires no additional training, as it amounts to running multiple *k*-means clustering on the patch embeddings.

Tokens CLS 7 8 16 CLS+1+2+4	ImageNet 78.2 66.8 - 11.4 67.6 - 10.6 69.2 - 9.0 71.4 - 42	SUN 397 66.9 59.7 -12 59.6 -13 60.3 -66	Places 205 56.0 51.1 -48 50.8 -52 51.5 -45	CLEVR 78.2 70.8 -74 73.7 -45 81.5 +33	COCO 17 70.6 67.6 - 3.0 67.7 - 2.9 67.7 - 2.9
CLS 7 8 16 CLS+1+2+4	78.2 66.8 -114 67.6 -106 69.2 -90 71 4 62	66.9 59.7 -72 59.6 -73 60.3 -66	56.0 51.1 -48 50.8 -52 51.5 -45	78.2 70.8 -74 73.7 -45 81.5 +33	70.6 67.6 - 3.0 67.7 - 2.9 67.7 - 2.9
7 8 16 CLS+1+2+4	66.8 -11.4 67.6 -106 69.2 -90 71.4 -69	59.7 -72 59.6 -73 60.3 -66	51.1 -48 50.8 -52 51.5 -45	70.8 -74 73.7 -45 81.5 +33	67.6 -30 67.7 -29 67.7 -29
8 16 CLS+1+2+4	67.6 - 106 69.2 - 10 71.4 - 68	59.6 -73 60.3 -66	50.8 -52 51.5 -45	73.7 -45 81.5 +33	67.7 .29 67.7 .29
16 CLS+1+2+4	69.2 71.4	60.3	51.5	81.5 +33	67.7 -29
CLS+1+2+4	714	50.4			
		JY.4 .15	54.6	75.9	69.4 -12
CLS+1+2+4+8	72.7 -55	59.2	55.5 .05	84.4 +62	70.8 +02
CLS+1+2+4	78.1	66.0	58.8 +28	85.6 +74	72.3 +17
CLS+1+2+4+8	78.2	66.0	59.2 +32	89.7 + 11.5	72.6 +20
CLS	83.9	77.4	64.4	74.6	77.0
CLS+1+2+4	84.3	77.9 + 0.5	65.9 +15	84.1 +95	81.3
	011	77 7	66 1	88 2	82.3
	CLS+1+2+4	CLS+1+2+4 84.3 +04	CLS 83.9 77.4 CLS+1+2+4 84.3 ••• 77.9 ••5 CLS+1+2+4+8 84.6 77.7	CLS 83.9 77.4 04.4 CLS+1+2+4 84.3 + 77.9 + 65.9 + 15 CLS+1+2+4+8 84.6 + 77.7 + 66.1 +	CLS+1+2+4 84.3 \cdot 64.7 \cdot 77.9 \cdot 65.9 \cdot 15 84.1 \cdot 85 CLS+1+2+4+8 84.6 \cdot 97.7 \cdot 96 66.1 \cdot 17 88.2 \cdot 196



Figure 2: Scaling trends w.r.t. the number of object tokens and backbone size. Using open-source DINOv2
 models of different sizes, we compare the classification performance when training on: the CLS token, all
 16 × 16 patch tokens, a square grid of average-pooled patches *e.g.* 3 × 3, or an increasing number of object
 tokens extracted by running multiple *k*-means on the dense features. The object-centric representation strikes a
 gradual trade-off between the lower performance of the CLS token and the top-line performance of all patches.
 At equivalent token counts, our clustering-based representation always outperforms the pooling-based one,
 highlighting the importance of capturing objects/parts over simply increasing the number of input vectors.

global representation is mostly sufficient. In contrary, for more complex scenes, our approach leads
 to sizable gains, efficiently bridging the gap between the global and dense top-line representation.
 Furthermore, the comparison with average-pooled patches highlights the importance of capturing
 objects/parts over simply increasing the number of input vectors. Unsurprisingly, our method scales
 well with the size of the backbone: Since object vectors are averaged from the patch embeddings, an
 improvement in the backbone will directly benefit the object-centric representation.

3463474.1.2 VIDEO ACTION RECOGNITION

Object-centric representations are particularly valuable in scenarios where large spatial and temporal
 resolutions are crucial to performance but processing a patch-based representation would be computa tionally expensive. Focusing on videos, we evaluate whether a small set of vectors corresponding to
 objects/parts can offer a compact and expressive representation for action recognition.

352 Datasets and features. We consider two datasets for action recognition: Kinetics-400 (Carreira 353 & Zisserman, 2018, K400) and Something-Something V2 (Goyal et al., 2017, SSv2). For each equally-spaced frame of a $T \times H \times W$ video, we apply a DINOv2 ViT-L/14 backbone (Darcet et al., 354 2023) and extract either: a) the CLS token alone, resulting in T tokens; b) our S_N representation, 355 resulting in TN tokens; c) all $THW/14^2$ patches. We then add temporal embeddings to identify 356 each frame and train a two-blocks transformer classifier on the token sequence. To demonstrate 357 the efficiency and flexibility of our method, we sweep over several resolutions: temporal with 358 $T \in \{8, 16\}$ frames, spatial with $H, W \in \{224, 448\}$ pixels, and "object" with $N \in \{8, 16\}$. With 359 this range of parameters, a video can be represented with as few as 8 tokens, up to 16384, with 360 obvious effects on performance. 361



Figure 3: Kinetics 400 and Something-Something v2 classification accuracy vs. number of tokens used to represent the input (log scale). We consider four options to represent each frame: using only the CLS token, using all tokens (maximal performance, expensive), or using our S_8 or S_{16} representation. We experiment with $T \in \{8, 16\}$ frames, and $\{224^2, 448^2\}$ resolution. Our S_N representation recovers most of the performance with orders of magnitude fewer tokens than the most expensive regimes, especially for Something-Something v2 where recognizing relevant objects is critical to the task. See Section 4.1.2 for more details.

