

Object-Aware DINO (Oh-A-Dino): Enhancing Self-Supervised Representations for Multi-Object Instance Retrieval

Stefan Sylvius Wagner¹ Stefan Harmeling²

Abstract

Object-centric learning is fundamental to human vision and crucial for models requiring complex reasoning. Traditional approaches rely on slot-based bottlenecks to learn object properties explicitly, while recent self-supervised vision models like DINO have shown emergent object understanding. However, DINO representations primarily capture global scene features, often confounding individual object attributes. We investigate the effectiveness of DINO representations and slot-based methods for multi-object instance retrieval. Our findings reveal that DINO representations excel at capturing global object attributes such as object shape and size but struggle with object-level details like color, whereas slot-based representations struggle at both global and object-level understanding. To address this, we propose a method that combines global and local features by augmenting DINO representations with object-centric latent vectors from a Variational Autoencoder trained on segmented image patches that are extracted from the DINO features. This approach improves multi-object instance retrieval performance, bridging the gap between global scene understanding and fine-grained object representation without requiring full model retraining.

1. Introduction

Object-centric learning is at the core of human vision and task understanding making it a critical problem in computer vision and key to developing models with complex reasoning (van Steenkiste et al., 2019; Schölkopf et al., 2021). Finding the right approach to learning object-centric representations is challenging since object properties need to be disentangled in a way that allows for generalisation.

Current methods for object-centric learning focus on disentangling and learning individual object features, however disentangling a scene into multiple factors of variation may hinder general scene understanding (Montero et al., 2024). An unexplored area in object-centric learning is the use of pre-trained representations from self-supervised models. Recently, large pre-trained vision models have been touted as vision foundation models due to the strong generalisation capabilities their learnt representations provide in various downstream computer vision tasks (Uelwer et al., 2023). For instance, it has been shown that certain self-supervised models like DINO (Caron et al., 2021; Oquab et al., 2023) learn emerging object properties thanks to the self-supervised learning objective. Thus using these models and their representations for object-centric learning is appealing since more common object-centric methods learn object properties via a slot-based bottleneck (Locatello et al., 2020) which generally requires retraining the complete model for different tasks. Although self-supervised models are effective on a broad range of vision tasks, their effectiveness on tasks that require object-centric understanding is not well understood. Recent work has shown that the self-supervised learning objective used to train these models can confound individual features in the learnt latent space (Bizeul et al., 2024).

In this work, we study the effectiveness of current pre-trained and object-centric representations on multi-object scenes by performing instance retrieval on object-centric datasets. Our work shows that the task of multi-object instance retrieval is hard, since its performance relies heavily on the quality of the representations that are used for retrieval. This is in contrast to traditional segmentation or classification tasks where an additional component is usually learnt for the downstream task. The datasets contain images of objects with complex attributes that require representations with both a specific and holistic understanding of the objects.

We study both DINO and slot-attention learnt representations where our work reveals the following: we find that slot-attention representations generally struggle with multi-object instance retrieval, faring weaker in predicting different object attributes compared to DINO representations.

¹Department of Computer Science, Heinrich Heine University Düsseldorf ²Department of Computer Science, Technical University of Dortmund. Correspondence to: Stefan Wagner <stefan.wagner@hhu.de>.

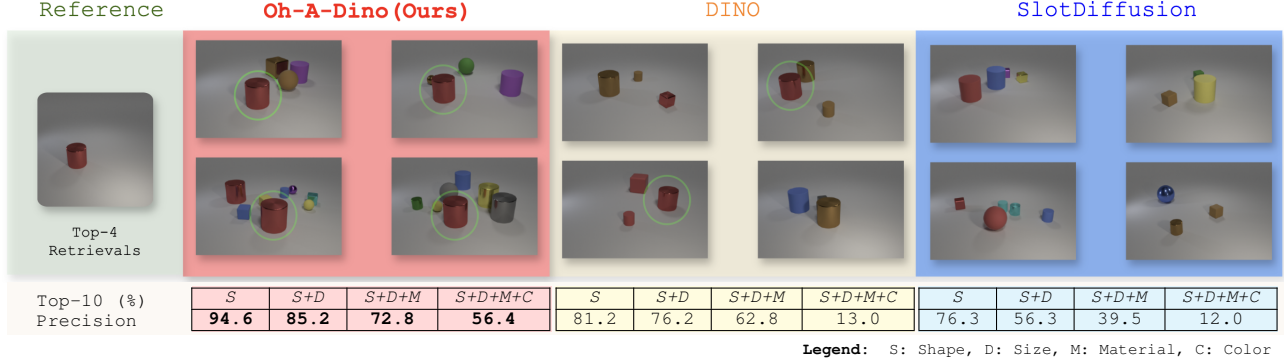


Figure 1: **DINO representations and slot-based representations struggle at multi-object instance retrieval, especially with fine-grained object features, such as material and color. Our method is able to retrieve more relevant images by combining global and local representations.** DINO representations excel at retrieving images that match global object features such as shape and size, while retrieval degrades for more fine-grained object-level features such as material and color. Slot-based representations lack specialisation, where retrievals sometimes match fine-grained features such as colour but often failing to retrieve similar object altogether. Our method performs the best being able to balance global and local object-level features.

While DINO’s representations perform better, they fail to retrieve color consistently. An issue that has been observed in previous work (Bizeul et al., 2024; Garrido et al., 2025).

To address this, we propose a method that enhances the DINO representations with object-centric latent vectors learnt with a Variation Autoencoder (VAE). We augment the DINO representations by concatenating latent vectors learnt from segmented image patches, where the segmentation masks are extracted directly from the DINO features via PCA. We observe that segmenting the image into object-level features improves DINO, while concatenating the latent representations from the VAE enforces the learning of object-level representations, which are dynamic and learnt implicitly through the VAE objective.

We summarise our contributions:

1. We identify the limitations of pre-trained DINO representations and slot-based representations in multi-object instance retrieval tasks. We reveal that that DINO representations are more effective at multi-object instance retrieval.
2. We propose a method to enhance pre-trained DINO representations with object-centric latent vectors learnt with a Variational Autoencoder. This approach separates the object-level understanding from general scene understanding, allowing for a better object-centric representation.
3. We show that our method improves multi-object instance retrieval for DINO representations by improving its poor ability to identify color. This advances

the understanding of how to effectively leverage pre-trained representations for complex downstream tasks that require multi-object understanding.

2. Background: Object-Centric Learning and Multi-Object Instance Retrieval

Slot-based object centric learning. Object-centric learning has emerged as a critical problem in computer vision, particularly in tasks requiring an understanding of complex visual scenes composed of multiple interacting objects. Recent advances in object-centric representation learning (e.g., SlotAttention (Locatello et al., 2020), MONet (Burgess et al., 2019), SlotDiffusion (Wu et al., 2023), Slot-VAE (Wang et al., 2023)) have enabled models to decompose images into object-level features, making it possible to reason about individual objects within a scene. These methods mostly focus on learning disentangled representations of individual objects by enforcing an object-level slot-based bottleneck, which then allows for faithful image reconstruction, generation of novel object configurations, or object segmentation. Slot-VAE similar to our work, uses a VAE to learn both object-level and global-level features. The main difference, however, is that Slot-VAE learns slot representations jointly for both the global and object-level features.

