# DeepCut: Unsupervised Segmentation using Graph Neural Networks Clustering

Amit Aflalo , Shai Bagon , Tamar Kashti , and Yonina Eldar

Faculty of Mathematics and Computer Science, Weizmann Institute of Science

Figure 1: **DeepCut:** We use graph neural networks with unsupervised losses from classical graph theory to solve various image segmentation tasks. Specifically, we employ deep features obtained from a pre-trained vision transformer to construct a graph representation for each image, and subsequently partition this graph to generate a segmentation.

## Abstract

*Image segmentation is a fundamental task in computer vision. Data annotation for training supervised methods can be labor-intensive, motivating unsupervised methods. Current approaches often rely on extracting deep features from pre-trained networks to construct a graph, and classical clustering methods like $k$-means and normalized-cuts are then applied as a post-processing step. However, this approach reduces the high-dimensional information encoded in the features to pair-wise scalar affinities. To address this limitation, this study introduces a lightweight Graph Neural Network (GNN) to replace classical clustering methods while optimizing for the same clustering objective function. Unlike existing methods, our GNN takes both the pair-wise affinities between local image features and the raw features as input. This direct connection between the raw features and the clustering objective enables us to implicitly perform classification of the clusters between different graphs, resulting in part semantic segmentation without the need for additional post-processing steps. We demonstrate how classical clustering objectives can be formulated as self-supervised loss functions for training an image segmentation GNN. Furthermore, we employ the Correlation-Clustering (CC) objective to perform clustering without defining the number of clusters, allowing for $k$-less clustering. We apply the proposed method for object localization, segmentation, and semantic part segmentation tasks, surpassing state-of-the-art performance on multiple benchmarks[1].*

## 1. Introduction

Object localization and segmentation play crucial roles in various real-world applications, such as autonomous cars, robotics, and medical diagnosis. These tasks have been longstanding challenges in computer vision, and significant efforts have been invested in improving their accuracy.

Presently, state-of-the-art performance in these tasks is achieved using supervised Deep Neural Networks (DNNs). However, the limited availability of annotated data restricts the applicability of these methods. Data annotation is labor-intensive and costly, particularly in specialized fields like medical imaging, where domain experts are required for accurate annotations. Several solutions have been proposed to address this challenge, including leveraging color data,

---

[1]Project page: https://sampl-weizmann.github.io/DeepCut/

1

adding priors such as boundaries or scribbles [21], semi-supervised learning[3, 20], weakly-supervised learning[24], and more. However, these approaches have limitations since they still rely on some form of annotations or prior knowledge about the image structure.

An alternative approach to tackle this problem is to explore *unsupervised* methods. Recent research on unsupervised Deep Neural Networks (DNNs) has yielded promising outcomes. For instance, self-DIstillation with No labels (DINO [9]) was employed to train Vision Transformers (ViTs), and the generated attention maps corresponded to semantic segments in the input image. The deep features extracted from these trained transformers demonstrated significant semantic meaning, facilitating their utilization in various visual tasks, such as object localization, segmentation, and semantic segmentation.[2, 23, 34].

Some recent work in this area has demonstrated promising results, especially using unsupervised techniques that combine deep features with classical graph theory for object localization and segmentation tasks [23, 34]. These methods are lightweight and rely on pre-trained unsupervised networks, unlike end-to-end approaches that require significant time and resources for training from scratch.

Our proposed method called *DeepCut*, introduces an innovative approach using Graph Neural Networks (GNNs) combined with classical graph clustering objectives as a loss function. GNNs are specialized neural networks designed to process graph-structured data, and they have achieved remarkable results in various domains, including protein folding with AlphaFold [19], drug discovery [37] and traffic prediction[18]. This work employs GNNs for computer vision tasks such as object localization, segmentation, and semantic part segmentation.

A key advantage of our method is the direct utilization of deep features within the clustering process, in contrast to previous approaches [23, 34] that discarded this valuable information and relied solely on correlations between features. This approach leads to improved object segmentation performance and enables semantic segmentation across multiple graphs, each corresponding to a separate image. To further enhance segmentation, we adopt a two-stage approach that overcomes shortcomings of the objective functions overlooked by previous methods. We first separate foreground and background and then perform segmentation individually for each part, leading to more accurate results. Moreover, we introduce a method to perform "k-less clustering," allowing data clustering without the need to predefine the number of clusters $k$, enhancing flexibility and adaptability.

Our contributions are:

- Utilizing a lightweight GNN with classical clustering objectives as unsupervised loss functions for image segmentation, surpassing state-of-the-art performance at various

unsupervised segmentation tasks (speed and accuracy).

- Optimizing the correlation clustering objective for deep features clustering, which is not feasible with classical methods, thus achieving $k$-less clustering.