 378

 379

 380

 381

 382

 383

Figure 4: Qualitative comparison on two images between the masks obtained with k-means (middle) and SPOT, a specialized object-centric method (right). SPOT mask don't overlap and have cleaner object boundaries. K-means mask have noisier boundaries, but capture objects at different granularity and overlapping parts.

Results. Figure 3 reports the classification accuracy versus the total number of input tokens. The accuracy increases almost monotonically as more tokens are used to represent the video. The best performance is obtained by using all patches from T = 16 frames at 448×448 resolution, *i.e.* the most computationally expensive option. However, the object-centric representation recovers most of the performance with as few as 256 tokens (T = 16, N = 16), yielding 82.8% on K400 and 56.5% on SSv2, close to the top-line performance of 84.7% and 61.2%, respectively. Another benefit of the set-based representation is that increasing the spatial resolution from 224^2 to 448^2 , boosts K400 performance with no additional compute in the classifier.

4.2 EVALUATING MASK QUALITY

Section 4.1 evaluates the quality of set-structured representations S_N for downstream tasks. However, 398 the set structure does not imply that the representation is necessarily *object-centric*, *i.e.* aligned with 399 a human notion of objects and parts. This defining property is the focus of this section. In previous 400 works (Locatello et al., 2020; Greff et al., 2022; Hénaff et al., 2022; Seitzer et al., 2023; Kakogeorgiou 401 et al., 2024), it is evaluated as unsupervised segmentation: Object masks \mathcal{M}_N are compared to ground-402 truth object segments in annotated datasets, irrespective of the associated embedding. Though the 403 task is framed as segmentation, the goal is not to achieve pixel-perfect masks, but rather to show that 404 object vectors are associated to meaningful entities in an interpretable way, as in Figure 4. 405

Datasets. As in previous works, we consider a variety of datasets with different levels of complexity and annotation quality. The PASCAL VOC 2012 (Everingham et al., 2010), COCO 2017 (Lin et al., 2014), and ADE20K (Zhou et al., 2017) datasets are standard benchmarks for real-world object segmentation, with annotations for both instance-level and category-level segmentation. ClevrTex (Karazija et al., 2021) is a synthetic dataset of textured 3D shapes. MOVI (Greff et al., 2022) is a video dataset in which rendered 3D shapes move around in a textured environment. Following standard practice, each frame is treated as an independent sample and results are reported for the MOVI-C and MOVI-E variants, which differ in the complexity and number of objects in a scene.

422

384

386

387

388

389

390

391

392

393

394

395 396

Table 2: Unsupervised segmentation on following Hénaff et al. (2022): mean best overlap (mBO) and detection rate (DRate). To capture objects at different scales and granularity, our method produces 255 masks by running multiple k-means with $K \in \{1, 2, 4, ..., 128\}$ on a pre-trained DINOv1 ViT-B/14 and DINOv2 ViT-B/16. ODIN applies k-means similarly on a ViT-B/14 backbone trained from scratch with an object-centric loss. For SPOT, we obtain 255 masks from an ensemble of models based on DINOv1 ViT-B/14 and fine-tuned on COCO.

-T day I															
428			COCO 2017		VOC 2012			ClevrTex		ADE20K		MOVI C		MOVI E	
400		mBO^i	mBO^c	$DRate^{i}$	mBO^i	mBO^{c}	$DRate^{i}$	mBO^i	$DRate^{i}$	mBO^i	$DRate^{i}$	mBO^i	$DRate^{i}$	mBO^i	$DRate^{i}$
429	ODIN (ViT-B/14)	49.7	54.7	48.2	-	-	-	-	-	-	-	-	-	-	-
430	SPOT (DINOv1 ViT-B/14)	46.3	51.9	46.0	56.3	59.5	64.8	59.9	77.1	51.7	52.7	52.9	59.7	45.4	45.6
404	K-means (DINOv1 ViT-B/14)	50.2	58.9	51.1	61.3	65.3	69.8	74.9	92.0	58.7	64.4	66.4	80.3	57.3	64.6
431	K-means (DINOv2 ViT-B/16)	51.3	61.8	55.0	61.8	66.3	71.0	70.2	89.4	57.7	66.1	64.5	81.1	56.9	68.6

both instance-level and category-level annotations, we mark the metrics with a superscript i or c to distinguish which ground-truth masks they refer to.

