Unsupervised Object Localization with Representer Point Selection

Yeonghwan Song¹

ong¹ Seokwoo Jang¹ ¹AI Graduate School, GIST Dina Katabi² ²MIT CSAIL

dina@csail.mit.edu

{yeonghwan.song, jangseokwoo}@gm.gist.ac.kr

Abstract

We propose a novel unsupervised object localization method that allows us to explain the predictions of the model by utilizing self-supervised pre-trained models without additional finetuning. Existing unsupervised and selfsupervised object localization methods often utilize classagnostic activation maps or self-similarity maps of a pretrained model. Although these maps can offer valuable information for localization, their limited ability to explain how the model makes predictions remains challenging. In this paper, we propose a simple yet effective unsupervised object localization method based on representer point selection, where the predictions of the model can be represented as a linear combination of representer values of training points. By selecting representer points, which are the most important examples for the model predictions, our model can provide insights into how the model predicts the foreground object by providing relevant examples as well as their importance. Our method outperforms the state-ofthe-art unsupervised and self-supervised object localization methods on various datasets with significant margins and even outperforms recent weakly supervised and few-shot methods. Our code is available at: https://github. com/yeonghwansong/UOLwRPS

1. Introduction

Object localization is one of the most fundamental problems in computer vision which aims to find a bounding box of a particular object category in a given image. Recent object localization methods achieve impressive performances thanks to advances in deep neural networks (DNNs) [20, 14] and large-scale datasets [35, 40]. Despite their successes, the biggest concern in this task is that collecting datasets with precise box-level annotations is laborintensive and time-consuming. Recently, object localization methods with less supervision, *i.e.* weakly-supervised methods [27, 55, 51, 13, 56, 49, 33, 47, 3], have been wellstudied to mitigate those problems, nevertheless, they require image-level class labels.



Jeany Son¹

jeany@gist.ac.kr

Figure 1. *GT-Known Loc* scores of various methods with respect to different levels of supervision on ImageNet-1K dataset. All methods use ViT-S [14] as a backbone and TC denotes TokenCut [42].

Beyond learning with less supervision, most recent works [45, 44, 54, 1, 23, 29, 38, 48, 37, 42] focus on the object localization task using self-supervised or unsupervised learning that does not require any human annotated labels. These works address the problem of identifying which regions are more likely to contain the foreground object, which is a salient object in an image. In order to discover foreground object regions, several methods [1, 23, 29, 38] attempt to use the magnitude of feature vectors as a clue for class-agnostic activation maps (CAAM) [1]. Most of these methods rely heavily on pre-trained models designed for the image classification task. However transferring these models to the object localization task is challenging, since foreground features are not easily distinguished from the background. To address this problem, these works have been focusing on learning to discriminate foreground and background representations from the pre-trained model.

Recent approaches for unsupervised object localization face a limitation in terms of the explainability of the predictions. DNNs are powerful predictors across various domains, but their complicated structure often makes them black-box functions. Class activation maps (CAM) [59] are frequently used to provide visually interpretable information about a specific class the model has learned, but a detailed explanation of how the model makes its predictions remains unclear. This limitation becomes even more severe for CAAM, which is commonly used in self-supervised object localization methods [1, 23, 29, 38].

In this paper, we propose a simple yet effective unsupervised object localization method based on the representer theorem using self-supervised representation learning models. The proposed method formulates the solution to the object localization task based on representer point selection [50]. This strategy is derived from the representer theorem, where the predictions of a neural network for a given test point find expression through a linear combination of the activations originating from training examples. The importance of each feature, which measures the significance of training examples in contributing to the final prediction, is computed by its norm. By selecting representer points that have the most influence on the model's predictions, our model gains valuable insights into how the model identifies foreground regions. In addition, we introduce a simple extension to weakly supervised object localization, demonstrating zero-shot transferability across different datasets without the need for additional training.

Our experiments are conducted on various datasets for the object localization task to prove the effectiveness of our method and show outstanding performances compared to the state-of-the-art methods with substantial margins. Furthermore, as shown in Figure 1, our method even outperforms some recent weakly supervised and few-shot object localization methods without any fine-tuning or postprocessing using additional refinement networks.

We summarize our contributions as follows:

- We propose a novel unsupervised object localization method that leverages the representer theorem on selfsupervised representation learning models. The proposed method helps to enhance the interpretability of predictions by providing relevant examples as well as their importance.
- We show that the representer point selection, which has not been extensively explored in previous research on deep neural networks, can be effectively applied to object localization.
- Our method is simple to implement and works well with a variety of self-supervised visual representation learning models. It is also an effective solution due to its good transferability across datasets and easy scalability to weakly supervised object localization.
- Our method shows outstanding results on various datasets for the object localization task and also achieves comparable performances to the state-of-the-art weakly-supervised and few-shot methods.

2. Related Works

Weakly Supervised Object Localization Weakly supervised object localization (WSOL) is a challenging task, aiming to localize an object with only image-level annotations. CAM [59] is the first work to introduce a method of generating a class activation map by projecting the weights of the classifier on the feature maps. However, CAM has a limitation in finding parts of the objects that are not relevant for classification. To address this problem, various CAM-based WSOL methods have been proposed [27, 55, 13, 56, 49, 33, 48, 47, 3]. HaS [27] used the random masking of image patches to encourage the model to observe more object regions. ACoL [55] adopted two adversarial classifiers to catch the relevant parts. ADL [13] and AE [46] progressively erased the discriminative parts on the feature map of each convolutional layer, ensuring that the focus was not solely on these discriminative areas. SPG [56] and I²C [57] utilized shallow convolutional feature maps and pixel-level correlations between two different images to adjust an activation map respectively. PSOL [53] proposed a new task of pseudo-supervised object localization which consists of two independent tasks: class-agnostic object localization and object classification. PSOL [53] and SEM [58] trained a localizer with an object bounding box generated by DDT [45] and edges of an input image extracted from the canny edge detector as pseudo labels, respectively. Recently, most CAM-based methods [49, 43, 33, 48, 47, 3] have focused on learning discriminative features for the foreground and background, while vision transformer-based methods [15, 18, 10, 2] have emerged to exploit long-range dependency.

Unsupervised Object Localization Unsupervised object discovery is strongly related to co-localization [45, 1, 23] and co-segmentation [52, 30] tasks. Early work for unsupervised object discovery [11] proposed a part-based region matching using off-the-shelf region proposals. With the advance of deep learning, DDT [45] and SCDA [44] introduced methods to mine the statistical properties of object regions in an image using features from pre-trained CNNs. MO [54] mined a frequently activated pattern from multiscale feature maps. These methods achieved impressive results not only in object discovery but also in co-localization, which aims to localize objects in a dataset that consists of only one coarse category. PsyNet [1] introduced a simple and practical framework for the co-localization task using self-supervised learning and class-agnostic activation mapping. Ki et al. [23] employed contrastive learning with the assumption of consistency between image-level and activation-level transformations. Su et al. [38] proposed a joint graph partition method splitting foreground and background regions in a paired image feature graph. $C^{2}AM$ [48] proposed to apply contrastive learning to disentangle the feature map into foreground and background with an additional trainable layer generating a class-agnostic activation map. LOST [37] and TokenCut [42] utilize a pre-trained vision transformer model [14] by DINO [5], which has the ability to produce fine segmentation masks using attention maps with [CLS] token. LOST [37] suggested a heuristic seed selection and expansion method, and TokenCut [42] employed a normalized cut algorithm on the graph using the similarity between patches.