- Performing semantic part segmentation on multiple images using test-time optimization. By applying the method to each image separately, we eliminate the need for any post-processing steps required by previous methods.

## 2. Background

In the era prior to the deep-learning surge, numerous classical approaches adopted quantitative criteria for segmentation based on principles from graph theory. These methods involved representing affinities between image regions as a graph and associating various image partitions with cuts in that graph (e.g., [26, 5, 4]). The quality of image segments was determined by the "optimality" of these cuts. Remarkably, these techniques operated without any supervision, relying solely on the provided affinities between image regions.

Different methods have defined optimal cuts in various ways, each presenting advantages. We provide a brief overview of the relevant approaches.

### 2.1. Graph Clustering

We leverage two graph clustering functionals from classical graph theory in our work, *normalized cut*[26] and *correlation clustering*[5].

**Notations** Let $G = (V, E)$ be an undirected graph induced by an image. Each node represents an image region, and the weights $w_{ij}$ represent the affinity between image regions $i$ and $j$, $i, j = 1 \ldots n$. Let $W$ be an $n \times n$ matrix whose entries are $w_{ij}$.

Our goal is to partition this graph into $k$ disjoint sets $A_1, A_2 ... A_k$ such that $\cup_i A_i = \mathcal{V}$ and $\forall_{j \neq i} A_i \cap A_j = \emptyset$. This partition can be expressed as a binary matrix $S \in \{0, 1\}^{n \times k}$ where $S_{ic} = 1$ iff $i \in A_c$.

**Normalized Cut (N-cut) [26]** A good partition is defined as one that maximizes the number of within-group connections, and minimizes the number of between-group connections. The number of between-group connections can be computed as the total weight of edges removed and described in graph theory as a *cut*:

$$cut(A, B) = \sum_{u \in A, v \in B} w(u, v). \tag{1}$$

Where $A, B$ are parts of two-way partition of $G$. This objective is formulated by the *normalized cut* (N-cut) functional:

$$Ncut(A, B) = \frac{cut(A, B)}{assoc(A, \mathcal{V})} + \frac{cut(A, B)}{assoc(B, \mathcal{V})}, \quad (2)$$

where $assoc(A, \mathcal{V}) = \sum_{i \in A, j \in \mathcal{V}} w_{ij}$ is the total affinities connecting nodes of $A$ to all nodes in the graph. The N-cut formulation can be easily extended to $K > 2$ segments. Shi and Malik [26] also suggested an approximated solution to Eq. (2) using spectral methods, known as *spectral clustering*, where the spectrum (eigenvalues) of a graph Laplacian matrix is used to get approximated solution to the N-cut problem.

**Correlation Clustering (CC) [5]** When the graph, derived from an image, contains *both negative and positive* affinities, N-cut is no longer applicable. In this case, a good image segment may be one that maximizes the *positive* affinities inside the segment and the *negative* ones across segments [5]. This objective can be formulated by the *correlation clustering* (CC) functional. Given a matrix $W$, an optimal partition $U$ minimizes:

$$CC(S) = -\sum_{ij} W_{ij} \sum_c S_{ic} S_{jc}. \quad (3)$$

Correlation clustering utilizes intra-cluster disagreement (repulsion) to *automatically* deduce the number of clusters $k$ [5] so that it does not need to know $k$ in advance. Further theoretical analysis on this property can be found in [4], along with classical optimization algorithms applying CC to image segmentation.

## 2.2. Graph Neural Networks

GNN (Graph Neural Network) is a category of neural networks designed to process graph-structured data directly. A GNN layer comprises two fundamental operations: **message passing** and **aggregation**. **Message passing** involves gathering information from the neighbors of each node and is applied to all nodes in the graph. The exchanged information can include a combination of node features, edge features, or any other data embedded in the graph. **Aggregation** refers to fusing all the messages obtained during the message passing phase into one message that updates the current node's state.

In our approach, we utilize a specific type of GNN called the Graph Convolutional Network (GCN) [35]. In a GCN layer $\ell$, each node $h_v$ is processed as follows:

$$h_v^{(\ell+1)} = \sum_{u \in \mathcal{N}(v)} \Theta \frac{h_u^{(\ell)}}{|\mathcal{N}(v)|}, \quad (4)$$

where $\Theta$ is a matrix of trainable parameters. $\mathcal{N}(v)$ denotes the neighbours of a node $h_v$, and $|\mathcal{N}(v)|$ is the number of neighbours. The GCN aggregation is performed through a summation operator, with the term inside the summation representing the message-passing operation. The objective during training is to optimize the network parameters $\Theta$ in a way that generates meaningful messages to be passed among nodes, optimizing the loss function.