Evaluation protocols. At a high level, both mBO and DetRate are recall-based metrics, *i.e.* they 435 measure the ability of a method to retrieve annotated ground-truth masks. For their nature, these 436 metrics are easily "cheated" by predicting more and more masks. Related works in the object-centric 437 literature typically follow one of two evaluation protocols. Total partitioning: Locatello et al. (2020) 438 and follow-up works constrain each pixel to belong to one and only one object, e.g. by taking an 439 arg max of the slot attention weights, effectively producing a total partitioning of the image. Under 440 this constraint, any overlap in the ground-truth masks can not be predicted, and therefore those 441 annotations are excluded from the metrics. Though appropriate for the synthetic images (Johnson 442 et al., 2017; Kabra et al., 2019) considered in early methods, we argue that this evaluation protocol is not aligned with the overarching goal of object-centric learning. Scoring a high mBO merely 443 indicates that the predictions match the annotation granularity of a dataset (e.g. bike vs. truck, wheel 444 vs. car) and that the number of predictions corresponds to the average number of annotations per 445 image. However, a model optimized for a specific dataset will hardly generalize outside its training 446 domain. Recall@N: the protocol of Hénaff et al. (2022) caps the number of predicted masks to a 447 large number, without enforcing a constraint on their spatial arrangement. This choice is better suited 448 to real-world scenes, in which three-dimensional objects can partially occlude each other, e.g. the 449 books on the table in the qualitative example of Figure 4. Furthermore, it allows to predict additional 450 masks for parts and sub-parts that may not be annotated in a test dataset, but are nonetheless relevant. 451 Overall, we find this second protocol to be more in line with the goal of object-centric learning: 452 obtaining representations that generalize well to a wide range of downstream tasks.

453 Results. For ease of comparison, we evaluate our method under both protocols. To match the 454 evaluation of Hénaff et al. (2022), we apply k-means multiple times with $K \in \{1, 2, 4, 8, \dots, 128\}$ 455 to extract a total of 255 masks per image. To compare with SPOT, we train eight SPOT models on 456 COCO with a number of slots in $\{2, 4, 8, \dots, 128\}$, and combine their independent predictions plus 457 a dummy mask. Note that it is highly impractical to train and run inference on multiple SPOT models 458 in parallel, nonetheless we provide the result for completeness. The results in Table 2 show that our 459 simple method significantly outperforms the more specialized baselines, especially for category-level segmentation. On the other hand, Table 3 follows the evaluation protocol of Locatello et al. (2020), 460 where each method predicts a small number of non-overlapping segments, and any overlap in the 461 ground-truth masks is ignored. We compare with the *state-of-the-art* methods LSD (Jiang et al., 462 2023), DINOSAUR (Seitzer et al., 2023), and SPOT (Kakogeorgiou et al., 2024), which fine-tune a 463 separate DINO ViT-B/16 with a fixed ad-hoc number of slots for each dataset. Instead, we run a single 464 k-means on a frozen DINOv2 ViT-S/14 backbone and choose K to match the number of slots of the 465 specialized methods on each dataset. Unsurprisingly, our method falls short of the state-of-the-art 466 because of the lack of dataset-specific tuning and the noisier masks produced by k-means. The 467 main takeaway is that our method is more suited for a high-recall no-tuning scenario where running 468 multiple k-means with different K values can capture objects/parts of different sizes and granularity. 469

470 4.3 Ablation of Image Encoder 471

The approach presented in this work is agnostic to the image representations used. We based our
analysis on DINOv2 but, in principle, similar object-centric representations could be obtained from
other pre-trained backbones. The ablation study in Table 4 compares classification and unsupervised

Table 3: Unsupervised segmentation following the protocol of Locatello et al. (2020): any overlap in the ground-truth masks is ignored for computing mean best overlap (mBO). The baseline methods fine-tune a separate model for each dataset with a number of slots fixed ad-hoc (in parentheses). Our method uses a frozen DINOv2 backbone, and run a single *k*-means with a matching number of clusters. As confirmed by the results, this setup is suboptimal for our method, which has no dataset-specific tuning and produces noisier masks.

481		COCO	2017 (7)	VOC 2	.012 (6)	MOVI-C (11)	MOVI-E (24)
482		mBO^i	mBO^{c}	mBO^i	mBO^c	mBO ⁱ	mBO^i
483	LSD	30.4	-	_	-	45.6	39.9
484	DINOSAUR	32.3	38.8	44.0	51.2	42.4	—
-0-	SPOT	35.0	44.7	48.3	55.6	47.3	40.1
485	K-means	30.4	38.8	44.0	50.1	41.8	36.0

486 **Table 4:** Ablation of the image encoder w.r.t. model family and size. For each backbone, we apply the same k-487 means clustering procedure and evaluate both classification and unsupervised segmentation tasks. Classification performance tends to increase with model size, while the trend for segmentation favors smaller models. 488

Backt	oone	Classificat	ion (Acc, mAP)	Segmentation (mBO ^{c} , mBO ^{i})		
Family	Size	SUN 397	COCO 2017	COCO 2017	MOVI-C	
MAE	ViT-B/16	60.8	69.6	42.7	48.1	
DINO	ViT-B/16	66.0	72.6	58.9	66.4	
DINOv2	ViT-B/14	77.7	82.3	61.8	64.5	
MAE	ViT-L/16	62.4	73.1	56.6	64.7	
CLIP	ViT-L/14	80.0	82.0	51.3	56.4	
DINOv2	ViT-L/14	79.2	84.0	60.1	63.6	
MAE	ViT-H/14	64.5	74.3	57.5	66.1	
AM-RADIO	ViT-H/16	80.9	83.8	61.9	63.5	
DINOv2	ViT-g/14	79.7	84.8	59.1	62.7	