Self-supervised Visual Representation Learning Selfsupervised visual representation learning (SSL) methods have been proposed to learn meaningful visual representation without manual annotations. The network learns visual representation by training pretext tasks rather than the task intended to be solved. For example, Noroozi et al. [34] trained CNN to solve a jigsaw puzzle, Larsson et al. [28] reconstructed color images from gray-scale images, and Gidaris et al. [16] predicted discrete rotation angles of rotated images. Recently, contrastive learning-based SSL methods [19, 6, 8, 17] have emerged with comparable performances to the supervised learning methods. They focus on modeling image similarity and dissimilarity between pairs of images [6, 19] or only consider image similarity [17, 8] based on data augmentation. DINO [5] proposed a knowledge distillation strategy to train the vision transformer model. These pre-trained SSL models are utilized in various downstream tasks of computer vision with finetuning. On the contrary, in this work, we focus on leveraging the pre-trained SSL models for the unsupervised object localization task without any fine-tuning.

3. Background

Representer Theorem. Representer theorem [36] provides a mathematical basis for traditional machine learning methods. Let $L(\mathbf{x}_n, y_n, \mathbf{w})$ be the loss function, and $\frac{1}{N} \sum_{n=1}^{N} L(\mathbf{x}_n, y_n, \mathbf{w})$ be the empirical risk over a reproducing kernel Hilbert space (RKHS) H_k . Then the representer theorem states that an optimal solution of a regularized empirical risk minimization (ERM) problem can be represented as a linear combination of a positive definite kernel k on the input set \mathcal{X} over H_k . The optimal parameter \mathbf{w}^* can be expressed through ERM as follows:

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}\in H_k} \{ \frac{1}{N} \sum_{n=1}^N L(\mathbf{x}_n, y_n, \mathbf{w}) + \lambda ||\mathbf{w}||^2 \}.$$
(1)

can be rewritten by the representer theorem

$$\mathbf{w}^{*}(\cdot) = \sum_{i=1}^{n} \alpha_{i} k(\cdot, \mathbf{x}_{i}) = \sum_{i=1}^{n} \alpha_{i} \langle \varphi(\cdot), \varphi(\mathbf{x}_{i}) \rangle, \qquad (2)$$

where $\alpha_i \in \mathbb{R}$, and φ denotes a mapping function into H_k . For further details for the Representer theorem, please refer to the work by [36, 50].

Representer Point Selection for DNNs. The representer theorem is originally developed for non-parametric predictors, where the parameters lie in a reproducing kernel Hilbert space. In deep neural networks, however, it is difficult to find a global solution for empirical risk minimization. To address this issue, [50] proposes a technique to decompose the pre-activation predictions into a linear combination of the activation values of the training points. The neural network can be represented as

$$\hat{y}_t = \sigma(\phi(\Phi(\mathbf{x}_t, \mathbf{W}_{\Phi}), \mathbf{w}_{\phi})) = \sigma(\mathbf{w}_{\phi}^{\top} \mathbf{f}_t), \qquad (3)$$

where σ is the activation function, $\mathbf{f}_n = \Phi(\mathbf{x}_n, \mathbf{W}_{\Phi})$ is the last intermediate layer feature of a training example \mathbf{x}_n in DNNs, ϕ denotes the last prediction layer for a specific task, and \mathbf{W}_{Φ} and \mathbf{w}_{ϕ} denote the parameters of Φ and ϕ , respectively. Let \mathbf{W}^* be a stationary point of the objective (1). Then we have decomposition as follows:

$$\mathbf{w}_{\phi}^{*\top}\mathbf{f}_{t} = \sum_{n=1}^{N} \alpha_{n} k(\mathbf{x}_{t}, \mathbf{x}_{n}) = \sum_{n=1}^{N} \alpha_{n} \mathbf{f}_{n}^{\top} \mathbf{f}_{t}, \qquad (4)$$

where
$$\alpha_n = \frac{1}{-2\lambda N} \frac{\partial L(\mathbf{x}_n, y_n, \mathbf{W})}{\partial \phi(\mathbf{f}_n, \mathbf{w}_{\phi})},$$
 (5)

and $\alpha_n \mathbf{f}_n^{\top} \mathbf{f}_t$ is a representer value of a training sample \mathbf{x}_n given a test sample \mathbf{x}_t . This representer value is assigned to each training point reflect their importance on the learned parameters and α_n is the global sample importance which is used to evaluate the importance of the training data \mathbf{x}_n . Consequently, this approach enables the selection of representer points – training instances with either large or small representer values – to enhance comprehension of the model's predictions.

4. Method

In this section, we discuss how we leverage the representer point selection [50] to solve unsupervised object localization (UOL) task, and our extension to weakly supervised object localization (WSOL) and zero-shot transferring.

4.1. Representer Point Selection for UOL

Inspired by the representer point selection, we formulate the object localization task as identifying the important regions in training examples, called representer points, using a pre-trained self-supervised model. Unlike other tasks such as classification, regression, and retrieval that are covered in [50], the object localization task requires dense predictions within a single image. To address this, we address all elements within image feature maps rather than relying solely on a single image-level feature vector.



Figure 2. Pipelines of our unsupervised object localization method, which consists of two separate stages: representer point selection and inference. In the representer point selection stage, we compute a foreground predictor \mathbf{w}^* using images in the training dataset. In the inference stage, the activation maps are computed by a dot product between the feature map of the test image and \mathbf{w}^* .

As our focus lies in unsupervised learning, the empirical risk or the loss function can be defined as $\frac{1}{N_D} \sum_{n=1}^{N_D} L(\mathbf{x}_n, \mathbf{W}_{\Phi})$ without a label y_n given pre-trained encoder Φ and training set D. Then the activation maps \mathbf{A}_t for foreground regions can be computed by aggregating feature maps using the global sample importance, which indicate how much each training example contributes to the prediction of being foreground:

$$\operatorname{vec}(\mathbf{A}_{t})_{j} = \underbrace{\sum_{n=1}^{N_{D}} \sum_{i=1}^{WH} \alpha_{n,i} \mathbf{\hat{f}}_{n,i}^{\top}}_{\mathbf{w}^{*}} \mathbf{\hat{f}}_{t,j} = \mathbf{w}^{*\top} \mathbf{\hat{f}}_{t,j}, \quad (6)$$

where $\mathbf{f}_{n,i} = \Phi(\mathbf{x}_n, \mathbf{W}_{\Phi})_i$ denotes i^{th} feature vector in the feature map of an image \mathbf{x}_n , which has been encoded using a self-supervised pre-trained encoder Φ , $\mathbf{\hat{f}} = \frac{\mathbf{f}}{\|\mathbf{f}\|}$ indicates the normalized feature vector, α is the global sample importance, and N_D , W, H denote the number of images in the training dataset, the width of the feature map, and the height of the feature map, respectively. Here, we refer to $\alpha_{n,i} \mathbf{\hat{f}}_{n,i}^{\top} \mathbf{\hat{f}}_{t,j}$ as the *representer value* for each training image patch $\mathbf{x}_{n,i}$ given a testing image patch $\mathbf{x}_{t,j}$, corresponding to receptive fields of $\mathbf{f}_{n,i}$ and $\mathbf{f}_{t,j}$, respectively. In Figure 2, we illustrate the entire structure of the proposed representer point selection method to compute a foreground predictor \mathbf{w}^* and inference pipelines using it.