## 3. Method

In our approach, we process each image through the pretrained network to extract deep features, which are then used for clustering. Subsequently, we construct a weighted graph based on these features. Employing a lightweight Graph Neural Network (GNN), we optimize our unsupervised graph partitioning loss functions (either $\mathcal{L}NCut$ or $\mathcal{L}CC$) separately for each image graph.

We adopt this methodology to achieve unsupervised object localization, segmentation, and semantic part segmentation.

### 3.1. From Deep Features to Graphs

As shown in Figure 2, given an image $M$ of size $m \times n$ and $d$ channels, we pass it through a transformer T. The transformer divides the image into $\frac{mn}{p^2}$ patches, where $p$ is the patch size of transformer T. To extract the transformer's internal representation for each patch, we utilize the *key* token from the last layer, as it has demonstrated superior performance across various tasks [2, 23]. The output is a feature vector $f_{(\frac{mn}{p^2} \times c)}$, where $c$ represents the token embedding dimension, containing all the extracted features from the different patches.

Consider a weighted graph $G = (V, E)$ with $W$ as the weight matrix. We construct a patch-wise correlation matrix from the features obtained by the Vision Transformer:

$$W = ff^T \in \mathbb{R}^{\frac{mn}{p^2} \times \frac{mn}{p^2}}. \quad (5)$$

In correlation clustering, the negative weights (representing reputations) carry valuable information utilized by the clustering objective. However, the normalized cut objective only accepts positive weights. As a result, we threshold at zero:

$$W = ff^T \cdot (ff^T > 0) \in \mathbb{R}^{\frac{mn}{p^2} \times \frac{mn}{p^2}}. \quad (6)$$

We introduce a hyper-parameter called *k sensitivity* denoted by $\alpha$ to adapt the cluster choosing process in correlation clustering. Since the number of clusters cannot be directly chosen in correlation clustering, this parameter allows us to control the sensitivity of the process, where a higher value of $\alpha$ corresponds to more clusters. We enforce this adjustment in the following manner:

Figure 2: **Method overview:** After extracting deep features from a pretrained ViT model, we construct a similarity matrix based on the patch-wise feature similarities, which becomes our adjacency matrix. We build a graph using this adjacency matrix and the deep features as node features. Next, we train a lightweight GNN using unsupervised graph partitioning loss functions (Sec. 2.1) to partition the graph into $k$ distinct clusters, which can be used for various downstream tasks.

$$W = ff^T - \frac{max(ff^T)}{\alpha}. \tag{7}$$

where $\alpha \in [1, \inf)$. Lower $\alpha$ value corresponds to higher repulsion forces between nodes (negative weights in $W$) and thus higher cluster count, as correlation clustering maximizes negative affinities between segments in the graph.

### 3.2. Graph Neural Network Clustering

Let $\hat{N}$ be a node feature matrix obtained by applying one or more layers of GNN convolution on a graph $G$ with an adjacency matrix $W$. In our case, we use a one-layer GCN and construct the graph using the patch-wise correlation matrix from ViT obtained features (Eq. (5), Eq. (6), Eq. (7)). Let $\mathbf{S}$ be the output of a Multi-Layer Perception (MLP) with a softmax function applied on $\hat{N}$:

$$
\begin{aligned}
\hat{N} &= GNN(N, W; \Theta_{GNN}), \\
S &= MLP(\hat{N}; \Theta_{MLP}),
\end{aligned}
\tag{8}
$$

where $\Theta_{MLP}$ and $\Theta_{GNN}$ are trainable parameters. The GNN output $S$, is the cluster assignment matrix containing vectors representing the node's probability of belonging to a particular cluster.

The GNN is optimized using either the normalized-cut relaxation proposed in [8] or our newly proposed method with the correlation clustering objective as the loss function.

The loss function in the case of normalized cut is:

$$\mathcal{L}_{NCuts} = \frac{Tr(S^T W S)}{Tr(S^T D S)} + \left\| \frac{S^T S}{\|S^T S\|_F} - \frac{\mathbb{I}_K}{\sqrt{K}} \right\|_F, \tag{9}$$

where $D = diag(\sum_j W_{i,j})$ is the row-wise sum diagonal matrix of $W$. $K$ denotes the number of disjoint sets we aim to partition the graph into, and $\mathbb{I}_K$ is the identity matrix. The

first term of the objective function promotes the clustering of strongly connected components together, while the second term encourages the cluster assignments to be orthogonal and have similar sizes.

The loss function for correlation clustering [5] is:

$$\mathcal{L}_{CC} = -Tr(WSS^T). \tag{10}$$

This term promotes intra-cluster agreement while encouraging repulsion (negative affinities) between clusters.

$W$ is defined as Eq. (6) and Eq. (7) for the N-cut and CC loss, respectively. To control the sensitivity of the clustering process for CC, we utilize a k-sensitivity value, as explained in Sec. 7.