In order to extract the objects, the image is usually tokenised and an attention mechanism is applied to the tokens (e.g. SLATE (Singh et al., 2021)). A slot decoder then reconstructs the original image from the slot representation. STEVE (Singh et al., 2022) reconstructs the image from the slot-representation autoregressively to learn dependencies between different image patches i.e. tokens. More recent

methods, have further exploited this by attempting to incorporate global scene level features into the representations (Chen et al., 2024) or by adding compositional reasoning into the slot-based representations (Jung et al., 2024).

Leveraging self-supervised pre-trained features. Self-supervised learning has emerged as a powerful paradigm for training vision models, where especially vision transformer models like DINO (Caron et al., 2021; Oquab et al., 2023) have shown strong generalization capabilities across a wide range of downstream tasks. The appeal of such pre-trained models is the ability to leverage the representations for various tasks without the need for extensive fine-tuning or retraining. Recent methods such as SLATE (Singh et al., 2021) and STEVE (Singh et al., 2022) have leveraged vision transformers for object-centric learning to learn improved slot representations that also capture global scene information via attention mechanisms. These methods have shown improved performance on object-centric tasks, but they require retraining the complete model, which can be computationally expensive and time-consuming. Similarly, DINOSAUR (Seitzer et al., 2023) and GOLD (Chen et al., 2024) additionally use DINO representations to provide better priors to learn object-centric representations, which allows for learning better slot-representations on natural images. In our work we leverage the representations directly to extract object-level features, which allows us to forego the slot-bottleneck by leveraging their inherent general scene understanding. We can thus segment the image into object-level features, which learn additional object-level features that are not present in the pre-trained model.

Assessing the quality of representations with instance retrieval. Instance retrieval is a challenging task that requires models to retrieve images with similar object attributes and configurations. Traditionally, this task has been best solved by methods that leverage representations from pre-trained models, such as ResNet (He et al., 2016) or vision transformers (Dosovitskiy et al., 2020). Generally, instance retrieval depends on learning a good latent space or representation of the image, where retrievals are ranked on a similarity measure (Chen et al., 2021). Naturally, instance retrieval has been used to evaluate the quality of representations in self-supervised models (Caron et al., 2021; Oquab et al., 2023). Since we are interested in evaluating the quality of the representations for multi-object understanding, we choose the task of multi-object instance retrieval.

Improving multi-object instance retrieval is particularly important for tasks like visual question answering, content-based image search, and scene understanding, where identifying complex object interactions is essential (Chen et al., 2021). Beyond traditional vision applications, multi-object instance retrieval plays a crucial role in goal-conditioned reinforcement learning (GCRL) and robotic tasks (Zheng

et al., 2024). In GCRL, agents often rely on goal images to specify desired configurations of objects in the environment, particularly in tasks requiring manipulation or navigation in multi-object settings (Le Paine et al., 2019). Similar challenges are present in hierarchical reinforcement learning (HRL), where sub-goals often involve interacting with specific object sets, robust multi-object retrieval can inform option selection and improve learning efficiency (Hafner et al., 2022).

3. Analysing Representation Quality with Object-Centric Instance Retrieval

A first step towards understanding the quality of object-centric representations is to assess the limitations of current methods. We compare the performance of DINO (Caron et al., 2021) and SlotDiffusion (Wu et al., 2023) representations in multi-object instance retrieval with regards to specific object attributes such as (shape, size, material and color). Specifically we want to understand how well combinations of these attributes can be retrieved. DINO is known to capture global scene information well, however it is unclear whether it can consistently capture fine-grained object-level features (Bizeul et al., 2024). Similarly, SlotDiffusion is a state-of-the-art slot-based object-centric method that learns slot representations to capture object-level features, however it is not exactly clear how well fixed slot representations can capture fine-grained object-level features and global features (Montero et al., 2024).

We perform a brief analysis in Figure 1 on the top-10 retrievals in order to understand the limitations of these representations in multi-object instance retrieval. On the top row we show the top-4 retrievals for a single object reference image, while on the bottom row we show the top-10 retrievals over the whole dataset for a combination of four different object attributes. We present our findings below.

Slot-based representations struggle with both global and object-level features. SlotDiffusion, struggles with both global and object-level features. In the top row, the top-4 retrievals are quite different from the query image, indicating that the slot-based representations are not able to capture the object-level features well. In the bottom row, analysing the top-10 retrievals over the whole dataset shows a steep decline in performance after a second attribute is added to the query image. We believe the fixed slot representations have most likely not separated the global and object-level properly, resulting in poor retrieval performance when more attributes are added to the query image.

DINO representations struggle with fine-grained object features. In the top row we see that DINO is able to retrieve the query image two out of four times. However, the DINO representations seem to be invariant to color, since the top

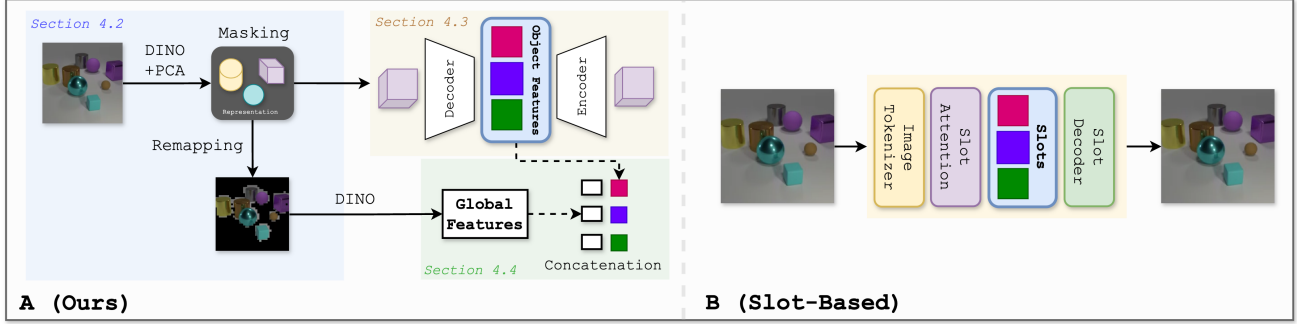


Figure 2: **Our method leverages the implicit global object understanding of self-supervised models such as DINO to extract latent features for multi-object instance retrieval.** **A:** Instead of learning slot-representations like traditional object-centric methods which try to compress global and object level features into a single latent space, we propose to learn a latent space over global and object-level features separately, which are then concatenated to form an improved latent representation for multi-object instance retrieval. **B:** Traditional slot-based object-centric methods learn slot-representations that compress global and object-level features into a single latent space. This can be understood as learning a latent space over both global and object-level features simultaneously. However, this can result in poorer representations since the bottleneck has to account for both global and object-level features.

retrieval contains the correct object, however with the wrong color. When looking at the top-10 retrievals over the whole dataset, we see that DINO is able to retrieve images concerning more global features like shape and size, however it struggles with capturing more fine-grained object features like color and material.