Global Sample Importance. α_n is referred to as the global sample importance since it is independent of the testing examples but reliant on the empirical risk of the training set. In order to tackle the unsupervised localization task, we first design the loss function which is suitable for an unsupervised setting to compute α . Let w be the parameters of a binary classifier for foreground predictions given feature vectors, then the loss function can be written as follows:

$$L(\mathbf{x}_{n,i}, y_{n,i}, \mathbf{W}_{\Phi}) = -y_{n,i}\mathbf{w}^{\top} \hat{\mathbf{f}}_{n,i},$$
(7)

where $y_{n,i} \in \{-1, 1\}$ denotes a label of $\mathbf{x}_{n,i}$, which is an image patch corresponding to a receptive field of $\mathbf{f}_{n,i}$. In the unsupervised object localization task, however, the goal is to find foreground regions in images without the use of labels. Thus, instead of using $y_{n,i}$, we use the norm of the feature vector to determine a soft pseudo label of the example for the objective function as follows:

$$L(\mathbf{x}_{n,i}, \mathbf{W}_{\Phi}) = -(\|\mathbf{f}_{n,i}\| - \tau)\mathbf{w}^{\top} \hat{\mathbf{f}}_{n,i}, \qquad (8)$$

where τ denotes a foreground threshold.

We estimate soft labels of inputs based on the norm of feature vectors because the norm can serve as an important indicator of the features that the network emphasized during training with the self-supervised contrastive loss. In other words, a higher norm of a feature vector indicates that the corresponding feature had more weight in the model training. These features are extracted from regions of the image that contain rich information to distinguish them from other instances, and these regions are more likely to be associated with significant objects. Likewise, prior self-supervised co-localization studies [1, 23, 29, 38] also use the magnitude of feature vectors as a clue to discover object regions.

Hence, the global sample importance is computed as:

$$\alpha_{n,i} = -\frac{1}{2\lambda N_D H W} \frac{\partial L(\mathbf{x}_{n,i}, \mathbf{W}_{\Phi})}{\partial \mathbf{w}^{\top} \mathbf{\hat{f}}_{n,i}}$$
(9)

$$= \frac{1}{2\lambda N_D H W} (||\mathbf{f}_{n,i}|| - \tau) = \frac{1}{C} (||\mathbf{f}_{n,i}|| - \tau).$$
(10)

Here, the sign of the global sample importance is related to whether to *representer points* – training examples that have large or small representer values – contribute to the foreground or background. In order to have large representer values, both α and similarity between two features, expressed as $\mathbf{f}_n^{\top} \mathbf{f}_t$, should have large values. By analyzing the representer values, we can gain insight into how the model predicts foreground regions by providing relevant examples



Figure 3. Illustration of how our method computes activation maps using two example points, yellow and green stars. To predict the point as a foreground region, two factors should be considered: the global sample importance and the similarity of features. the representer values are obtained by multiplying these two factors, and a high representer value indicates a stronger contribution to the foreground prediction.

Algorithm 1: Representer Point Selection
model: pre-trained encoder Φ
const: the constant value for $lpha$
th: foreground threshold of norm $ au$
for x in data loader:
f = model(x) # f: [N, C, H, W]
v += sum(f,d=(N,H,W))
u += sum(normalize(f,d=C),d=(N,H,W))
w = v/const - th*u/const
return w

as well as their importance. Figure 3 shows an example of global sample importance, the similarity between features, and representer values for given points in a given test image.

Straightforward Computation of w^{*} We then compute a foreground predictor w^{*} using $\alpha_{n,i}$ in Eq (10) as follows:

$$\mathbf{w}^* = \frac{1}{C} \underbrace{\sum_{n=1}^{N_D} \sum_{i=1}^{HW} \mathbf{f}_{n,i}}_{\mathbf{v}} - \frac{\tau}{C} \underbrace{\sum_{n=1}^{N_D} \sum_{i=1}^{HW} \hat{\mathbf{f}}_{n,i}}_{\mathbf{u}}, \qquad (11)$$

The optimal \mathbf{w}^* can be computed by taking the difference between the mean feature vector of all images in the dataset and the mean of the normalized feature vectors of all images in the dataset, multiplied by two constant values. Algorithm 1 illustrates the Python-style pseudo code of the proposed representer point selection method to compute the foreground predictor \mathbf{w}^* and it clearly shows the simplicity and reproducibility of our method.

4.2. Towards WSOL and Zero-shot Transferring

Our method can be easily extended to weakly supervised and zero-shot transferring settings with a simple modification, since activation maps are computed by the linear combination of activations of training points directly without learning. Given a class label c, the activation map \mathbf{A}^c is computed as follows:

$$\operatorname{vec}(\mathbf{A}_{t}^{c})_{j} = \sum_{n \in D_{c}} \sum_{i=1}^{WH} \alpha_{n,i} \mathbf{\hat{f}}_{n,i}^{\top} \mathbf{\hat{f}}_{t,j} = \mathbf{w}_{c}^{*\top} \mathbf{\hat{f}}_{t,j}, \quad (12)$$

where D_c denotes the subset of the training dataset D that contains only the images assigned to class c. Note that, although WSOL allows for the use of image-level class labels, we utilize an off-the-shelf classification model to determine the class of each image in D during our experiments in order to ensure a fair comparison with other methods [48, 53]. Extending our method to zero-shot transferring scenarios can be simply done by computing the weight vector \mathbf{w}^* on the different datasets in which the classes do not overlap with the test dataset.

5. Experiments

5.1. Implementation Details

We adopt ResNet50 [20] and ViT-S [14] as backbones for SSL pre-trained model, utilizing publicly available pretrained weights. We use the input image size of 224×224 for the representer point selection. To evaluate object localization performances on fine-grained classification datasets, the input image is resized to 480×480 and then centercropped to a size of 448×448 . In the case of ImageNet-1K, the input image is resized to 256×256 and then centercropped to a size of 224×224 . We use $\|\mathbf{v}\| / \|\mathbf{u}\|$ for τ , and min-max normalization to estimate the bounding box from the generated activation map. We also set $\lambda = 0.001$ as in [50], but changing this value does not have an impact on the final results due to the use of min-max normalization.

5.2. Datasets and Evaluation Metrics

We conduct experiments on the four fine-grained datasets for unsupervised object localization to evaluate our

Table 1. Comparison between our method and the state-of-the-art self-supervised/unsupervised methods in terms of *GT-Known Loc* performance on ImageNet-1K and four fine-grained datasets: CUB200-2011, Stanford Cars, FGVC-Aircraft, and Stanford Dogs. * denotes our reproduced results from their official code. † and ‡ indicate a backbone initialized by MoCo v2 and DINO pre-trained weights, respectively.

Method	Backbone	\mathcal{T}	CUB-200-2011	Stanford Cars	FGVC-Aircraft	Stanford Dogs	ImageNet-1K
Self-supervised methods							
ORE .22 [29]	VGG16	1	87.72	97.43	97.86	85.45	-
JGP _{'22} [38]	VGG16	1	88.83	97.73	96.72	-	-
PsyNet .20 [1]	SE-ResNet50	1	85.10	98.81	97.81	77.84	-
Ki et al. _{'21} [23]	SE-ResNet50	1	85.93	98.95	98.75	80.32	-
PSOL ·20 [53]	ResNet50	1	90.00	-	-	-	65.44
C ² AM _{'22} [48]	ResNet50 [†]	1	89.90	99.44*	98.65*	89.91*	66.51
C ² AM (+PSOL) _{'22} [48]	ResNet50 [†]	1	91.54	-	-	-	68.07
w/o finetuning							
DDT _{'19} [45]	VGG16	X	82.26	71.33	92.53	-	-
MO ₂₀ [54]	VGG16	X	80.45	92.51	94.94	80.70	-
LOST _{'21} [37]	ViT-S [‡]	X	89.70	-	-	-	60.00
TokenCut ·22 [42]	ViT-S [‡]	X	91.80	-	-	-	65.40
Ours	VGG16	X	90.47	98.31	96.88	89.67	59.07
Ours	ResNet50 [†]	X	96.67	99.69	98.71	95.07	66.89
Ours	ViT-S [‡]	X	91.61	99.84	99.22	94.71	73.65

Table 2. Comparison of localization performances in terms of *Top-1*, *Top-5* and *GT-known Loc* on CUB-200-2011 test set and ImageNet-1K validation set. Self-supervised and unsupervised methods including ours employ an additional classifier to evaluate in WSOL setup following PSOL [53]. \mathcal{T} and \mathcal{S} indicate whether each method uses training and image-level supervision (B for bounding box annotations), respectively. \dagger , \ddagger , and \circ indicate the use of MoCo v2, BYOL, and DINO pre-trained weights, respectively.