### 3.3. Graph Neural Network Segmentation

In Sec. 3.1, we propose constructing a graph from deep features extracted from unsupervised trained ViTs, where each node represents a patch of the original image. The nodes in the graph are then clustered into disjoint sets, representing different image segments. For this clustering process, we employ graph neural network clustering as described in Sec. 3.2, utilizing either the CC or N-cut losses. However, unlike previous approaches that used the N-cut loss defined strictly for positive weights [8, 23], we advocate using the correlation clustering functional as a loss. Correlation clustering enables us to utilize negative weights for graph building (see Eq. (5) and Eq. (7)), facilitating clustering without predefining the number of clusters.

Additionally, we incorporate deep features as node features for graph building, which is absent in previous methods that solely used the correlation matrix of the deep features [8, 23, 34]. As a result, our method allows for feature classification while clustering and implicitly facilitates semantic part segmentation without necessitating postprocessing steps, as seen in previous methods [23].

Figure 3: **Proposed two-stage clustering:** First, we cluster the image into two disjoined sets, then apply clustering to the foreground and background separately.

| Method | VOC-07 | VOC-12 | COCO-20k |
|---|---|---|---|
| Selective Search[28] | 18.8 | 20.9 | 6.0 |
| EdgeBoxes[38] | 31.1 | 31.6 | 28.8 |
| DINO-[CLS][9] | 45.8 | 46.2 | 42.1 |
| LOST[27] | 61.9 | 64.0 | 50.7 |
| Spectral Methods[23] | 62.7 | 66.4 | 52.2 |
| TokenCut[34] | 68.8 | 72.1 | 58.8 |
| DeepCut: CC loss | 68.8 | 67.9 | 57.6 |
| DeepCut: N-cut loss | **69.8** | **72.2** | **61.6** |

Table 1: **Object localization results.** CorLoc metric (percentage of images with $IOU > 0.5$). CC and N-cut denotes correlation clustering and Normalized Cut respectively.

Our approach is versatile and applicable to various image partitioning-related tasks, including:

**Object Localization**   Object localization involves identifying the primary object in an image and enclosing it with a bounding box. To perform localization, we follow the steps below: (1) Use our GCN clustering method with $k = 2$. (2) Examine the edges of the clustered image and identify the cluster that appears on more than two edges as the background, while the other cluster becomes our main object. (3) Apply a bounding box around the identified main object.

**Object Segmentation**   Object segmentation involves the separation of the foreground object in an image, commonly known as foreground-background segmentation. The method employed for object segmentation is identical to object localization, with the inclusion of stage (3), where the bounding box is applied.

**Semantic Part Segmentation**   DeepCut achieves semantic part segmentation through a test-time optimization paradigm, where the model is sequentially exposed to each image, optimizing the model weights based on the previous image. This means the model does not require training on all images beforehand, eliminating the need for co-segmentation or post-processing steps.

This advantage is derived from our approach to learning deep features with GNN, which differs from previous methods [8, 34] that solely rely on correlations and discard the high-dimensional data of deep features. By leveraging the intricate semantic information embedded within deep features, our model implicitly performs the classification of the clusters between different image graphs, enabling semantic part segmentation without the need for explicit post-processing or additional training steps.

Our segmentation process involves two steps, as shown in Figure 3: foreground-background segmentation (k=2) followed by semantic part segmentation on the foreground object (k=4). This two-stage process addresses the bias of the clustering functions towards larger clusters (e.g., background-foreground), which limits the level of detail in foreground object segmentation. The exact process can also be applied to improve background segmentation, as depicted in Figure 1 and Figure 5.

## 4. Training and performance

For all experiments, We use DINO [9] trained ViT-S/8 transformer for feature extraction, with pre-trained weights from the DINO paper authors (trained on ImageNet[25]). **No training is conducted on the tested datasets.** We employ a test-time training paradigm for all experiments, where each image is trained using a proposed graph neural network (GNN) segmentation approach for ten epochs. Since the proposed losses are unsupervised and involve solving an optimization problem, generalizing the model to the entire dataset does not improve accuracy. This training method achieves a performance (speed) that is **more than two times** that of the current state-of-the-art TokenCut[34] on the same hardware. DeepCut efficiency stems from its lightweight architecture, consisting of only 30k trainable parameters. Implementation details and performance analyses are provided in the supplementary material.

## 5. Results

Our method is evaluated on three unsupervised tasks: single object localization, single object segmentation, and semantic part segmentation. We compare our approach with other unsupervised methods published on these tasks using widely used benchmarks. The results are presented for both the normalized-cut and correlation clustering GNN objectives.

| Method | CUB | DUTS | ECSSD |
|---|---|---|---|
| OneGAN[7] | 55.5 | - | - |
| Voynov et al.[31] | 68.3 | 49.8 | - |
| Spectral Methods[23] | 76.9 | 51.4 | 73.3 |
| TokenCut[34] | - | 57.6 | 71.2 |
| DeepCut: CC loss | 77.7 | 56.0 | 73.4 |
| DeepCut: N-cut loss | **78.2** | **59.5** | **74.6** |

Table 2: **Single object segmentation results.** mIOU (mean intersection-over-union) CC denotes correlation clustering and N-cut Normalized Cut.