498 499 500

501

504

505

506

507

segmentation performance of several models from other "families" including MAE (He et al., 2022), CLIP (Radford et al., 2021), DINO (Caron et al., 2021) and AM-RADIO (Ranzinger et al., 2024). Among base models, DINOv2 is a clear winner for its performance on SUN 397, COCO multi-label 502 classification and COCO category-level unsupervised segmentation. Moving to the large models, we observe that classification performance increases consistently with model size, while the trend for segmentation is less clear. CLIP shows strong classification performance, but its dense features are noisy and do not cluster well into object-centric segments. Surprisingly, among DINOv2 models, the smallest size works best for segmentation, likely due to clustering issues when the embeddings grow larger or due to pre-training artifacts. Therefore, for the large-size class, we retain the DINOv2 508 ViT-L/14 because of its strong performance on both task types compared to CLIP and MAE. At the 509 largest model size, AM-RADIO is the best performer, though it can not be compared directly to the other models as it is distilled from DINOv2, CLIP, and SAM. 510

511 512

513

5 CONCLUSION

514 We presented a clustering-based method for extracting object-centric representations of images using 515 any pre-trained SSL backbone. The strength of the approach is in its simplicity, speed and flexibility: 516 by applying k-means multiple times with different number of clusters, we are able to assemble a rich 517 yet compact set-structured representation of objects. Our findings show that the clustering approach 518 outperforms state-of-the-art object-centric learning methods in terms of representation quality and is a strong competitor in terms of unsupervised segmentation, though its masks are not pixel-perfect. 519

520 **Previous work.** Compared to the existing literature, our approach improves on several aspects: 521 a) it does not require lengthy re-training runs to swap the backbone or adjust the desired number 522 of objects, allowing for fast experimentation, b) it avoids learnable projections and preserves the 523 quality of the backbone embeddings, c) object semantics are learned during the pre-training of the 524 backbone on large-scale datasets, which avoids narrow domain biases, and d) at inference time, 525 the only computation overhead is the clustering algorithm, which is negligible w.r.t. the backbone. With this work, we also aim to refocus the research on object-centric *learning* to the evaluation 526 of its representations. The results in Section 4.1 shall serve as a baseline for more extensive and 527 comprehensive benchmarks. 528

529 Future directions. The outcome of this work stands on the figurative shoulders of powerful SSL 530 backbones. Any improvement in SSL pre-training that yields smoother, richer, and more expressive dense features will directly improve the clustering the resulting object-centric representations. As 531 more powerful SSL backbones are developed, discovering and representing objects from their dense 532 features using simple clustering methods may become the dominant approach. This work chooses 533 k-means for its simplicity and similarity to slot attention (Locatello et al., 2020), but other methods, 534 e.g. agglomerative clustering (Sibson, 1973) or DBSCAN (Ester et al., 1996), may be applicable 535 too. Last, we focused our evaluations on scene and video classification tasks, but the potential of 536 object-centric representations is broader. A small set of tokens that represent objects can alleviate the 537 computational cost of training vision-language models, and serve as a compact representation for 538 world models in reinforcement learning. This calls for standardized models and extended benchmarks, so that the community can iterate forward.

540 REFERENCES 541

542 543 544 545	Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat, Yann LeCun, and Nicolas Ballas. Self-Supervised Learning from Images with a Joint-Embedding Predictive Architecture. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and</i> <i>Pattern Recognition</i> , pp. 15619–15629, 2023.
546 547 548	Christopher P. Burgess, Loic Matthey, Nicholas Watters, Rishabh Kabra, Irina Higgins, Matt Botvinick, and Alexander Lerchner. MONet: Unsupervised Scene Decomposition and Representation. <i>arXiv:1901.11390 [cs, stat]</i> , January 2019.
550 551 552	Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging Properties in Self-Supervised Vision Transformers. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , October 2021.
553 554	Joao Carreira and Andrew Zisserman. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. <i>arXiv:1705.07750 [cs]</i> , February 2018.
555 556 557 558	Michael Chang, Thomas L. Griffiths, and Sergey Levine. Object Representations as Fixed Points: Training Iterative Refinement Algorithms with Implicit Differentiation. In <i>NeurIPS 2022</i> . arXiv, October 2022. doi: 10.48550/arXiv.2207.00787.
559 560 561	Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework for Contrastive Learning of Visual Representations. In <i>International Conference on Machine Learning</i> , pp. 1597–1607. PMLR, July 2020.
562 563 564	Minsu Cho, Suha Kwak, Cordelia Schmid, and Jean Ponce. Unsupervised object discovery and localization in the wild: Part-based matching with bottom-up region proposals. In <i>CVPR</i> , 2015.
565 566	Subhabrata Choudhury, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. Unsupervised Part Discovery from Contrastive Reconstruction. In <i>Neural Information Processing Systems</i> , 2021.
567 568 569	Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision Transformers Need Registers, September 2023.
570 571 572	Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In <i>IEEE Conference on Computer Vision and Pattern Recognition</i> , June 2009. doi: 10.1109/CVPR.2009.5206848.
573 574 575 576	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In <i>International Conference on Learning Representations</i> , September 2020.
577 578 579 580	Alaaeldin El-Nouby, Michal Klein, Shuangfei Zhai, Miguel Angel Bautista, Alexander Toshev, Vaishaal Shankar, Joshua M. Susskind, and Armand Joulin. Scalable Pre-training of Large Autoregressive Image Models, January 2024.
581 582 583	Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In <i>Knowledge Discovery and Data Mining</i> , pp. 226–231. AAAI Press, 1996.
584 585 586 587	Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The Pascal Visual Object Classes (VOC) Challenge. <i>International journal of computer vision</i> , 88(2): 303–338, 2010.
588 589 590	Ke Fan, Zechen Bai, Tianjun Xiao, Tong He, Max Horn, Yanwei Fu, Francesco Locatello, and Zheng Zhang. Adaptive Slot Attention: Object Discovery with Dynamic Slot Number. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 23062–23071, 2024.
591 592 593	Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised Representation Learning by Predicting Image Rotations. In <i>International Conference on Learning Representations</i> , February 2018.