Method	Backbone	S	τ		CUB-200-201	1		ImageNet-1k	K
			,	Top-1 Loc	Top-5 Loc	GT-Known	Top-1 Loc	Top-5 Loc	GT-Known
Few-shot method									
SPNet '19 [31]	ResNet50	√ (B)	1	-	-	95.80	-	-	71.50
Weakly supervised methods									
FAM _{'21} [33]	ResNet50	1	1	73.74	-	85.73	54.46	-	64.56
CREAM ·22 [49]	ResNet50	1	1	76.03	-	89.88	55.66	-	69.31
BGC [,] 22 [24]	ResNet50	1	1	73.16	86.68	91.60	53.76	65.75	69.89
DAOL •22 [61]	ResNet50	1	1	66.65	-	81.83	55.84	-	70.27
BAS ₂₂ [47]	ResNet50	1	1	77.25	90.08	95.13	57.18	68.44	71.77
BagCAM [,] 22 [60]	ResNet50	1	1	69.67	-	94.01	44.24	-	72.08
TS-CAM _{'21} [15]	ViT-S	1	1	71.30	83.80	87.70	53.40	64.30	67.60
LCTR ₂₂ [10]	ViT-S	1	1	79.20	89.90	92.40	56.10	65.80	68.70
ViTOL '22 [18]	ViT-S	1	1	-	-	80.90	58.64	-	72.51
SCM '22 [2]	ViT-S	1	1	76.40	91.60	96.60	56.10	66.40	68.80
Self-supervised methods									
PSOL '20 [53]	ResNet50	×	1	70.68	86.64	90.00	53.98	63.08	65.44
$C^{2}AM_{22}$ [48]	ResNet50 [†]	×	1	-	-	89.90	-	-	66.51
C ² AM (+PSOL) '22 [48]	ResNet50 [†]	×	1	74.76	87.37	91.54	54.65	65.05	68.07
w/o finetuning									
LOST _{'21} [37]	ViT-S°	X	X	71.30	-	89.70	49.00	-	60.00
TokenCut _{'22} [42]	ViT-S°	X	X	72.90	-	91.80	52.30	-	65.40
Ours	ResNet50 [†]	×	X	79.57	92.60	96.67	55.60	66.05	69.10
Ours	ResNet50 [‡]	×	X	79.55	92.58	96.72	56.03	66.78	69.93
Ours	ViT-S°	×	×	74.97	84.24	91.61	62.03	71.96	74.44

method: CUB-200-2011 [40], Stanford Cars [25], FGVC-Aircraft [32], and Stanford Dogs [22]. We also evaluate our method on ImageNet-1K to compare with other state-ofthe-art methods in unsupervised and WSOL setups. In addition, we evaluate the segmentation quality of our method by comparing it with the WSOL state-of-the-art methods on CUB-200-2011 and OpenImages-30K [4] datasets. Following [59], we use Top-1 (*Top-1 Loc*), Top-5 (*Top-5 Loc*) and GT-known localization (*GT-known Loc*) accuracy for evaluation. *GT-known Loc* computes the ratio of the images where the Intersection over Union (IoU) between the estimated bounding box and the known ground-truth boxes is greater than 50%. *Top-1 Loc* and *Top-5 Loc* compute the ratio of the images correctly classified with Top-1 and Top-

Table 3. Comparison between the proposed method and the stateof-the-art weakly supervised object localization methods in terms of PxAP and PIoU on CUB-200-2011 and OpenImages-30K. * denotes that the results are obtained from their official code.

Method	CUB-2	00-2011	OpenIma	ages-30K
	PIoU	PIoU PxAP		PxAP
WSOL methods				
BGC _{'22} [24]	-	-	-	63.70
CREAM ·22 [49]	-	-	-	64.70
DAOL •22 [61]	56.18	74.70	49.68	65.42
BagCAM [,] 22 [60]	74.51	90.38	52.17	67.68
Self-supervised methods				
C ² AM _{'22} [48]	69.65*	88.74*	47.75*	58.28*
Ours with SSL/WSOL pre	-trained ba	ickbone		
MoCo v2 [7] + Ours	71.28	87.03	50.20	62.18
DenseCL [41] + Ours	70.11	87.74	54.68	68.32
DINO ViT-S [5] + Ours	78.97	94.37	61.24	73.14
C ² AM [48] + Ours	71.54	90.22	56.15	69.79
DAOL [61] + Ours	67.37	85.99	53.25	67.62

5 predictions, respectively, with the condition of *GT-Known Loc*. For segmentation evaluation, *Pixel-wise average precision* (*PxAP*) computes the area under the pixel precision-recall curve and *Peak-IoU* (*PIoU*) is the best IoU score with various thresholds.

5.3. Results

Unsupervised Object Localization. Our method is primarily designed for the unsupervised object localization (UOL) task which does not focus on localizing objects for specific categories. In this task, w* is computed using the entire dataset, regardless of the number of categories in the dataset. We report the performances of our method and the state-of-the-art methods in terms of GT-known Loc on five commonly used datasets in Table 1. We compare our method with recent self-supervised object localization methods including ORE [29], PsyNet [1], Ki et al. [23], JGP [38] and C²AM [48], as well as object localization methods without finetuning including DDT [45], MO [54], LOST [37], and TokeneCut [42]. As shown in Table 1, our method significantly surpasses other methods on all benchmark datasets. Our method outperforms not only unsupervised methods but also the self-supervised methods where the models are finetuned on the training split of the target datasets. Even with a shallower architecture, such as VGG-16, our method shows better performance compared to the methods using ResNet50 on CUB-200-2011. For the VGG-16 backbone, we employ ImageNet pre-trained weights.

Weakly Supervised Object Localization. In addition, as mentioned in Section 4.2, our method can be extended to the weakly supervised object localization (WSOL) setting by computing \mathbf{w}_c^* for each class *c* using a pre-trained classifier. To evaluate on *Top-1/5* Loc metrics, we follow [53] which divides the WSOL task into two sub-tasks: object localiza-

Table 4. *GT-Known Loc* results of our method on various SSL models with ResNet50 backbones except DINO ViT.

Model	BYOL	MoCo v2 / v3	DINO Res / ViT	SimSiam
CUB	96.72	96.67 / 97.03	96.15 / 91.61	94.99
ImageNet	67.05	66.89 / 66.31	65.15 / 73.65	64.10
ImageNet (w*)	69.93	69.10 / 69.34	69.15 / 74.44	67.57

Table 5. *GT-Known Loc* results of models with and without our inference using the representer point selection.