These findings inspire us to enhance the already strong performing DINO representations. Essentially, we propose a method that realigns the DINO representations to focus more on object-level features, while still retaining the global scene information. We show in the following how we achieve this by learning object-level features with a Variational Autoencoder (VAE) and combining them with the DINO representations.

4. Enhancing DINO with Self-Supervised Object-Centric Representations

Enhancing DINO representations with self-supervised object attributes is a multi-step process which is described in Figure 2. In Section 4.2, we first leverage the global features of the DINO representations to segment the image around important objects. For this we apply PCA on the DINO features of t consecutive images and use the first principal component to extract the object-level features. This allows us to extract a segmented image that is then split up into image patches. In Section 4.3, we learn a latent space over these image patches using a Variational Autoencoder. The VAE is trained to reconstruct the segmented image patches to learn a latent space at the object level. This latent space can be seen to capture the object-level attributes of the image, which we then seek to combine with the global features

of the DINO representations. In Section 4.4, we then concatenate the global features of the DINO representations with the latent vectors learnt from the VAE to create a new and more informative representation.

4.1. Preliminaries

Given an image $x \in \mathbb{R}^{h \times w \times c}$, the first step is to retrieve the DINO representations from a pre-trained DINO model, which we define as a function $g : \mathbb{R}^{h \times w \times c} \rightarrow \mathbb{R}^{p_h \times p_w \times n_y}$, where $y \in \mathbb{R}^{p_h \times p_w \times n_y}$ are the $p_h \times p_w$ resulting patch embeddings each with an embedding dimension of n_y . The number of patches $p_h \cdot p_w$ depends on the chosen patch size $s_h \times s_w$. In the following, we either view the output of g as $y \in \mathbb{R}^{p_h \times p_w \times n_y}$ or as $y \in \mathbb{R}^{p \times n_y}$ where $p = p_h \cdot p_w$ depending on whether a list of patches is sufficient.

The DINO embeddings are quite high-dimensional. Curiously, the most important PCA components capture the most relevant object-level features (see supplementary material A.2). For our purposes, three components are sufficient to obtain a low-dimensional representation that corresponds to the image space. For stability, we compute the PCA components on a batch of t different images,

$$X_{1:t} = [x_1, x_2, \dots, x_t] \in \mathbb{R}^{t \times h \times w \times c}, \quad (1)$$

where $x_i \in \mathbb{R}^{h \times w \times c}$ is the i -th image in the batch. We sample $X_{1:t-1}$ from the dataset and append the image of interest x as the last image x_t we want to extract the object-level features from. In our experiments, we set $t = 50$.

Finally, for the t images we have

$$Y_{1:t} = g(X_{1:t}), \quad (2)$$

where $Y_{1:t} \in \mathbb{R}^{t \times p \times n_y}$ are the patch embeddings of the t images. We omit the indices and write Y where appropriate. In the following, we describe how we use PCA to extract the object-level features from y_t , i.e., the DINO patch embeddings of a given image.

4.2. Extracting Object-Level Features with PCA

Extracting object-level features from the DINO patch embeddings consists of three steps: (i) Separating background and foreground by applying PCA to the DINO patch embeddings and thresholding the first component to obtain a segmentation mask. (ii) Reapplying PCA to the foreground patch embeddings to obtain a segmentation mask that is consistent across objects. (iii) Remapping the segmentation mask to the original image size and splitting the image into p image patches. All steps are described in detail below.

(i) Separating foreground and background. Since we want to apply PCA to the collected patches of t images, we flatten Y over the first two dimensions and get

$$\hat{Y} \in \mathbb{R}^{(t \cdot p) \times n_y}. \quad (3)$$

Then PCA applied to \hat{Y} provides the first principal component

$$Z' = \text{PCA}(\hat{Y}) \in \mathbb{R}^{(t \cdot p)}, \quad (4)$$

of the DINO patch embeddings.

Interestingly, the first principal component carries the information what pixels belong to foreground objects. This insight can be exploited to separate foreground objects from the background: to obtain the segmentation mask we threshold the first principal component Z' against the median of its components (i.e., $\text{median}(Z') \in \mathbb{R}$). More precisely, we define the segmentation mask to extract the foreground (abbr. as “fg”) as follows

$$\hat{M}^{\text{fg}} = \mathbb{1}(Z' > \text{median}(Z')) \in \{0, 1\}^{(t \cdot p)} \quad (5)$$

where the function $\mathbb{1}$ translates a boolean vector into a zero/one vector. In the following, we use the segmentation mask \hat{M}^{fg} to extract the foreground objects and perform a second PCA to obtain a consistent segmentation mask across objects.

(ii) Reapplying PCA to the foreground patch embeddings. While the previous step already delivers a decent segmentation mask, we found that reapplying PCA to the foreground patch embeddings delivers a more consistent segmentation mask across objects, i.e., one that does not cut across objects. We thus define the foreground patch embeddings as

$$Y^{\text{fg}} = \hat{M}^{\text{fg}} \odot \hat{Y} \in \mathbb{R}^{(t \cdot p) \times n_y}, \quad (6)$$

where the \odot is the component-wise multiplication that is broadcasted according to the size of \hat{Y} . The result Y^{fg} contains the foreground patch embeddings with the background patches zeroed out. We then reapply Equations (4) and (5) to Y^{fg} and obtain a new segmentation mask $M^{\text{fg}} \in \{0, 1\}^{(t \cdot p)}$ that improves \hat{M}^{fg} by being more faithful with respect to objects.

(iii) Remapping the masked patch embeddings to the original image to get image patches that contain objects.

One of the benefits of using a vision transformer architecture is that the patch embeddings can be remapped to the original image. We thus remap the segmentation mask M^{fg} to the original image to extract a segmented image to calculate global features and to extract image patches to learn object-level features with a VAE.