594 595 596 597	Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The" something something" video database for learning and evaluating visual common sense. In <i>Proceedings of the IEEE International Conference on Computer Vision</i> , pp. 5842–5850, 2017.
598 599 600 601 602	Klaus Greff, Raphaël Lopez Kaufman, Rishabh Kabra, Nick Watters, Christopher Burgess, Daniel Zoran, Loic Matthey, Matthew Botvinick, and Alexander Lerchner. Multi-Object Representation Learning with Iterative Variational Inference. In <i>International Conference on Machine Learning</i> , pp. 2424–2433. PMLR, May 2019.
603 604	Klaus Greff, Sjoerd van Steenkiste, and Jürgen Schmidhuber. On the Binding Problem in Artificial Neural Networks. <i>arXiv:2012.05208 [cs]</i> , December 2020.
605 606 607 608 609 610 611 612	Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J Fleet, Dan Gnanapragasam, Florian Golemo, Charles Herrmann, Thomas Kipf, Abhijit Kundu, Dmitry Lagun, Issam Laradji, Hsueh-Ti (Derek) Liu, Henning Meyer, Yishu Miao, Derek Nowrouzezahrai, Cengiz Oztireli, Etienne Pot, Noha Radwan, Daniel Rebain, Sara Sabour, Mehdi S. M. Sajjadi, Matan Sela, Vincent Sitzmann, Austin Stone, Deqing Sun, Suhani Vora, Ziyu Wang, Tianhao Wu, Kwang Moo Yi, Fangcheng Zhong, and Andrea Tagliasacchi. Kubric: A scalable dataset generator. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , 2022.
613 614 615 616	Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised Learning. arXiv:2006.07733 [cs, stat], June 2020.
618 619 620	Mark Hamilton, Zhoutong Zhang, Bharath Hariharan, Noah Snavely, and William T. Freeman. Unsupervised Semantic Segmentation by Distilling Feature Correspondences. In <i>International</i> <i>Conference on Learning Representations</i> , September 2021.
621 622 623	Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked Autoencoders Are Scalable Vision Learners. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , 2022.
624 625 626 627	Olivier J. Hénaff, Skanda Koppula, Evan Shelhamer, Daniel Zoran, Andrew Jaegle, Andrew Zisserman, João Carreira, and Relja Arandjelović. Object discovery and representation networks. <i>arXiv:2203.08777 [cs]</i> , March 2022.
628 629	Geoffrey Hinton. How to represent part-whole hierarchies in a neural network. <i>arXiv:2102.12627</i> [cs], February 2021.
630 631 632	Geoffrey E. Hinton, Sara Sabour, and Nicholas Frosst. Matrix capsules with EM routing. In <i>International Conference on Learning Representations</i> , February 2018.
633 634	Jindong Jiang, Fei Deng, Gautam Singh, and Sungjin Ahn. Object-Centric Slot Diffusion. In <i>NeurIPS</i> . arXiv, July 2023. doi: 10.48550/arXiv.2303.10834.
635 636 637 638 639	Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross Girshick. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> , 2017.
640 641	Philip N. Johnson-Laird. Mental models and human reasoning. <i>Proceedings of the National Academy of Sciences</i> , 107(43):18243–18250, October 2010. doi: 10.1073/pnas.1012933107.
642 643 644	Rishabh Kabra, Chris Burgess, Loic Matthey, Raphael Lopez Kaufman, Klaus Greff, Malcolm Reynolds, and Alexander Lerchner. Multi-object datasets. DeepMind, 2019.
645 646 647	Ioannis Kakogeorgiou, Spyros Gidaris, Konstantinos Karantzalos, and Nikos Komodakis. SPOT: Self-Training with Patch-Order Permutation for Object-Centric Learning with Autoregressive Transformers., April 2024.