Method	CUB	ImageNet
MoCo v2 + CAAM	86.93	61.19
MoCo v2 + Ours	96.67 (+9.74)	69.10 (+7.91)
C ² AM [48]	89.90	66.51
$C^{2}AM$ [48] + PSOL	91.54 (+1.64)	68.07 (+1.56)
C ² AM [48] + Ours	96.86 (+6.96)	69.80 (+3.29)
DAOL [61]	81.83	70.27
DAOL [61] + Ours	89.44 (+7.61)	70.97 (+0.70)

tion and classification. We first localize objects, and then we determine classes for localized objects using the pre-trained classifier, *i.e.* ResNet50 pre-trained on ImageNet.

Table 2 shows the performance of our method and the state-of-the-art weakly supervised object localization methods, including FAM [33], CREAM [49], BGC [24], BAS [47], DAOL [61], BagCAM [60], TS-CAM [15], LCTR [10], ViTOL [18], and SCM [2] on CUB-200-2011 and ImageNet-1K datasets. We also compare our method with PSOL [53] and C²AM [48] which are self-supervised methods. Our method shows outstanding results on CUB-200-2011 and comparable results on ImageNet-1K, compared to the advanced weakly supervised methods without any supervision and training. As shown in Table 2, our method outperforms all self-supervised methods with the additional classifier using the same architecture with localization backbone [20, 14]. For a fair comparison, we report $C^{2}AM$ initialized with MoCo v2 pre-trained weights as well as its refinements with PSOL. Note that, without any training and post-processing as PSOL, our method significantly surpasses C²AM (+PSOL) and PSOL on CUB-200-2011. For a fine-grained dataset such as CUB-200-2011, we utilize a class-agnostic foreground predictor w* for comprehensive object localization rather than a class-specific one.

Our method shows outstanding performance on not only localization but also segmentation, as shown in Table 3, particularly with the model trained by self-supervised pixel-wise contrastive learning (*i.e.* DenseCL [41]) on the OpenImages-30K dataset. All the methods compared in this table employ ResNet50 as their backbone, with the exception of DINO. Upon utilizing the DINO ViT model as our backbone, we observe a noteworthy and substantial enhancement in performance. Exploiting our method as an inference method for feature extractors trained by the WSOL methods [48, 61] further enhances performance, yielding

Table 6. *GT-Known Loc* results of C^2AM and ours using the model trained/computed on CIFAR₁₀, CIFAR₁₀₀, and CUB-200-2011 to show the zero-shot transferability of the methods.

Method	Training	CUB	Cars	Aircraft	Dogs	ImgNet
C^2AM	CIFAR ₁₀	73.18	96.74	98.35	80.54	56.51
Ours	CIFAR ₁₀	94.24	99.04	98.26	92.39	63.87
C^2AM	CIFAR ₁₀₀	65.12	93.40	95.68	78.75	57.29
Ours	CIFAR ₁₀₀	92.51	98.13	97.66	89.67	63.72
C ² AM	CUB	89.90	90.85	89.80	90.44	59.08
Ours	CUB	96.67	98.20	98.32	93.14	65.24



Figure 4. *GT-Known* scores with respect to different thresholds using various pre-trained models on the CUB-200-2011 dataset. The proposed threshold τ consistently achieves near-optimal scores.

significant improvements. It clearly shows that our method demonstrates consistently good performances on various datasets and diverse metrics.

5.4. Ablation Study

We conduct comprehensive ablation studies to verify the effectiveness of our method. We adopt the *GT-Known Loc* metric as a performance evaluation for all experiments.

Robustness across various SSL pre-trained models. We report *GT-known Loc* scores of our method using popular self-supervised pre-trained models including BYOL [17], MoCo v2 [7], MoCo v3 [9], DINO [5] and SimSiam [8] to support compatibility of our method in Table 4. Our method consistently shows good performance with various pre-trained models. We believe that our method generates more precise activation map of the object, as long as good representations are provided by pre-trained models no matter how they are trained.

Effectiveness of representer point selection. Table 5 reports the results of our inference method using representer point selection on several models including C^2AM [48] and DAOL [61] as well as class-agnostic activation map (CAAM) proposed in [1]. The backbones for our model and C^2AM are initialized by MoCo v2 [7] pre-trained weights. Even with the state-of-the-art models trained for the object localization task, applying our inference improves their performance with significant margins on both CUB-200-2011 and ImageNet-1K datasets. In the case of C^2AM , our inference significantly improves its performance more than us-

Table 7. Ablation study on the dataset sampling ratio for the representer point selection. \dagger denotes the use of a CIFAR₁₀ for the foreground predictor \mathbf{w}^* , and \ddagger indicates the use of ResNet50 as the backbone; otherwise, ViT-s is employed.

Method (sampling ratio)	CUB	ImageNet
w/o finetuning		
Ours (100%)	91.62	73.65
Ours (10%)	91.70 ± 0.11	73.64 ± 0.03
Ours (1%)	91.65 ± 0.42	73.63 ± 0.03
Ours (0.1%)	91.08 ± 1.02	73.67 ± 0.06
Ours $(0\%^{\dagger})$	92.73	73.10
TokenCut [42] (100%)	91.80	65.40
Weakly supervised method		
ViTOL [18] (100%)	80.90	72.51
Few-shot method		
SPNet [‡] [31] (~10%/~1%)	95.80	71.50

ing post-processing by PSOL [53] which requires additional training using a ResNet50 network.

Zero-shot transferability across datasets. We also examine the transferability of our model across different datasets in Table 6. We compare our method's performance with C^2AM [48], using our foreground predictor and C^2AM 's trained model from CIFAR₁₀, CIFAR₁₀₀ [26] and CUB-200-2011 datasets, and then testing on other datasets without additional training. In this comparison, we ensure fairness by using class-agnostic weights and initializing the ResNet50 model with pre-trained weights from MoCo v2 [7]. The results in Table 6 highlight that our model, leveraging representer points from these datasets, can be transferred to various datasets with strong performance, while C^2AM 's performance notably drops.

Impact of sampling rate on a dataset D. In order to show the impact of the number of samples for selecting representer points, we randomly sample images on the training set. In Table 7, to illustrate the effect of sampling rate, we report the results of our method with DINO ViT-s using 100%, 10%, 1%, 0.1% and 0% of the dataset, where 0% uses CIFAR₁₀ dataset to establish representer points. For 10%, 1%, and 0.1% sampling rates, we compute mean and standard deviation with 10 random trials. As shown in Table 7, the number of samples to compute the foreground predictor has no significant effect on performance drops. In fact, in certain instances, it even produces better results compared to using the entire dataset. We also report results of ViT-s based models, such as TokenCut [42], ViTOL [18], and a few-shot method, SPNet [31]. SPNet uses 5-shots for 200 fine-grained classes on CUB-200-2011 and 10-shots for 1K classes on ImageNet-1K, with object localization annotations, while ViTOL and TokenCut uses a full dataset to train the model. Even with an extremely small amount of samples in the dataset, our method surpass SPNet, ViTOL,



Figure 5. Qualitative results of our method compared to C^2AM . Ground truth bounding boxes are denoted in red, whereas predicted object localization bounding boxes are denoted in green or blue. Best viewed in color.



Figure 6. Failure cases of our method on ImageNet-1K (left) and CUB-200-2011 (right) datasets.

and TokenCut on ImageNet.

Impact of threshold τ . We conduct an empirical analysis to show the impact of a threshold $\tau = \|\mathbf{v}\| / \|\mathbf{u}\|$. As demonstrated by the results in Figure 4, we found that the proposed threshold τ consistently achieved nearly optimal scores, even when tested with various thresholds. Our defined threshold τ also consistently performs well across diverse models. Further detailed theoretical analysis for τ is provided in *our supplementary material*.