## 5.1. Object Localization

In Table 1, We evaluate our unsupervised object localization performance on three datasets: PASCAL VOC 2007 [13], PASCAL VOC 2012 [14], and COCO20K (20k images chosen from MS-COCO dataset[22] introduced in previous work[29]). We compare our unsupervised approach to state-of-the-art unsupervised methods. We report our results in the Correct Localization metric (CorLoc), defined as the percentage of images whose intersection-over-union with the ground truth label is greater than 50%. DeepCut with N-cut losses gives the best results on all data sets, and DeepCut with correlation clustering loss surpasses all previous methods except TokenCut[34].

## 5.2. Single Object Segmentation

In Table 2, We evaluate our unsupervised single object segmentation performance on three datasets: CUB (widely-used dataset of birds for fine-grained visual categorization task)[32], DUTS (the largest saliency detection benchmark)[33] and ECSSD (Extended Complex Scene Saliency Dataset)[36]. We report our results in mean intersection-over-union (mIoU). DeepCut with N-cut losses archives the best results on all data sets. Visual comparison of segmentation with our two losses is presented in Figure 4; note how the CC loss segments the foreground better than the Ncut loss when background for specific images that include different objects and shadows. As this is not the case in most images, the N-cut usually outperforms the CC loss.

## 5.3. Semantic Part Segmentation

In Table 3, We evaluate our approach on the CUB dataset [32] and report results using the Normalized Mutual Info score (NMI) and Adjusted Rand Index score (ARI) on the entire test set. A comparison between our technique and deep spectral method that uses classical graph theory for deep features-based segmentation[8] is presented in Figure 6. As seen at Table 3, DeepCut with normalized-cut loss surpasses all other methods, including the top three methods[17, 16, 11] that uses ground-truth foreground

| Method | NMI | ARI |
|---|---|---|
| SCOPS[17] (model) | 24.4 | 7.1 |
| Huang and Li[16] | 26.1 | 13.2 |
| Choudhury et al.[11] | 43.5 | 19.6 |
| DFF[12] | 25.9 | 12.4 |
| Amir et al.[12] | 38.9 | 16.1 |
| DeepCut: N-cut loss | **43.9** | **20.2** |

Table 3: **Semantic part segmentation results.** ARI and NMI over the entire CUB-200 dataset are used to evaluate cluster quality. Predictions were performed with $k = 4$ with our method using the N-cut objective. First three methods use ground truth foreground masks as supervision.



Figure 4: **DeepCut segmentation: N-cut vs. CC losses.** Top: original image, middle: DeepCut with N-cut loss, bottom: DeepCut with correlation clustering loss. Note that for images with complex background DeepCut with CC loss outperforms DeepCut with N-cut loss.



Figure 5: **Semantic part segmentation: N-cut vs. CC losses.** Top: original image, middle: normalized-cut loss with our proposed two-step clustering; bottom: correlation clustering with one-step $k$-less clustering.

Figure 6: **Semantic part segmentation. deep spectral method vs. DeepCut with N-cut loss.** Top: original image; middle: deep spectral method[23]; bottom: our segmentation with N-cut loss. The deep spectral method failed to preform semantic segmentation across all three images.

masks as supervision.

| Pretraining | DeepCut N-cut | DeepCut CC | Negative percentage |
|---|---|---|---|
| DINO[9] | 59.4 | 59.8 | 33.6 |
| MoCo-v3[10] | 62.3 | 43.1 | 20 |
| MAE [15] | 47.2 | 31.2 | 5.4 |

Table 4: **Pretraining.** Object-localization performance on PASCAL VOC07[13] with different unsupervised training methods for acquiring deep features. All methods use ViT-B-16 architecture, and were trained on ImageNet[25]. We provide mIOU (mean Intersection-Over-Union) for correlation clustering (CC) and normalized-cut (N-cut) objective. We also provide the mean percentage (across all the dataset) of negative weights in the corresponding affinity matrix. Note the correlation: the higher the percentage of negative weights of the transformer, the better the mIOU of the DeepCut with CC loss; for DINO DeepCut CC loss outperforms N-cut loss for this dataset.

## 6. Deep Features Selection

Our work and previous unsupervised segmentation methods heavily rely on the correlation between deep features from different patches. However, in this study, we proposed a novel approach by utilizing the negative correlations (correlation clustering loss) that were discarded in previous works. By doing so, we leverage all available information to enhance performance in unsupervised segmentation. It is crucial to recognize that the performance of unsupervised segmentation methods that rely on deep features depends on the quality of these features.