668

669

- Laurynas Karazija, Iro Laina, and Christian Rupprecht. ClevrTex: A Texture-Rich Benchmark for Unsupervised Multi-Object Segmentation. *arXiv:2111.10265 [cs]*, November 2021.
 Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *International Conference*
- 652 on Learning Representations, May 2014.
- Thomas Kipf, Gamaleldin F. Elsayed, Aravindh Mahendran, Austin Stone, Sara Sabour, Georg Heigold, Rico Jonschkowski, Alexey Dosovitskiy, and Klaus Greff. Conditional Object-Centric Learning from Video. arXiv:2111.12594 [cs, stat], November 2021.
- Alexander Kirillov. Segment Anything. In *Proceedings of the IEEE/CVF International Conference* on Computer Vision, 2023.
- Brenden M. Lake, Tomer D. Ullman, Joshua B. Tenenbaum, and Samuel J. Gershman. Building machines that learn and think like people. *Behavioral and Brain Sciences*, 40, 2017. ISSN 0140-525X, 1469-1825. doi: 10.1017/S0140525X16001837.
- Yann LeCun. A Path Towards Autonomous Machine Intelligence. Technical report, Meta FAIR, June 2022.
- 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 *European Conference on Computer Vision*, pp. 740–755, 2014.
 - S. Lloyd. Least squares quantization in PCM. *IEEE Transactions on Information Theory*, 28(2), March 1982. ISSN 1557-9654. doi: 10.1109/TIT.1982.1056489.
- Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold,
 Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-Centric Learning with Slot
 Attention. In *Advances in Neural Information Processing Systems*, October 2020.
- Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. Deep Spectral Methods:
 A Surprisingly Strong Baseline for Unsupervised Semantic Segmentation and Localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8364–8375, 2022.
- Mehdi Noroozi and Paolo Favaro. Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling (eds.), *Computer Vision ECCV 2016*, Lecture Notes in Computer Science, pp. 69–84, Cham, 2016. Springer International Publishing. ISBN 978-3-319-46466-4. doi: 10.1007/978-3-319-46466-4_5.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov,
 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas
 Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael
 Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut,
 Armand Joulin, and Piotr Bojanowski. DINOv2: Learning Robust Visual Features without
 Supervision, April 2023.
- 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 *Proceedings of the* 38th International Conference on Machine Learning, pp. 8748–8763. PMLR, July 2021.
- Sai Saketh Rambhatla, Ishan Misra, Rama Chellappa, and Abhinav Shrivastava. MOST: Multiple
 Object localization with Self-supervised Transformers for object discovery. In *ICLR*, August 2023. doi: 10.48550/arXiv.2304.05387.
- Mike Ranzinger, Greg Heinrich, Jan Kautz, and Pavlo Molchanov. AM-RADIO: Agglomerative
 Vision Foundation Model Reduce All Domains Into One, April 2024.
- Bryan Russell, William Freeman, Alexei Efros, Josef Sivic, and Andrew Zisserman. Using multiple segmentations to discover objects and their extent in image collections. In *CVPR*, 2006.

702	Mehdi S. M. Saijadi, Daniel Duckworth, Aravindh Mahendran, Sjoerd van Steenkiste, Filin Pavetić
700	Wend S. W. Sujjadi, Damer Dackworth, Thavinan Manendran, Sjoerd van Steenkiste, Tinp Tavete,
703	Mario Lučić, Leonidas J. Guibas, Klaus Greff, and Thomas Kipf. Object scene representation
704	transformer NeurIPS 2022
	ualistoffici. Wewill 5, 2022.
705	

- Maximilian Seitzer, Max Horn, Andrii Zadaianchuk, Dominik Zietlow, Tianjun Xiao, Carl-Johann Simon-Gabriel, Tong He, Zheng Zhang, Bernhard Schölkopf, Thomas Brox, and Francesco Locatello. Bridging the Gap to Real-World Object-Centric Learning. In *The Eleventh International Conference on Learning Representations*, February 2023.
- R. Sibson. SLINK: An optimally efficient algorithm for the single-link cluster method. *The Computer Journal*, 16(1):30–34, January 1973.
- Oriane Siméoni, Gilles Puy, Huy V. Vo, Simon Roburin, Spyros Gidaris, Andrei Bursuc, Patrick Pérez, Renaud Marlet, and Jean Ponce. Localizing Objects with Self-Supervised Transformers and no Labels. *arXiv preprint arXiv:2109.14279*, 2021.
- Gautam Singh, Fei Deng, and Sungjin Ahn. Illiterate DALL-E Learns to Compose. In *International Conference on Learning Representations*, March 2022a.
- Gautam Singh, Yi-Fu Wu, and Sungjin Ahn. Simple Unsupervised Object-Centric Learning for Complex and Naturalistic Videos. In *NeurIPS*. arXiv, May 2022b. doi: 10.48550/arXiv.2205.14065.
- Josef Sivic, Bryan Russell, Alexei Efros, Andrew Zisserman, and William Freeman. Discovering
 objects and their location in images. In *ICCV*, 2005.
- Elizabeth S. Spelke and Katherine D. Kinzler. Core knowledge. *Developmental Science*, 10(1):89–96, January 2007. ISSN 1363-755X. doi: 10.1111/j.1467-7687.2007.00569.x.
- Karl Stelzner, Kristian Kersting, and Adam R. Kosiorek. Decomposing 3D Scenes into Objects via
 Unsupervised Volume Segmentation, April 2021.
- Chen Sun, Calvin Luo, Xingyi Zhou, Anurag Arnab, and Cordelia Schmid. Does visual pretraining help end-to-end reasoning? In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 21524–21540. Curran Associates, Inc., 2023.
- Jiayu Tang and Paul H Lewis. Non-negative matrix factorisation for object class discovery and image auto-annotation. In *CIVR*, 2008.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
 - Wouter Van Gansbeke, Simon Vandenhende, and Luc Van Gool. Discovering Object Masks with Transformers for Unsupervised Semantic Segmentation, June 2022.
- Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th International Conference on Machine Learning*, ICML '08, pp. 1096–1103, New York, NY, USA, July 2008. Association for Computing Machinery. ISBN 978-1-60558-205-4. doi: 10.1145/1390156.1390294.
 - Huy V. Vo, Francis Bach, Minsu Cho, Kai Han, Yann LeCun, Patrick Pérez, and Jean Ponce. Unsupervised image matching and object discovery as optimization. In *CVPR*, 2019.
- Huy V. Vo, Patrick Pérez, and Jean Ponce. Toward Unsupervised, Multi-object Discovery in LargeScale Image Collections. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael
 Frahm (eds.), *Computer Vision ECCV 2020*, pp. 779–795, Cham, 2020. Springer International
 Publishing. ISBN 978-3-030-58592-1. doi: 10.1007/978-3-030-58592-1_46.
- Huy V. Vo, Elena Sizikova, Cordelia Schmid, Patrick Pérez, and Jean Ponce. Large-Scale Unsupervised Object Discovery. In *Advances in Neural Information Processing Systems*, volume 34, pp. 16764–16778. Curran Associates, Inc., 2021.