Initially, we reshape the masking vector M^{fg} to a $t \times p$ matrix. Since we are only interested in the last image x_t (see Section 4.1), we choose the last row of M^{fg} which we call $m \in \mathbb{R}^p$. This vector m is the segmentation mask of the last image $x_t \in \mathbb{R}^{h \times w \times c}$ which was our initial image x . More precisely, the mask determines what patches are relevant for the objects, i.e., an entry in m tells us whether a particular patch is masked or not. To obtain a masked image, we expand m to the size of x , by creating patches with zeros or ones, respectively, for each mask position. We thus reverse the operations in Section 4.1 and upscale the segmentation mask to the original image size. More precisely, we insert patches of size $s_h \times s_w$ at each patch position $p_h \times p_w$ to upscale the mask to the original image size. We define a function g' for this purpose

$$g' : \mathbb{R}^p \rightarrow \mathbb{R}^{(p_h \times p_w) \times (s_h \times s_w)}, \quad (7)$$

where $p = p_h \cdot p_w$. Applying g' we get

$$m^o = g'(m), \quad (8)$$

where $m^o \in \mathbb{R}^{(p_h \times p_w) \times (s_h \times s_w)}$ is the remapped segmentation mask. In other words, each entry in m^o is a patch of size $s_h \times s_w$ filled either with ones or with zeroes and all patches correspond to non-overlapping regions in the original image.

We can now obtain the segmented image and the image patches by applying element-wise multiplication between the original image splitted into patches $x \in \mathbb{R}^{(p_h \times p_w) \times (s_h \times s_w) \times c}$ and the segmentation mask m^o

$$x^o = x \odot m^o. \quad (9)$$

(again broadcasting over the extra dimension of the color channels) where $x^o \in \mathbb{R}^{(p_h \times p_w) \times (s_h \times s_w) \times c}$ is the segmented image. Finally, x^o can easily be reshaped to the original image dimensions $h \times w \times c$.

For convenience, we also keep a list of all image patches $\rho \in \mathbb{R}^{p \times (s_h \times s_w) \times c}$ where each element of size $s_h \times s_w \times c$ is an image patch containing an object segment.

In the next steps, we use the **segmented image** x^o to (i) calculate the global features of the image in Section 4.4 and (ii) learn the object-level features with a VAE using the **image patches** ρ .

4.3. Learning Object-Level Latent Representations with a Variational Autoencoder

To capture the object-level attributes of the segmented image patches, we employ a Variational Autoencoder (VAE). The VAE is trained to reconstruct the extracted image patches, hereby learning a compact latent space that encodes meaningful object attributes.

VAE Architecture. The VAE consists of an encoder $q_\phi(z|\rho)$ that maps a list $\rho \in \mathbb{R}^{p \times s_h \times s_w \times c}$ of image patches (list index from 1 to p) to a list of latent vectors (arranged as a matrix) $z \in \mathbb{R}^{p \times n_z}$, and a decoder $f_\theta(\rho|z)$ that reconstructs the input patches. Here, n_z is the embedding dimension of the VAE. The training objective is given by:

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_{z \sim q_\phi(z|\rho)} [\log f_\theta(\rho|z)] - \beta D_{\text{KL}}(q_\phi(z|\rho) \| f(z)), \quad (10)$$

where $f(z) = \mathcal{N}(z; 0, 1)$ is a Gaussian prior, and β controls the regularization strength.

Training with object patches. To train the VAE we convert our images x_1, \dots, x_t into t lists of patches

$$\rho_1, \rho_2, \dots, \rho_t \in \mathbb{R}^{p \times s_h \times s_w \times c} \quad (11)$$

(again each with list index from 1 to p , see also Section 4.2). Each patch is resized to a fixed resolution and normalized before being fed into the encoder (see supplementary material C). The model is trained using the reparameterization trick and the Adam optimizer. After training, the encoder q_ϕ is used to map the segmented patches to a list of *latent vectors* (arranged as a matrix):

$$z \in \mathbb{R}^{p \times n_z}. \quad (12)$$

These vectors provide a structured representation of object attributes, which can be combined with the global DINO features to enhance downstream tasks. The final representation is obtained by concatenating the latent vector with the global DINO representation, which we will describe in the next section.

4.4. Combining Global and Local Representations

So far, we have created the *object-level features* in form of a matrix $z \in \mathbb{R}^{p \times n_z}$. The global information can be found in the CLS token $\epsilon \in \mathbb{R}^{n_g}$ resulting from the DINO model

after passing the segmented image x^o . We call ϵ the *global latent feature vector*.

To combine the global and local features, we need to repeat ϵ exactly p times to match the number of object-level latent vectors:

$$z \in \mathbb{R}^{p \times n_z}, \quad e = [\epsilon_1, \dots, \epsilon_p] \in \mathbb{R}^{p \times n_g}.$$

we now have all the necessary components to create a new representation that combines global and local features. For this we concatenate z and e to obtain the final representation

$$v = [e, z] \in \mathbb{R}^{p \times (n_g + n_z)} = \mathbb{R}^{p \times n_v}.$$

Summarizing, v is a structured representation that combines global and local features with p embedding vectors of dimension $n_v = n_g + n_z$.

Calculating similarities. We now define how to calculate the similarity on patch level between a query image x and a candidate image x' using the representations v and v' . Since v and v' share the same embedding dimension, we can calculate the similarity between the two patches $v_i \in \mathbb{R}^{n_v}$ and $v'_j \in \mathbb{R}^{n_v}$ by computing the pairwise cosine similarity between the two representations:

$$S_{ij} = \frac{v_i \cdot v'_j}{\|v_i\| \cdot \|v'_j\|} \in \mathbb{R}, \quad (13)$$

for all $i, j \in \{1, \dots, p\}$.

Then, for each query patch v_i we retrieve the patch with the highest cosine similarity, i.e., the most similar patch from v'

$$s_i^{\max} = \max_{j \in \{1, \dots, p\}} S_{ij}. \quad (14)$$

The list $S^{\max} = [s_1^{\max}, \dots, s_p^{\max}]$ then contains the cosine similarities for the most similar patches from v' for each patch in v . Finally, our *similarity score* of the query image x with another image x' is the average of the similarities in S^{\max} , that is $\bar{s} = \text{average}(S^{\max})$. Note, if we compare x against n candidate images x' , we will have n similarity scores for the given query image x .

5. Experiments

To test and compare the representations learnt with our method we propose a multi-object instance retrieval task. As shown in Caron et al. (2021), instance retrieval is a task that relies heavily on the standalone quality of the representations. The more common benchmarks however, rarely contain multi-object scenes, also the objects are usually very varied and seldomly share attributes with which a detailed analysis of the representations can be made. We thus propose to use two object-centric datasets, CLEVR

and CLEVRText, which contain images of objects with both general scene attributes (position, orientation, number of objects) and individual object attributes (colour, shape, material, size). These datasets enable us to study the pre-trained representations in a challenging multi-object setting that remains simple to interpret. In the following, we describe the datasets, evaluation metrics and experimental setup used to assess the quality of the representations.

5.1. Datasets

The CLEVR dataset (Johnson et al., 2017) is a synthetic dataset of images of 3D-rendered objects. Each image contains a set of objects with different attributes such as shape, colour, material, and size, in addition to general scene attributes like object position and orientation. The dataset is designed to test reasoning capabilities of models on complex visual scenes, making it suitable for evaluating object-centric representations. The CLEVRText (Karazija et al., 2021) dataset is an extension of CLEVR, where the objects are textured with real-world textures. This dataset is more challenging than CLEVR, as the textures introduce additional complexity to the object attributes. Both datasets contain images with multi-object scenes, making them suitable for evaluating multi-object instance retrieval.