Furthermore, no universal solution fits all scenarios, as different methods require specific types of information to achieve optimal results. For instance, our paper presents two techniques—one based on correlation clustering and the other on normalized cut. Correlation clustering depends on both positive and negative correlations between features, whereas normalized cut only accepts positive correlations.

In Table 4, we compare three unsupervised ViT training approaches: DINO[9], MoCo[10], and MAE[15]. We observe that correlation clustering performs best with a method that supplies deep features with the largest amount of negative correlation between features (DINO 33.6%). We also observe that even a small amount of negative correlation can be counterproductive, even with normalized cut (MAE 5.4%) that relies solely on positive correlation.

Considering these factors, it becomes essential to identify the best combinations of unsupervised training methods and segmentation techniques. This insight enables us to improve unsupervised training methods, allowing us to extract more valuable information for segmentation purposes.

## 7. $k$-less Clustering

In this section, we explore clustering, using image classification as an example. Traditional clustering methods usually require a predefined number of clusters (or classes) $k$ for classifying data, which can be a disadvantage as it demands prior knowledge about the data. Our DeepCut method introduces a GNN model with correlation clustering loss, enabling $k$-less clustering, where the number of clusters is derived from the data. To demonstrate this valuable property, we use the Fashion Product Images Dataset[1], sampling a subset of 500 images from 5 classes: Top-wear, Shoes, Bags, Eye-wear, and Belts.

We conducted three classification experiments, each repeated ten times: one with three classes, one with four, and one with all classes. We employ our GNN approach with correlation clustering loss in each experiment to cluster the images into different groups. We set the parameter $\alpha$ ($k$ *sensitivity* defined in Eq. (7)) to 3, using class tokens extracted from the unsupervised trained DINO[9] ViT base transformer with a patch size of 16 as our features. We report the clustering result as classification purity in Tab. 5. To validate the algorithm's robustness, we apply class permutation and random subset sampling of each class.

The versatility of correlation clustering's $k$-less nature extends to segmentation tasks, as illustrated in Figure 5. We compare our proposed DeepCut method with two baselines: connected component and spectral clustering. For spectral clustering, we use eigen-gap[30] to select the num-

| Number of classes | 3 classes | 4 classes | 5 classes |
|---|---|---|---|
| Connected Components | 33.3 ±0 | 25.0 ±0 | 20.0 ±0 |
| Spectral Clustering | 85.5 ±14.1 | 77.9 ±6.6 | 82.8 ±6.9 |
| DeepCut: cc loss | 98.3 ±0.9 | 97.1 ±1.7 | 99.27 ±0.6 |

Table 5: $k$-**less clustering.** We apply DeepCut with correlation clustering loss on the Fashion Product Images Dataset[1], using a subset of images from 5 classes: Top-wear, Shoes, Bags, Eye-wear, and Belts. We employ our $k$-less method and report the results as classification purity. The experiments are conducted with *k sensitivity = 3* as defined in Eq. (7). The presented results are the mean and standard deviation (*std*) obtained from multiple experiments.



Figure 7: **k-sensitivity.** The CC segmentation results demonstrate that similar objects, such as people and bikes, cluster together consistently across different $\alpha$ values. Smaller $\alpha$ values lead to a finer partition of objects, but even then, CC still assigns similar objects to the same clusters.

ber of clusters. The connected components method proves ineffective, clustering all images into the same group for all experiments, resulting in zero standard deviation (STD). On the other hand, spectral clustering performs better but exhibits high variance. In contrast, our DeepCut with correlation clustering loss achieves the highest accuracy and lowest variance. This powerful $k$-less property of correlation clustering is demonstrated through an image clustering example in terms of quantitative results. Furthermore, it can be effectively applied to image segmentation, as depicted in Fig. 7.

**Choosing $\alpha$:** The CC functional determines the number of clusters by considering repulsion forces between clusters, as explained in Sec. 3. These forces are calculated based on the amount of negative affinities in the adjacency matrix defined at Eq. (7). The affinities can vary when using differently trained ViT models (as observed in Tab. 4) or different datasets. To address this issue, we propose introducing a k-sensitivity variable $\alpha$ (as defined in Eq. (7)) to artificially control the amount of repulsion forces and thereby regulate the number of clusters. The choice of $\alpha$ should be tailored to the specific task at hand. As shown in Fig. 7, different $\alpha$ values yield different solutions for various tasks.

## 8. Conclusion

The study presents DeepCut, an innovative unsupervised segmentation technique that combines classical graph the-ory, Graph Neural Networks (GNNs), and self-supervised pre-trained networks. DeepCut's effectiveness is demonstrated by achieving superior performance in object localization, object segmentation, and semantic part segmentation tasks, surpassing existing state-of-the-art methods in accuracy and speed.