738

739

744

745

- Huy V. Vo, Vasil Khalidov, Timothée Darcet, Théo Moutakanni, Nikita Smetanin, Marc Szafraniec, Hugo Touvron, Camille Couprie, Maxime Oquab, Armand Joulin, Herve Jegou, Patrick Labatut, and Piotr Bojanowski. Automatic Data Curation for Self-Supervised Learning: A Clustering-Based Approach. *Transactions on Machine Learning Research*, May 2024. ISSN 2835-8856.
- Xinlong Wang, Zhiding Yu, Shalini De Mello, Jan Kautz, Anima Anandkumar, Chunhua Shen, and Jose M. Alvarez. FreeSOLO: Learning To Segment Objects Without Annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14176–14186, 2022.
- Xudong Wang, Rohit Girdhar, Stella X. Yu, and Ishan Misra. Cut and learn for unsupervised object
 detection and instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3124–3134, 2023a.
- Yangtao Wang, Xi Shen, Yuan Yuan, Yuming Du, Maomao Li, Shell Xu Hu, James L. Crowley, and Dominique Vaufreydaz. TokenCut: Segmenting Objects in Images and Videos With Self-Supervised Transformer and Normalized Cut. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(12):15790–15801, December 2023b. ISSN 0162-8828, 2160-9292, 1939-3539. doi: 10.1109/TPAMI.2023.3305122.
- Markus Weber, Max Welling, and Pietro Perona. Towards automatic discovery of object categories. In *CVPR*, 2000.
- Alfred North Whitehead. Symbolism: Its Meaning and Effect. Fordham University Press, 1985.
 ISBN 978-0-8232-1138-8.
- Ziyi Wu, Nikita Dvornik, Klaus Greff, Thomas Kipf, and Animesh Garg. SlotFormer: Unsupervised Visual Dynamics Simulation with Object-Centric Models. In *The Eleventh International Conference on Learning Representations*, September 2022.
- Ziyi Wu, Jingyu Hu, Wuyue Lu, Igor Gilitschenski, and Animesh Garg. SlotDiffusion: Object-Centric Generative Modeling with Diffusion Models. In *Thirty-Seventh Conference on Neural Information Processing Systems*, November 2023.
- Jianxiong Xiao, James Hays, Krista A. Ehinger, Aude Oliva, and Antonio Torralba. SUN database:
 Large-scale scene recognition from abbey to zoo. In 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pp. 3485–3492, June 2010. doi: 10.1109/CVPR.2010.
 5539970.
- Andrii Zadaianchuk, Matthaeus Kleindessner, Yi Zhu, Francesco Locatello, and Thomas Brox.
 Unsupervised Semantic Segmentation with Self-supervised Object-centric Representations. In *The Eleventh International Conference on Learning Representations*, September 2022.
- Quanshi Zhang, Ying Nian Wu, and Song-Chun Zhu. Mining and-or graphs for graph matching and object discovery. In *ICCV*, 2015.
- Richard Zhang, Phillip Isola, and Alexei A. Efros. Colorful Image Colorization. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling (eds.), *Computer Vision ECCV 2016*, Lecture Notes in Computer Science, pp. 649–666, Cham, 2016. Springer International Publishing. ISBN 978-3-319-46487-9. doi: 10.1007/978-3-319-46487-9_40.
- Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. Learning Deep
 Features for Scene Recognition using Places Database. In *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene
 parsing through ADE20K dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. iBOT: Image BERT Pre-Training with Online Tokenizer. In *ICLR*. arXiv, January 2022. doi: 10.48550/ arXiv.2111.07832.
- ⁸⁰⁹ Jun-Yan Zhu, Jiajun Wu, Yan Xu, Eric Chang, and Zhuowen Tu. Unsupervised object class discovery via saliency-guided multiple class learning. In *CVPR*, 2012.

810 A ADDITIONAL EXPERIMENTS

A.1 ABLATION STUDY: WHICH FEATURES TO CLUSTER FOR UNSUPERVISED SEGMENTATION?

For all experiments in the main text, we use the output features of each model as the input to the clustering algorithm. In the case of DINOv1 and DINOv2, these are the final patch features, after the last transformer block (self attention and MLP), and after the final LayerNorm operation. Other works in the unsupervised segmentation literature have proposed to use the queries, keys, and values of the last attention block instead (Siméoni et al., 2021; Wang et al., 2022; 2023b;a). In Table 5, we compare the performance of clustering-based unsupervised segmentation when using the queries, keys, and values of the last attention block, or the final patch features at the output of the model. These results extend the ones presented in Table 2 in Section 4.2. In the setting considered in this work, the patch features prove to be the best choice for clustering-based unsupervised segmentation, which we attribute to the clustering-like effect of the IBOT loss applied to patch features during DINOv2 pre-training (Oquab et al., 2023).