5.5. Qualitative Results and Failure Cases

We illustrate the qualitative results of our method and C^2AM in Figure 5. Both methods use the MoCo v2 model with the ResNet50 backbone. Without any training on the

dataset unlike C^2AM , our method successfully localizes objects on various datasets. Compared with the other method, our activation map shows the more precise area of the object, and activations in the background are suppressed.

We also demonstrate some failure cases in Figure 6. Since our method only relies on features from the pretrained networks without additional training, it has limited ability to find the small target object particularly in scenarios where multiple objects coexist within the image. We observe that our method tends to find objects with a dominant size in the image. Moreover localizing the whole object is challenging in cases where the object is either occluded or fragmented. These failure cases, however, are not only our limitations but also the limitations of current unsupervised and weakly supervised object localization methods.

6. Conclusion

In this paper, we propose a simple and effective unsupervised object localization method without using additional training and supervision. We leverage the representer point selection in the object localization task, which allows us to explain the decision in terms of the activations of data points. The proposed method outperforms not only the state-of-the-art self-supervised and unsupervised object localization methods by a significant margin but also surpasses recent weakly supervised and few-shot methods.

Acknowledgement. This work was supported by the IITP grants (No.2019-0-01842, No.2021-0-02068, No.2022-0-00926) funded by MSIT, the ISTD program (No.20018334) funded by MOTIE, and the GIST-MIT Research Collaboration grant funded by GIST, Korea.

References

- Kyungjune Baek, Minhyun Lee, and Hyunjung Shim. Psynet: Self-supervised approach to object localization using point symmetric transformation. In *AAAI*, 2020. 1, 2, 4, 6, 7, 8
- [2] Haotian Bai, Ruimao Zhang, Jiong Wang, and Xiang Wan. Weakly supervised object localization via transformer with implicit spatial calibration. In *ECCV*, 2022. 2, 6, 7, 13
- [3] Soufiane Belharbi, Aydin Sarraf, Marco Pedersoli, Ismail Ben Ayed, Luke McCaffrey, and Eric Granger. F-cam: Full resolution class activation maps via guided parametric upscaling. In WACV, 2022. 1, 2
- [4] Rodrigo others Benenson. Large-scale interactive object segmentation with human annotators. In CVPR, 2019. 6
- [5] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021. 3, 7, 8, 13
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 3
- [7] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. arXiv preprint arXiv:2003.04297, 2020. 7, 8, 13
- [8] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In CVPR, 2021. 3, 8
- [9] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *ICCV*, 2021. 8
- [10] Zhiwei Chen, Changan Wang, Yabiao Wang, Guannan Jiang, Yunhang Shen, Ying Tai, Chengjie Wang, Wei Zhang, and Liujuan Cao. Lctr: On awakening the local continuity of transformer for weakly supervised object localization. In AAAI, 2022. 2, 6, 7
- [11] 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 *CVPR*, 2015. 2
- [12] Junsuk Choe, Seong Joon Oh, Seungho Lee, Sanghyuk Chun, Zeynep Akata, and Hyunjung Shim. Evaluating weakly supervised object localization methods right. In *CVPR*, 2020. 12
- [13] Junsuk Choe and Hyunjung Shim. Attention-based dropout layer for weakly supervised object localization. In *CVPR*, 2019. 1, 2, 13
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020. 1, 3, 5, 7
- [15] Wei Gao, Fang Wan, Xingjia Pan, Zhiliang Peng, Qi Tian, Zhenjun Han, Bolei Zhou, and Qixiang Ye. Ts-cam: Token semantic coupled attention map for weakly supervised object localization. In *ICCV*, 2021. 2, 6, 7

- [16] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. arXiv preprint arXiv:1803.07728, 2018. 3
- [17] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. In *NeurIPS*, 2020. 3, 8
- [18] Saurav Gupta, Sourav Lakhotia, Abhay Rawat, and Rahul Tallamraju. Vitol: Vision transformer for weakly supervised object localization. In *CVPR*, 2022. 2, 6, 7, 8, 13
- [19] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In CVPR, 2020. 3
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In CVPR, 2016. 1, 5, 7, 12
- [21] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In CVPR, 2017. 12
- [22] Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Fei-Fei Li. Novel dataset for fine-grained image categorization: Stanford dogs. In *FGVC*, 2011. 6
- [23] Minsong Ki, Youngjung Uh, Junsuk Choe, and Hyeran Byun. Contrastive attention maps for self-supervised colocalization. In *ICCV*, 2021. 1, 2, 4, 6, 7
- [24] Eunji Kim, Siwon Kim, Jungbeom Lee, Hyunwoo Kim, and Sungroh Yoon. Bridging the gap between classification and localization for weakly supervised object localization. In *CVPR*, 2022. 6, 7, 13
- [25] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *ICCV workshops*, 2013. 6
- [26] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, Univ. of Toronto, 2009. 8
- [27] Krishna Kumar Singh and Yong Jae Lee. Hide-and-seek: Forcing a network to be meticulous for weakly-supervised object and action localization. In *ICCV*, 2017. 1, 2, 13
- [28] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Colorization as a proxy task for visual understanding. In CVPR, 2017. 3
- [29] Huifang Li, Yidong Li, Yi Jin, and Tao Wang. Object representation enhancement for self-supervised colocalization. *IJIS*, 37:8277–8290, 2022. 1, 2, 4, 6, 7
- [30] Weihao Li, Omid Hosseini Jafari, and Carsten Rother. Deep object co-segmentation. In ACCV, 2018. 2
- [31] Tie Liu, Zejian Yuan, Jian Sun, Jingdong Wang, Nanning Zheng, Xiaoou Tang, and Heung-Yeung Shum. Learning to detect a salient object. *IEEE TPAMI*, 33(2):353–367, 2010. 6, 8
- [32] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
 6
- [33] Meng Meng, Tianzhu Zhang, Qi Tian, Yongdong Zhang, and Feng Wu. Foreground activation maps for weakly supervised object localization. In *ICCV*, 2021. 1, 2, 6, 7

- [34] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*, 2016. 3
- [35] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015. 1
- [36] Bernhard Schölkopf, Ralf Herbrich, and Alex J Smola. A generalized representer theorem. In COLT, 2001. 3
- [37] 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. In *BMVC*, 2021. 1, 3, 6, 7
- [38] Yukun Su, Guosheng Lin, Yun Hao, Yiwen Cao, Wenjun Wang, and Qingyao Wu. Self-supervised object localization with joint graph partition. In AAAI, 2022. 1, 2, 4, 6, 7
- [39] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019. 12, 13
- [40] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. 1, 6
- [41] Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. Dense contrastive learning for self-supervised visual pre-training. In *CVPR*, 2021. 7
- [42] Yangtao Wang et al. Self-supervised transformers for unsupervised object discovery using normalized cut. In CVPR, 2022. 1, 3, 6, 7, 8
- [43] Jun Wei, Qin Wang, Zhen Li, Sheng Wang, S Kevin Zhou, and Shuguang Cui. Shallow feature matters for weakly supervised object localization. In *CVPR*, 2021. 2, 12, 13
- [44] Xiu-Shen Wei, Jian-Hao Luo, Jianxin Wu, and Zhi-Hua Zhou. Selective convolutional descriptor aggregation for fine-grained image retrieval. *IEEE TIP*, 26(6):2868–2881, 2017. 1, 2
- [45] Xiu-Shen Wei, Chen-Lin Zhang, Jianxin Wu, Chunhua Shen, and Zhi-Hua Zhou. Unsupervised object discovery and colocalization by deep descriptor transformation. *PR*, 88:113– 126, 2019. 1, 2, 6, 7
- [46] Yunchao Wei, Jiashi Feng, Xiaodan Liang, Ming-Ming Cheng, Yao Zhao, and Shuicheng Yan. Object region mining with adversarial erasing: A simple classification to semantic segmentation approach. In CVPR, 2017. 2, 13
- [47] Pingyu Wu, Wei Zhai, and Yang Cao. Background activation suppression for weakly supervised object localization. In *CVPR*, 2022. 1, 2, 6, 7
- [48] Jinheng Xie, Jianfeng Xiang, Junliang Chen, Xianxu Hou, Xiaodong Zhao, and Linlin Shen. C2am: Contrastive learning of class-agnostic activation map for weakly supervised object localization and semantic segmentation. In *CVPR*, 2022. 1, 2, 5, 6, 7, 8, 12, 13
- [49] Jilan Xu, Junlin Hou, Yuejie Zhang, Rui Feng, Rui-Wei Zhao, Tao Zhang, Xuequan Lu, and Shang Gao. Cream: Weakly supervised object localization via class re-activation mapping. In *CVPR*, 2022. 1, 2, 6, 7, 13