CLEVR/CLEVRText Dataset Splits

Train	Validation	Test	Query Size
2000	500	5000	50

Table 1: **Data splits for the CLEVR and CLEVRText datasets.** The train set is used to train the VAE, while the validation set is used to select query images for the retrieval task. The test set contains the candidate images for the retrieval task.

Train-validation-test splits. In the *training phase*, we only train the VAE with the segmented image patches. Thus for training the VAE we use the multi-object samples from the training split of the CLEVR and CLEVRText datasets. We found that a training set size of 2000 images was sufficient to learn the latent space of the VAE.

In the *testing phase*, we evaluate the representations on the test-split of the datasets. We have query images and candidate images for the retrieval task. The query images are chosen from a validation set of the datasets, containing 500 images. The candidate images are chosen from the test set of the datasets, containing 5000 images. As attributes to retrieve, we define the attributes **shape** and **size** as global attributes and **material** and **colour** as object-level attributes. The former are attributes which require general object understanding, while the latter are more fine-grained object attributes. We discuss experimental details in the following.

5.2. Experimental Setup

Experimental details. To evaluate the representations, we define two retrieval tasks: *single object retrieval* and *multi-object retrieval*.

In *single object retrieval*, we evaluate the ability of the representations to retrieve a single object from a multi-object scene. In *multi-object retrieval*, we evaluate the ability of the representations to retrieve multiple objects in a multi-object scene.

In both setups, we always have a set of query images and a set of candidate images, where we use our similarity score as in Section 4.4 to calculate the similarity between the query and candidate images. We show the size of the sets for each setup in Table 1. Next, we describe the evaluation metrics used for the retrieval task.

Evaluation metrics. To evaluate the representations we check whether the attributes of the best retrievals match the attributes of the reference image. We use three evaluation metrics:

1. The **Top-10 Precision** measures the percentage of the Top-10 best retrievals that match the attributes of the reference image. That is, within the ten first retrieved images, we check whether the attributes of the images match the attributes of the reference image. This is averaged over the set of query images.
2. The **Weighted Precision** is the same as the Top-10 Precision, but we weight the precision by the rank of the retrieval. That is, the precision of the i -th retrieval is weighted by $1/i$. This metric gives more weight to the best retrievals, which is important for the retrieval task.
3. The **Error Rate** measures whether our method is able to find at least one relevant image for a given query image. That is, the error rate is the percentage of reference images for which we do not find at least one relevant image in the Top-10 retrievals.

Baselines. We test our method against the following baselines:

- **DINOseg:** We calculate the cosine similarity between DINO CLS-token, where the image input for DINO (Caron et al., 2021) is masked as shown in Section 4.2.
- **SlotDiffusion (Wu et al., 2023):** We calculate the pairwise cosine similarity as in Section 4.4 between the slot-based representations.
- **SlotDiffusion-DINOseg:** Similar to our method, we concatenate the slot representations with the DINO representations.

For each baseline we run seven trials, where each trial consists of a set of 50 query images and 5000 candidate images. The set of query images for each trial is chosen randomly without replacement from the validation split of the datasets, to ensure comparability between methods, we set seven different seeds. We present and discuss the results of our experiments in the following section.

6. Results

6.1. Single Object Retrieval

The results of single object retrieval are presented in Figure 3 (top-left). For the CLEVR dataset, we observe that our method outperforms the baselines by a large margin, achieving a Top-10 Precision of 56.4% compared to 13.0% for DINOseg and 12.0% for SlotDiffusion. Moreover, the weighted precision and error rate present a similar trend. This indicates that our method ranks the most similar images consistently higher, while also finding new images that the baselines are not able to retrieve. The results for the CLEVR-Text dataset show a similar trend, with our method achieving a Top-10 Precision of 20.3% compared to 12.8% and 1.7% for DINOseg and SlotDiffusion, respectively. Surprisingly, SlotDiffusion performs very poorly on the CLEVR-Text dataset, which we attribute to the increased complexity of the textures in the dataset. We also see a reduced relative improvement over the baselines in CLEVR-Text due to the increased complexity of the materials, that is some the differences between materials can be subtle and hard to distinguish. Furthermore, this increased number of materials in CLEVR-Text reduces the number of suitable candidates for retrieval, making the task more challenging. However, we are still able to improve the performance over the baselines, showing the effectiveness of our method.

6.2. Multi-Object Retrieval

We show the results of multi-object retrieval in Figure 3 (bottom-left). The query images contain 3 objects with different attributes and we measure whether at least two objects can be retrieved. The results show that retrieving multiple objects is a more challenging task than single object retrieval. Our method is able to improve performance upon the baselines, achieving a Top-10 Precision of 21.0% compared to 5.0% and 2.3% for DINO and SlotDiffusion, respectively. Similar to single object retrieval our method also improves the weighted precision and error rate. We attribute the effectiveness of our method to the DINO representations being able to retrieve images with a similar number of objects, while the object-level features are able to capture the individual object attributes. The results of the CLEVR-Text dataset show that multi-object retrieval is too challenging in this setting even for our method, where the general issues mentioned in the single object retrieval are

exacerbated in the multi-object retrieval task.

6.3. Object Attribute Ablation

To measure the effectiveness of our method at capturing individual object attributes, we perform an attribute ablation study where we evaluate the retrieval performance adding one attribute at a time. We show the results in Figure 3 (top-right and bottom-right). Generally, we observe that our method especially improves retrieval conditioned on object-specific features such as material and color. This is both the case for the CLEVR and CLEVR-Text datasets. The DINOseg representations are strong at understanding more global attributes such as shape and size compared to the slot-based representations. This also validates our approach of learning a specialised latent space for object-level features, which can capture fine-grained object attributes. Interestingly, while concatenating the slot-based representations with the DINOseg representations improved performance slightly, it was not able to outperform our method. This indicates that the slot-based representations are not as specialised as our patch-based representations.

7. Discussion

Limitations and Future Work. While we presented an effective method, there are a few limitations that should be considered. Our work only uses a single pre-trained model, DINO, to demonstrate the effectiveness of our method. Potentially, other larger pre-trained models like CLIP, which have shown strong generalization capabilities will prove effective as well, possibly even better. For now, we chose DINO as it has shown emergent object understanding, which is crucial for our task of multi-object instance retrieval. In this sense, our method relies on the quality of the pre-trained model to extract object-level features. On the one hand this is an advantage as it allows us to adapt our method to multiple settings without the need for extensive re-training (only the VAE). On the other hand, it might limit the ability to fine-tune to specific datasets. However, we believe that latent object representations take the role of learning fine-grained features in such a case. Future work should explore the effectiveness of our method on real-world tasks that require multi-object understanding. For instance, our method could be used to improve the performance of goal-conditioned reinforcement learning agents that rely on multi-object goal images to specify desired configurations of objects in the environment.