Table 5: Ablation study: which DINOv2 ViT-B/16 (Darcet et al., 2023) features work best for clustering-based unsupervised segmentation? We compare taking the queries, keys, and values of the last attention block, or taking the final patch features at the output of the model. We follow the evaluation protocol of Hénaff et al. (2022), namely clustering the spatial features multiple times with *k*-means using $k \in \{1, 2, 4, ..., 128\}$ to obtain 255 overlapping masks, and repording mean best overlap (mBO) and detection rate (DRate).

	COCO 2017		VOC 2012		ClevrTex		AD	E20K	MOVI C		MOVI E			
	mBO^i	mBO^{c}	DRate^i	mBO^i	mBO^{c}	DRate^i	mBO^i	DRate^i	mBO^i	DRate^i	mBO^i	DRate^i	mBO^i	$DRate^{i}$
Query	33.8	39.1	22.7	43.7	46.5	35.6	52.5	53.5	44.8	38.7	50.2	47.5	41.0	31.3
Key	30.7	36.0	19.9	41.8	45.1	34.5	50.8	50.8	42.2	36.1	48.7	45.1	38.7	28.4
Value	26.5	32.1	16.2	37.6	41.8	29.4	39.4	33.5	35.9	27.6	38.8	28.9	27.5	13.5
Output	51.3	61.8	55.0	61.8	66.3	71.0	70.2	89.4	57.7	66.1	64.5	81.1	56.9	68.6

A.2 ABLATION STUDY: VIDEO ACTION RECOGNITION ON A BUDGET

In Section 4.1.2, we demonstrate the effectiveness of using k-means centroids to represent video frames for action recognition. The hypothesis is that k-means clustering is able to isolate and represent objects, offering a compact and structured representation of each frame. We demonstrate that given a small budget of 8 or 16 tokens per frame, the k-means centroids are able to approximate the performance of a more expensive model that takes all patches of each frame as input. Given this result, one may wonder if other choices of representation can yield the same performance, while satisfying the budget constraint. In other words, we wish to investigate whether it is only a matter of token count or whether the content of these tokens is also important. We set up an additional experiment using videos at 224×224 pixel resolution, where we consider:

- dropping patches at random to fit into a given budget, *e.g.* keeping 16 patches out of 256;
- average-pooling the spatial features to fit into a given budget, *e.g.* pooling 16×16 patches to 4×4 with bicubic interpolation.

In Figure 5, we compare the performance of these different choices of representation for video action recognition on Kinetics-400 and Something-Something-V2. For both datasets and all model sizes considered, ViT-S/16, ViT-B/16, ViT-L/16, we observe that k-means centroids consistently yield better performance than other choices of representation at a given budget. We conclude that the object-centric representation provided by k-means centroids is beneficial for video action recognition, and that the performance gains can not be solely attributed to the number of tokens used.



Under review as a conference paper at ICLR 2025

Figure 5: Ablation study: video action recognition on Kinetics-400 and Something-Something-V2 with a budget constraint on the number of tokens. Starting from 8 frames per video at 224×224 resolution, we compare the performance of taking k-means centroids as object representations, dropping patches at random, or average-pooling the spatial features on a square grid, *e.g.* 1×1 , 2×2 , 3×3 , ..., 16×16 . As a baseline, we also include the performance of a model that represents each frame only with the CLS token.

918 A.3 TRAINING SPOT WITH A DINOV2 BACKBONE

For the experiments in Section 4, we train several SPOT models using the author's code (Kakogeorgiou et al., 2024). To confirm that our reproduction matches the official results, we train a DINOv1 ViT-B/16 model on COCO with 7 slots, as done in the original work. The unsupervised segmentation results on COCO reported in Table 6 are consistent with the official results.

Table 6: Unsupervised segmentation results on COCO for a SPOT model with 7 slots fine-tuned on COCO using a DINOv1 ViT-B/16 backbone.

	FG-ARI	mBO^{c}	mBO^{i}
Official	37.8	44.7	35.0
Reproduction	36.2	45.0	35.0

We then attempted to train a SPOT model with more recent DINOv2 ViT-B/16 backbones of different sizes. For all models, we use COCO as the training set, but train different models with 7, 14, and 28 slots to be tested on COCO, VOC 2012, MOVI C, and MOVI E. Each model takes around 24 hours to train, both stage 1 and stage 2, on a single GPU. In this case, however, we observe a significant performance drop when using DINOv2 as the backbone, therefore we do not use these models for the experiments in the main text.

Table 7: SPOT models with 7, 14, and 28 slots based on either a DINOv1 or DINOv2 backbone. All models are fine-tuned on COCO and evaluated on the other datasets.

	COC	CO 2017	(7)	VO	C 2012	(6)	MOVI-	C (11)	MOVI-	E (24)
	FG-ARI	mBO ^c	mBO^i	FG-ARI	mBO ^c	mBO^i	FG-ARI	mBO ⁱ	FG-ARI	mBO
DINOv1 ViT-B/16	36.2	45.0	35.0	22.0	58.4	50.3	64.2	38.8	60.6	29.5
DINOv2 ViT-S/14	41.6	35.1	29.3	27.6	43.8	41.0	63.3	26.1	60.1	21.4
DINOv2 ViT-B/14	45.1	33.8	28.9	36.6	44.2	41.8	59.8	22.4	53.9	18.3
DINOv2 ViT-L/14	39.6	32.7	26.6	39.7	44.2	38.8	56.9	18.9	56.0	16.5