- [50] Chih-Kuan Yeh, Joon Kim, Ian En-Hsu Yen, and Pradeep K Ravikumar. Representer point selection for explaining deep neural networks. In *NeurIPS*, 2018. 2, 3, 5
- [51] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *ICCV*, 2019. 1
- [52] Chi Zhang, Guankai Li, Guosheng Lin, Qingyao Wu, and Rui Yao. Cyclesegnet: Object co-segmentation with cycle refinement and region correspondence. *IEEE TIP*, 30:5652– 5664, 2021. 2
- [53] Chen-Lin Zhang, Yun-Hao Cao, and Jianxin Wu. Rethinking the route towards weakly supervised object localization. In *CVPR*, 2020. 2, 5, 6, 7, 8, 12, 13
- [54] Runsheng Zhang, Yaping Huang, Mengyang Pu, Jian Zhang, Qingji Guan, Qi Zou, and Haibin Ling. Object discovery from a single unlabeled image by mining frequent itemsets with multi-scale features. *IEEE TIP*, 29:8606–8621, 2020. 1, 2, 6, 7
- [55] Xiaolin Zhang, Yunchao Wei, Jiashi Feng, Yi Yang, and Thomas S Huang. Adversarial complementary learning for weakly supervised object localization. In CVPR, 2018. 1, 2
- [56] Xiaolin Zhang, Yunchao Wei, Guoliang Kang, Yi Yang, and Thomas Huang. Self-produced guidance for weaklysupervised object localization. In *ECCV*, 2018. 1, 2
- [57] Xiaolin Zhang, Yunchao Wei, and Yi Yang. Inter-image communication for weakly supervised localization. In ECCV, 2020. 2
- [58] Xiaolin Zhang, Yunchao Wei, Yi Yang, and Fei Wu. Rethinking localization map: Towards accurate object perception with self-enhancement maps. arXiv preprint arXiv:2006.05220, 2020. 2
- [59] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *CVPR*, 2016. 1, 2, 6
- [60] Lei Zhu, Qian Chen, Lujia Jin, Yunfei You, and Yanye Lu. Bagging regional classification activation maps for weakly supervised object localization. In *ECCV*, 2022. 6, 7, 13
- [61] Lei Zhu, Qi She, Qian Chen, Yunfei You, Boyu Wang, and Yanye Lu. Weakly supervised object localization as domain adaption. In CVPR, 2022. 6, 7, 8, 13

Unsupervised Object Localization with Representer Point Selection

-Supplementary Material-

A. Further Analysis of Threshold au

In this supplementary section, we offer theoretical and empirical support to determine the threshold value, τ . As defined in Eq. (5) of our original manuscript, \mathbf{w}^* , which is used as a foreground predictor, can be rewritten with the sample global importance α from Eq. (9) as follows:

$$\mathbf{w}^* = \sum_{n=1}^{N_D} \sum_{i=1}^{HW} \alpha_{n,i} \hat{\mathbf{f}}_{n,i} = \sum_{i=1}^{N} \alpha_i \hat{\mathbf{f}}_i$$
(13)

$$= \frac{1}{C} \sum_{i=1}^{N} (\|\mathbf{f}_i\| - \tau) \mathbf{\hat{f}}_i$$
(14)

$$= \frac{1}{C} \sum_{i=1}^{N} \mathbf{f}_{i} - \frac{\tau}{C} \sum_{i=1}^{N} \hat{\mathbf{f}}_{i}, \qquad (15)$$

where \mathbf{f}_i is a feature vector, $\hat{\mathbf{f}}_i = \frac{\mathbf{f}_i}{||\mathbf{f}_i||}$, $N = N_D H W$ for simplicity and C denotes a constant. In Eq. (15), τ is a threshold that is used to determine soft pseudo labels for the training examples using the norm of the feature vector.

A straightforward approach to determine the threshold is to calculate the expected value of feature vector norms across the training set, given by $\tau = E[||\mathbf{f}_i||] = \sum_i^N ||\mathbf{f}_i||/N$. However, using a uniform probability distribution for expected value calculations does not provide information on the directions and similarities between feature vectors. Therefore, we propose a joint probability distribution that considers the correlations among feature vectors to compute the expected value. Let us denote two independent random variables, X and X', which share the sample space of feature vector norms, and XX' is a joint random variable. To utilize the relationships between all pairs of feature vectors, we employ the cosine similarity to compute the joint probability mass function of XX'. The joint probability mass function of XX' is then expressed as follows:

$$P(X = \|\mathbf{f}_i\|, X' = \|\mathbf{f}_j\|) \propto \mathbf{\hat{f}}_i^\top \mathbf{\hat{f}}_j, \tag{16}$$

and the expectation of XX' is given by

$$E[XX'] = \sum_{i=1}^{N} \sum_{j=1}^{N} \frac{\|\mathbf{f}_i\| \|\mathbf{f}_j\| P(X = \|\mathbf{f}_i\|, X' = \|\mathbf{f}_j\|)}{\sum_{i'=1}^{N} \sum_{j'=1}^{N} P(X = \|\mathbf{f}_i'\|, X' = \|\mathbf{f}_j'\|)}$$
(17)

$$=\sum_{i=1}^{N}\sum_{j=1}^{N}\|\mathbf{f}_{i}\|\|\mathbf{f}_{j}\|\frac{\mathbf{\hat{f}}_{i}^{\top}\mathbf{\hat{f}}_{j}}{\sum_{i'=1}^{N}\sum_{j'=1}^{N}\mathbf{\hat{f}}_{i'}^{\top}\mathbf{\hat{f}}_{j'}},$$
(18)

where $\hat{\mathbf{f}}_i^{\top} \hat{\mathbf{f}}_j > 0, \forall i, j$, because the last layer of pre-trained encoder $\Phi(\cdot)$ contains a ReLU operation.