Conclusion. We presented a new method that enhances pre-trained representations using specialised object-level latent vectors. Our method leverages the inherent object understanding of pre-trained models to extract global object-level features hereby enabling us to learn the object-level latent

Main Results							Object Attribute Ablation			
Single-Object Retr. (%)	CLEVR			CLEVRText			CLEVR			
	Top-10 Precision \uparrow	Weighted Precision \uparrow	Error Rate \downarrow	Top-10 Precision	Weighted Precision	Error Rate	S	S+D	S+D+M	S+D+M+C
Oh-A-Dino (Ours)	56.4\pm2.0	49.0\pm2.2	1.0\pm0.0	20.3\pm2.6	11.2\pm1.4	41.0\pm1.0	94.6\pm2.2	85.2\pm2.1	72.8\pm1.4	56.4\pm1.6
DINOseg	13.0 \pm 1.1	13.8 \pm 1.2	28.0 \pm 1.0	12.8 \pm 2.1	8.0 \pm 1.3	55.6 \pm 0.7	81.2 \pm 2.8	76.2 \pm 3.1	62.8 \pm 1.2	13.0 \pm 1.1
SlotDiffusion-DINOseg	15.1 \pm 1.3	16.6 \pm 1.4	28.0 \pm 1.0	1.7 \pm 0.4	1.4 \pm 0.3	91.6 \pm 0.3	89.1 \pm 2.6	68.5 \pm 2.4	53.6 \pm 1.4	15.1 \pm 1.3
SlotDiffusion	12.0 \pm 1.2	12.7 \pm 1.3	33.6 \pm 1.5	1.4 \pm 0.4	0.9 \pm 0.2	93.4 \pm 0.2	76.3 \pm 2.7	56.3 \pm 0.4	39.5 \pm 1.3	12.0 \pm 1.2
Multi-Object Retr. (%)	CLEVR			CLEVRText			CLEVRText			
	Top-10 Precision \uparrow	Weighted Precision \uparrow	Error Rate \downarrow	Top-10 Precision	Weighted Precision	Error Rate	S	S+D	S+D+M	S+D+M+C
Oh-A-Dino (Ours)	21.0\pm1.6	11.0\pm0.8	14.4\pm1.5	0.8\pm0.2	0.4\pm0.1	96.6\pm0.2	100.0\pm0.0	80.0\pm1.4	20.3\pm2.6	
DINOseg	5.0 \pm 0.6	3.7 \pm 0.5	61.0 \pm 1.0	0.0 \pm 0.0	0.0 \pm 0.0	100 \pm 0.0	100.0 \pm 0.0	75.1 \pm 1.3	12.8 \pm 2.1	
SlotDiffusion-DINOseg	3.5 \pm 0.6	2.9 \pm 0.5	72.2 \pm 1.0	0.1 \pm 0.0	0.0 \pm 0.0	99.0 \pm 0.0	100.0 \pm 0.0	53.6 \pm 0.5	1.7 \pm 0.4	
SlotDiffusion	2.3 \pm 0.5	1.8 \pm 0.4	82.0 \pm 2.0	0.0 \pm 0.0	0.0 \pm 0.0	100 \pm 0.0	94.4 \pm 2.7	46.7 \pm 0.4	1.4 \pm 0.4	

Figure 3: **Our method outperforms the baselines on both the CLEVR and CLEVRText datasets by improving retrieval on all attributes.** For the CLEVR dataset, the largest performance improvements comes from the improved retrieval of the color attribute, while all other attributes are also improved. Our method is still effective with the more complex textures in the CLEVRText dataset, where the DINO representations are able to generalise better compared to the slot-based representations.

vectors. We showed that combining both these representations outperform traditional slot-based and self-supervised methods by combining the fine-grained object-level features learnt by a VAE with the global scene understanding of pre-trained models. We believe our work validates the approach of learning distinct specialised latent spaces for an image. As self-supervised models continue to advance, this approach unlocks new possibilities for multi-object understanding in computer vision.

References

- Bizeul, A., Schölkopf, B., and Allen, C. A Probabilistic Model Behind Self-Supervised Learning. *arXiv e-prints*, art. arXiv:2402.01399, February 2024. doi: 10.48550/arXiv.2402.01399.
- Burgess, C. P., Matthey, L., Watters, N., Kabra, R., Higgins, I., Botvinick, M., and Lerchner, A. MONet: Un-supervised Scene Decomposition and Representation. *arXiv e-prints*, art. arXiv:1901.11390, January 2019. doi: 10.48550/arXiv.1901.11390.
- Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., and Joulin, A. Emerging Properties in Self-Supervised Vision Transformers. *arXiv e-prints*, art. arXiv:2104.14294, April 2021. doi: 10.48550/arXiv.2104.14294.
- Chen, T., Huang, Y., Shen, Z., Huang, J., Li, B., and Xue, X. Learning Global Object-Centric Representations via Disentangled Slot Attention. *arXiv e-prints*, art. arXiv:2410.18809, October 2024. doi: 10.48550/arXiv.2410.18809.
- Chen, W., Liu, Y., Wang, W., Bakker, E., Georgiou, T., Fieguth, P., Liu, L., and Lew, M. S. Deep Learning for Instance Retrieval: A Survey. *arXiv e-prints*, art. arXiv:2101.11282, January 2021. doi: 10.48550/arXiv.2101.11282.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Garrido, Q., Ballas, N., Assran, M., Bardes, A., Najman, L., Rabbat, M., Dupoux, E., and LeCun, Y. Intuitive physics understanding emerges from self-supervised pretraining on natural videos. *arXiv e-prints*, art. arXiv:2502.11831, February 2025. doi: 10.48550/arXiv.2502.11831.
- Hafner, D., Lee, K.-H., Fischer, I., and Abbeel, P. Deep Hierarchical Planning from Pixels. *arXiv e-prints*, art. arXiv:2206.04114, June 2022. doi: 10.48550/arXiv.2206.04114.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Johnson, J., Hariharan, B., Van Der Maaten, L., Fei-Fei, L., Lawrence Zitnick, C., and Girshick, R. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2901–2910, 2017.