The proposed graph structure incorporates node features, allowing more information to be utilized during the clustering process. Consequently, implicit semantic part segmentation is achievable, eliminating the need for post-processing methods. The versatility of this methodology extends to various downstream tasks, such as video segmentation and image matting. Additionally, the study emphasizes the importance of selecting an appropriate combination of unsupervised training method and clustering method, as discussed in Section Sec. 6.

We demonstrate that our method is capable of optimizing various loss functions derived from classical graph theory, including those that are challenging to optimize using conventional tools (e.g. correlation clustering).

## References

[1] Param aggarwal. Fashion Product image dataset. https://www.kaggle.com/datasets/paramaggarwal/fashion-product-images-dataset. Accessed: 2020-27-10. 7, 8

[2] Shir Amir, Yossi Gandelsman, Shai Bagon, and Tali Dekel. Deep vit features as dense visual descriptors. *arXiv preprint arXiv:2112.05814*, 2021. 2, 3

[3] Mahmoud Assran, Mathilde Caron, Ishan Misra, Piotr Bojanowski, Armand Joulin, Nicolas Ballas, and Michael Rabbat. Semi-supervised learning of visual features by non-parametrically predicting view assignments with support samples. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8443–8452, 2021. 2

[4] Shai Bagon and Meirav Galun. Large scale correlation clustering optimization. *arXiv preprint arXiv:1112.2903*, 2011. 2, 3

[5] Nikhil Bansal, Avrim Blum, and Shuchi Chawla. Correlation clustering. *Machine learning*, 56(1):89–113, 2004. 2, 3, 4

[6] Jonathan T Barron and Ben Poole. The fast bilateral solver. In *European conference on computer vision*, pages 617–632. Springer, 2016. 11

[7] Yaniv Benny and Lior Wolf. Onegan: Simultaneous unsupervised learning of conditional image generation, foreground segmentation, and fine-grained clustering. In *European Conference on Computer Vision*, pages 514–530. Springer, 2020. 6

[8] Filippo Maria Bianchi, Daniele Grattarola, and Cesare Alippi. Spectral clustering with graph neural networks for graph pooling. In *International Conference on Machine Learning*, pages 874–883. PMLR, 2020. 4, 5, 6

[9] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021. 2, 5, 7, 11

[10] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9640–9649, 2021. 7

[11] Subhabrata Choudhury, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. Unsupervised part discovery from contrastive reconstruction. *Advances in Neural Information Processing Systems*, 34:28104–28118, 2021. 6

[12] Edo Collins, Radhakrishna Achanta, and Sabine Susstrunk. Deep feature factorization for concept discovery. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 336–352, 2018. 6

[13] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results. http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html. 6, 7

[14] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html. 6

[15] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 7

[16] Zixuan Huang and Yin Li. Interpretable and accurate fine-grained recognition via region grouping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8662–8672, 2020. 6

[17] Wei-Chih Hung, Varun Jampani, Sifei Liu, Pavlo Molchanov, Ming-Hsuan Yang, and Jan Kautz. Scops: Self-supervised co-part segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 869–878, 2019. 6

[18] Weiwei Jiang and Jiayun Luo. Graph neural network for traffic forecasting: A survey. *Expert Systems with Applications*, page 117921, 2022. 2

[19] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021. 2

[20] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023. 2

[21] Di Lin, Jifeng Dai, Jiaya Jia, Kaiming He, and Jian Sun. Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3159–3167, 2016. 2

[22] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 6

[23] Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. Deep spectral methods: A surprisingly strong baseline for unsupervised semantic segmentation and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8364–8375, 2022. 2, 3, 4, 5, 6, 7

[24] Zhongzheng Ren, Zhiding Yu, Xiaodong Yang, Ming-Yu Liu, Yong Jae Lee, Alexander G Schwing, and Jan Kautz. Instance-aware, context-focused, and memory-efficient weakly supervised object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10598–10607, 2020. 2

[25] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 5, 7, 11

[26] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on pattern analysis and machine intelligence*, 22(8):888–905, 2000. 2, 3

[27] 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-British Machine Vision Conference*, 2021. 5

[28] Jasper RR Uijlings, Koen EA Van De Sande, Theo Gevers, and Arnold WM Smeulders. Selective search for ob-

ject recognition. *International journal of computer vision*, 104(2):154–171, 2013. 5

[29] Huy V Vo, Patrick Pérez, and Jean Ponce. Toward unsupervised, multi-object discovery in large-scale image collections. In *European Conference on Computer Vision*, pages 779–795. Springer, 2020. 6

[30] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007. 7

[31] Andrey Voynov, Stanislav Morozov, and Artem Babenko. Object segmentation without labels with large-scale generative models. In *International Conference on Machine Learning*, pages 10596–10606. PMLR, 2021. 6