Let us denote $\tau = E[X]$, and then the expected value of a jointly distributed discrete random variables of two independent random variables is given by the product of the expected values of

two random variables as follows:

$$\tau^{2} = (E[X])^{2} = E[X]E[X'] = E[XX']$$
(19)

$$=\sum_{i=1}^{N}\sum_{j=1}^{N} \|\mathbf{f}_{i}\| \|\mathbf{f}_{j}\| \frac{\mathbf{f}_{i}^{\top}\mathbf{f}_{j}}{\sum_{i'=1}^{N}\sum_{j'=1}^{N}\hat{\mathbf{f}}_{i'}^{\top}\hat{\mathbf{f}}_{j'}}$$
(20)

$$=\frac{\sum_{i=1}^{N}\sum_{j=1}^{N}\mathbf{f}_{i}^{\top}\mathbf{f}_{j}}{\sum_{i'=1}^{N}\sum_{j'=1}^{N}\hat{\mathbf{f}}_{i'}^{\top}\hat{\mathbf{f}}_{j'}},$$
(21)

where the denominator and numerator in Eq. (21) can be expressed by u and v as in Eq. (15) as follows:

$$\sum_{i=1}^{N} \sum_{j=1}^{N} \hat{\mathbf{f}}_{i}^{\top} \hat{\mathbf{f}}_{j} = \left(\sum_{i=1}^{N} \hat{\mathbf{f}}_{i}\right)^{\top} \sum_{i=1}^{N} \hat{\mathbf{f}}_{i}$$
(22)

$$= \|\sum_{i=1}^{N} \hat{\mathbf{f}}_{i}\|^{2} = \|\mathbf{u}\|^{2}, \qquad (23)$$

$$\sum_{i=1}^{N} \sum_{j=1}^{N} \mathbf{f}_{i}^{\top} \mathbf{f}_{j} = \|\sum_{i=1}^{N} \mathbf{f}_{i}\|^{2} = \|\mathbf{v}\|^{2}.$$
 (24)

Therefore, τ is computed as follows:

$$\tau = \frac{\|\mathbf{v}\|}{\|\mathbf{u}\|}.$$
(25)

B. Additional WSOL Results

Advanced Classifier In comparing our method with other weakly supervised object localization methods [43, 53, 48] that utilize more advanced classifiers, such as EfficientNetB7 [39], we also evaluate our method using EfficientNetB7 for classification in a weakly supervised setting. As shown in Table 8, our method outperforms other self-supervised methods which utilize much deeper networks, such as DenseNet161 [21], instead of ResNet50 [20], by significant margins. Furthermore, our method exhibits superior performances compared to the weakly supervised method [43] which relies on explicit class labels for training, while our method and self-supervised methods solely use pre-trained classification networks for classification.

MaxBoxAccv2 In catering to various demands for localization accuracy, [12] proposed evaluating WSOL methods through *MaxBoxAccV2*. *MaxBoxAccV2* is calculated by averaging the *MaxBoxAcc* performance across various IoU threshold $\delta \in$ {0.3, 0.5, 0.7}. As shown in Table 9, our method surpasses other self-supervised and weakly supervised methods on ImageNet. In the CUB-200-2011 dataset as well, our approach achieves performance with negligible differences from the state-of-the-art, independent of the architecture.

C. Advantage of SSL Pre-trained Backbone

We present the results of both supervised and self-supervised pre-trained models in Table 10. Interestingly, we found that

Table 8. Comparison between our method and several object localization methods that use the additional classifier, EfficientNetB7 [39], in terms of *Top-1*, *Top-5* and *GT-known Loc* on CUB-200-2011 test set and ImageNet-1K validation set. Loc. and Cls. denote the localization and classification backbones, respectively. † indicates MoCo v2 pre-trained backbone.

Method Loc		Cls		τ	(CUB-200-2011			ImageNet-1K		
u	2001		2	,	Top-1 Loc	Top-5 Loc	GT-Known	Top-1 Loc	Top-5 Loc	GT-Known	
Weakly supervise	Weakly supervised method										
SPOL '21 [43]	ResNet50	EfficientNetB7	1	1	80.12	93.44	96.46	59.14	67.15	69.02	
Self-supervised n	nethods										
PSOL ·20 [53]	DenseNet161	EfficientNetB7	X	1	80.89	89.97	91.78	58.00	65.02	66.28	
C ² AM _{'22} [48]	DenseNet161	EfficientNetB7	X	1	81.76	91.11	92.88	59.56	67.05	68.53	
w/o finetuning											
Ours	ResNet50 [†]	EfficientNetB7	X	X	84.90	94.74	96.67	60.17	67.87	69.30	

Table 9. Comparison between the proposed method and the state-of-theart weakly supervised object localization methods in terms of *MaxBox*-*AccV2* on CUB-200-2011 and ImageNet-1K.

Methods	Backbone	CUB-200-2011	ImageNet-1K
WSOL methods			
BGC '22 [24]	ResNet50	75.90	68.70
CREAM '22 [49]	ResNet50	73.50	67.40
DAOL '22 [61]	ResNet50	69.87	68.23
BagCAM ·22 [60]	ResNet50	84.88	69.97
ViTOL '22 [18]	ViT-S	73.17	70.47
SCM ₂₂ [2]	ViT-S	89.90	-
Self-supervised metho	ds		
$C^{2}AM_{22}$ [48]	ResNet50	83.80	66.80
w/o finetuning			
MoCo v2 [7] + Ours	ResNet50	87.26	66.38
DINO [5] + Ours	ViT-S	88.83	73.04

Table 10. Comparison of the performance of our method between Moco v2 and supervised pre-trained ResNet50 on UOL setup.

	Pre-trai	ining	CUB	Cars	Aircraft	Dogs	ImageNet	
-	Supervi MoCo	ised v2	88.45 96.67	96.98 99.69	98.47 98.71	89.53 95.07	63.22 66.89	
Su	pervised	MoC	:o v2	Supervised	MoCo v2	Superv	vised MoCo v	v2

Figure 7. Visualization of activation maps using the supervised and selfsupervised (MoCo v2) pre-trained models.

the results of the supervised model were inferior to those of the self-supervised model. This disparity can be linked to a wellestablished challenge in object localization with class-level supervision [13, 27, 46]. Class-level supervised models often concentrate mainly on the most discriminative parts, as they are trained to learn features that have a substantial impact on classification. In the context of our method, which does not involve fine-tuning the model, this issue becomes more pronounced. To further illustrate this phenomenon, we include examples in Fig. 7. Here, the supervised model activates only the most visually prominent features, while the self-supervised model demonstrates a more comprehensive ability to localize the entire object.

D. More Qualitative Results

We also include further qualitative results to illustrate the operating process of our method for selecting representer points, as depicted in Figure 3 of the main manuscript. As shown in Figure 8, we present examples that highlight global sample importance, the similarity between features, and representer values for given points within an image. These examples reveal that representer values tend to escalate when both the feature similarity and the importance α of each training example are pronounced. In the visualizations of representer value maps, red regions indicate excitatory points for the foreground prediction, while blue regions indicate inhibitory points. By fostering a comprehensive understanding of the model's predictions, it provides valuable insights into the reasoning behind the model's specific predictions and exclusions.

E. Limitations and Social Impacts

Our method does not rely on ground-truth annotations, which reduces the risk of bias, but it increases the likelihood of errors in object localization when compared to supervised methods. In addition, our method shares common limitations with other unsupervised, self-supervised, or weakly supervised methods, such as difficulties in detecting and recognizing rare, small, or complexly appearing objects including objects with similar textures or shapes and those set against cluttered backgrounds. Additionally, since our approach utilizes training examples, it carries a potential risk of privacy violations if the dataset is not meticulously curated. However, despite these limitations, we believe our method offers a unique advantage: it provides explainability about how it discovers objects. This ability sets our approach apart from other methods and adds to its appeal.



Figure 8. Illustration of how our method computes activation maps using two example points, yellow and green stars. Both importance maps and representer value maps are normalized to be centered at zero, unlike similarity maps' min-max normalization. Hence, red and blue regions denote positive and negative representer points, respectively. Green or yellow colored regions indicate very small absolute values of the representer value.