- Jung, W., Yoo, J., Ahn, S., and Hong, S. Learning to Compose: Improving Object Centric Learning by Injecting Compositionality. *arXiv e-prints*, art. arXiv:2405.00646, May 2024. doi: 10.48550/arXiv.2405.00646.
- Karazija, L., Laina, I., and Rupperecht, C. Clevrtex: A texture-rich benchmark for unsupervised multi-object segmentation. *arXiv preprint arXiv:2111.10265*, 2021.
- Le Paine, T., Gulcehre, C., Shahriari, B., Denil, M., Hoffman, M., Soyer, H., Tanburn, R., Kapturowski, S., Rabinowitz, N., Williams, D., Barth-Maron, G., Wang, Z., de Freitas, N., and Worlds Team. Making Efficient Use of Demonstrations to Solve Hard Exploration Problems. *arXiv e-prints*, art. arXiv:1909.01387, September 2019. doi: 10.48550/arXiv.1909.01387.
- Locatello, F., Weissenborn, D., Unterthiner, T., Mahendran, A., Heigold, G., Uszkoreit, J., Dosovitskiy, A., and Kipf, T. Object-Centric Learning with Slot Attention. *arXiv e-prints*, art. arXiv:2006.15055, June 2020. doi: 10.48550/arXiv.2006.15055.
- Montero, M. L., Bowers, J. S., and Malhotra, G. Successes and Limitations of Object-centric Models at Compositional Generalisation. *arXiv e-prints*, art. arXiv:2412.18743, December 2024. doi: 10.48550/arXiv.2412.18743.
- Oquab, M., Darcet, T., Moutakanni, T., Vo, H., Szafraniec, M., Khalidov, V., Fernandez, P., Haziza, D., Massa, F., El-Nouby, A., Assran, M., Ballas, N., Galuba, W., Howes, R., Huang, P.-Y., Li, S.-W., Misra, I., Rabbat, M., Sharma, V., Synnaeve, G., Xu, H., Jegou, H., Mairal, J., Labatut, P., Joulin, A., and Bojanowski, P. DINOv2: Learning Robust Visual Features without Supervision. *arXiv e-prints*, art. arXiv:2304.07193, April 2023. doi: 10.48550/arXiv.2304.07193.
- Schölkopf, B., Locatello, F., Bauer, S., Ke, N. R., Kalchbrenner, N., Goyal, A., and Bengio, Y. Towards Causal Representation Learning. *arXiv e-prints*, art. arXiv:2102.11107, February 2021. doi: 10.48550/arXiv.2102.11107.
- Seitzer, M., Horn, M., Zadaianchuk, A., Zietlow, D., Xiao, T., Simon-Gabriel, C.-J., He, T., Zhang, Z., Schölkopf, B., Brox, T., and Locatello, F. Bridging the gap to real-world object-centric learning. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=b9tUk-f_aG.
- Singh, G., Deng, F., and Ahn, S. Illiterate DALL-E Learns to Compose. *arXiv e-prints*, art. arXiv:2110.11405, October 2021. doi: 10.48550/arXiv.2110.11405.
- Singh, G., Wu, Y.-F., and Ahn, S. Simple Unsupervised Object-Centric Learning for Complex and Naturalistic Videos. *arXiv e-prints*, art. arXiv:2205.14065, May 2022. doi: 10.48550/arXiv.2205.14065.
- Uelwer, T., Robine, J., Sylvius Wagner, S., Höftmann, M., Upschulte, E., Konietzny, S., Behrendt, M., and Harmeling, S. A Survey on Self-Supervised Representation Learning. *arXiv e-prints*, art. arXiv:2308.11455, August 2023. doi: 10.48550/arXiv.2308.11455.
- van Steenkiste, S., Greff, K., and Schmidhuber, J. A Perspective on Objects and Systematic Generalization in Model-Based RL. *arXiv e-prints*, art. arXiv:1906.01035, June 2019. doi: 10.48550/arXiv.1906.01035.
- Wang, Y., Liu, L., and Dauwels, J. Slot-VAE: Object-Centric Scene Generation with Slot Attention. *arXiv e-prints*, art. arXiv:2306.06997, June 2023. doi: 10.48550/arXiv.2306.06997.
- Wu, Z., Hu, J., Lu, W., Gilitschenski, I., and Garg, A. SlotDiffusion: Object-Centric Generative Modeling with Diffusion Models. *arXiv e-prints*, art. arXiv:2305.11281, May 2023. doi: 10.48550/arXiv.2305.11281.
- Zheng, Y., Yao, L., Su, Y., Zhang, Y., Wang, Y., Zhao, S., Zhang, Y., and Chau, L.-P. A Survey of Embodied Learning for Object-Centric Robotic Manipulation. *arXiv e-prints*, art. arXiv:2408.11537, August 2024. doi: 10.48550/arXiv.2408.11537.

A. Learning Object-Level Features with PCA

A.1. Detailed View of PCA

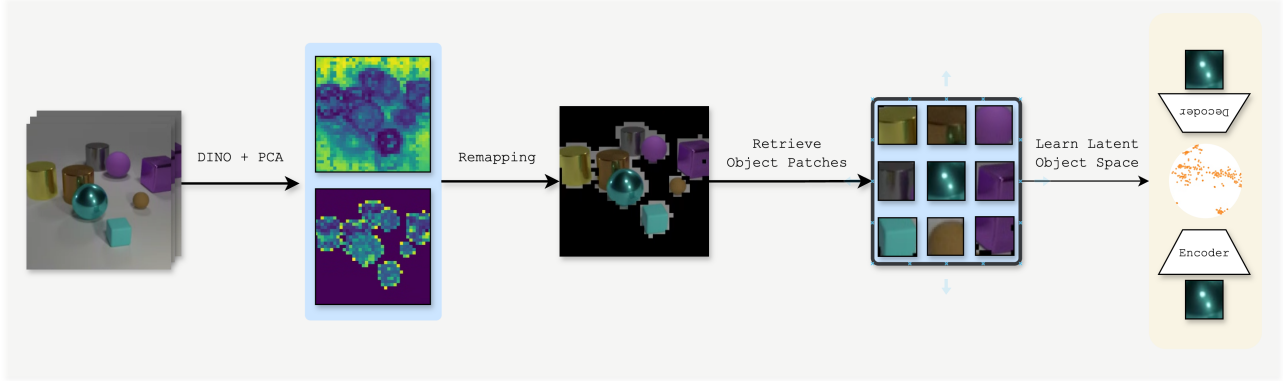


Figure 4: Detailed view of extracting object-level features using PCA as shown in Section 4.2.

A.2. The Benefit of Reapplying PCA

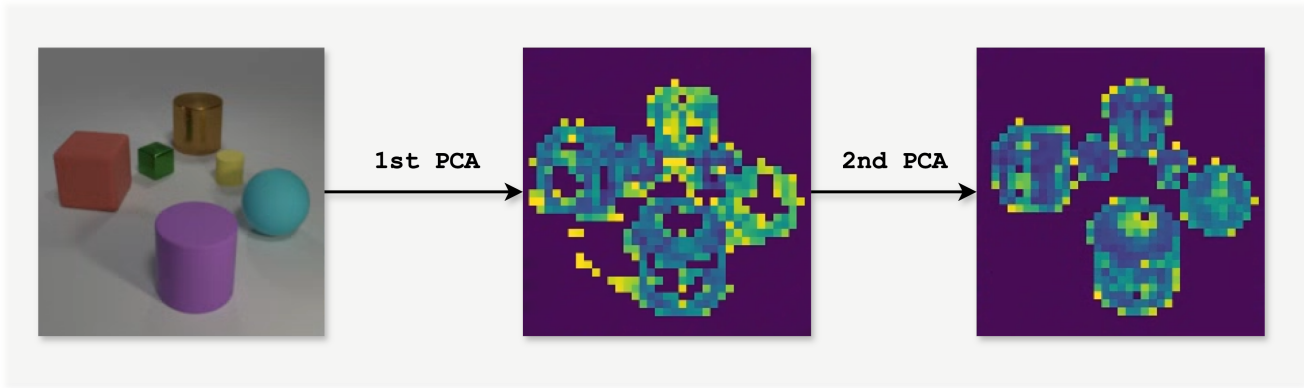


Figure 5: We show the benefit of reapplying PCA to the foreground object-level features. The first PCA passthrough to separate the background and foreground features, while the second PCA pass-through helps to separate the individual objects and get a continuous mask across the objects.

B. Additional Dataset Details

B.1. Dataset Statistics

<i>Dataset</i>	<i>#Images</i>	<i>#Objects</i>	<i>#Shapes</i>	<i>#Colors</i>	<i>#Materials</i>	<i>#Backgrounds</i>
CLEVR	100k	3-10	3	8	2	1
CLEVRTex	50k+10k	3-10	4	-	60	60

Table 2: Dataset characteristics for CLEVR and CLEVRTex. The CLEVRTex dataset has additional textures for the objects and backgrounds compared to CLEVR. It also has no color information for the objects, which is replaced by the material information. In general, instance retrieval with the CLEVRTex dataset is much harder due to the increased number of possible object configurations.

VAE Hyperparameters

<i>Parameter</i>	<i>Value</i>
Input size	64×64
Latent size	32
Learning rate	$1e-4$

C. Hyperparameter Details

C.1. VAE Hyperparameters

C.2. DINO Hyperparameters

DINO Hyperparameters

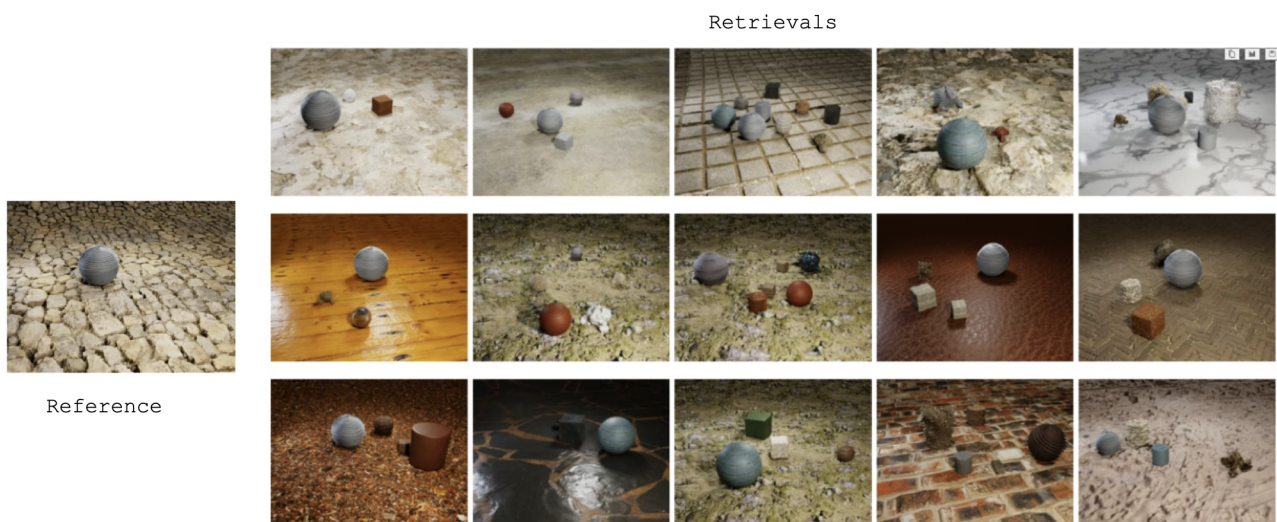
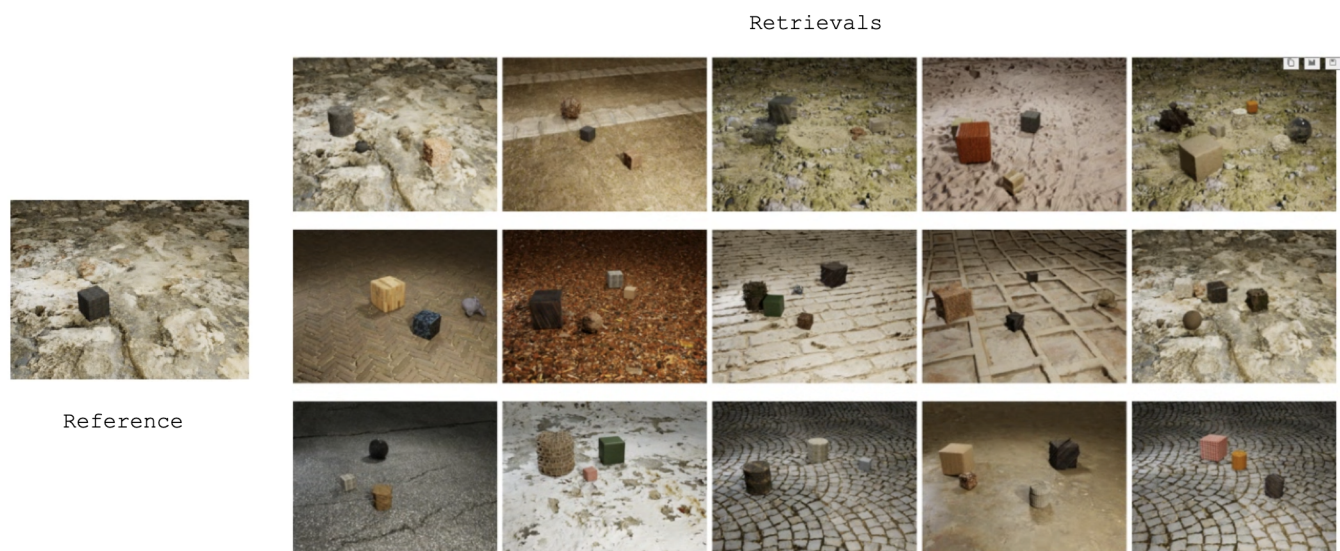
<i>Parameter</i>	<i>Value</i>
Input size	518×518
Embedding Dimension	384
Patch Size	14

D. Additional Visualisations

D.1. CLEVRTex

D.1.1. OURS





D.1.2. DINO