[32] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. Cub-dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. 6

[33] Lijun Wang, Huchuan Lu, Yifan Wang, Mengyang Feng, Dong Wang, Baocai Yin, and Xiang Ruan. Learning to detect salient objects with image-level supervision. In *CVPR*, 2017. 6

[34] Yangtao Wang, Xi Shen, Shell Xu Hu, Yuan Yuan, James L Crowley, and Dominique Vaufreydaz. Self-supervised transformers for unsupervised object discovery using normalized cut. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14543–14553, 2022. 2, 4, 5, 6

[35] Max Welling and Thomas N Kipf. Semi-supervised classification with graph convolutional networks. In *J. International Conference on Learning Representations (ICLR 2017)*, 2016. 3, 11

[36] Qiong Yan, Li Xu, Jianping Shi, and Jiaya Jia. Hierarchical saliency detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1155–1162, 2013. 6

[37] Jiaxuan You, Bowen Liu, Zhitao Ying, Vijay Pande, and Jure Leskovec. Graph convolutional policy network for goal-directed molecular graph generation. *Advances in neural information processing systems*, 31, 2018. 2

[38] C Lawrence Zitnick and Piotr Dollár. Edge boxes: Locating object proposals from edges. In *European conference on computer vision*, pages 391–405. Springer, 2014. 5

# Supplementary Material



Figure 8: **Segmentation refinement.** Single object segmentation with our method and bilateral solver[6] as a complementary step to refine the obtained segmentation.



Figure 9: **Segmentation refinement.** Resolution manipulation using smaller size stride. ViT image input resolution is $280 \times 280$ for both images.

## 1. Resolution Manipulation

In this work, we utilize deep features extracted from ViTs. Those features represent the image patches corresponding to the ViT patch size $p \times p$. We perform our segmentation patch-wise thus, for image size $(m \times n)$ our segmentation resolution is $(\frac{m}{p} \times \frac{n}{p})$. For higher resolution segmentation, we change the ViT stride value to $\frac{p}{2}$ instead of $p$ as it doubles the number of patches the ViT uses, doubling the resolution of our segmentation to $(\frac{2m}{p} \times \frac{2n}{p})$. We found this method to yield better than changing input image resolutions as the transformer in use in this paper wasn't trained on high resolution images. Example at Fig. 9

In order to improve the segmentation resolution further, a bilateral solver[6] can be added as a complementary step to refine the boundaries of the obtained segmentation: see Figure 8. All reported results int the paper is without any post-processing methods.

## 2. Implementation details

For all experiments, we use DINO [9] trained ViT-S/8 transformer for feature extraction; specifically, we use the *keys* features from the last layer of the DINO trained "student" transformer. We use pre-trained weights provided by the DINO paper authors (trained on ImageNet[25]). **We do not conduct any training of the transformer on any of the tested datasets.** Input images resized to a resolution of

$280 \times 280$, images resized using *Lanczos* interpolation. All expirements where conducted using *Tesla V100* GPU.

---

**Algorithm 1** DeepCut
---
1: $x \leftarrow Input\ image$
2: $x \leftarrow ViT(x)$               ▷ Deep features from ViT
3: $G \leftarrow Build\ graph(x)$
4: **for** Each training epoch **do**
5:      $s \leftarrow GCN(G)$         ▷ Single layer of GCN
6:      $s \leftarrow ELU(s)$
7:      $s \leftarrow MLP(s)$          ▷ Two layer MLP
8:      $s \leftarrow softmax(s)$
9:      $Loss \leftarrow \mathcal{L}NCut$ or $\mathcal{L}CC$
10: **end for**
11: $Output\ segmentation \leftarrow argmax(s)$

---

**ViT**    Evaluation mode, frozen weights.

**Build graph**    Create a graph from deep features as depicted at Sec. 3.1.

**GCN**    Graph Convolutional Network[35]. *Input size* = Deep features size, *hidden size* = 64.

**ELU**    The Exponential Linear Unit activation function.

**MLP**    Consists from 2 linear layers, **layer 1:** from GCN *hidden size* to $\frac{hidden\ size}{2}$. **layer 2:** from $\frac{hidden\ size}{2}$ to $k$ the number of desired clusters. For k-less usage with correlation clustering, the output will be set to a maximum of desired clusters. Between the layers, there is an *elu* activation function and 0.25 dropout.

**Loss**    We suggest two loss function derived from classical graph theory; NCut and CC.

**Output segmentation**    At step 8, in order to obtain the final segmentation, the vector $s$ is extracted. Each entry in this vector corresponds to a patch of the image and contains a probability vector that describes the likelihood of the patch belonging to a specific cluster. We chose the most likely cluster assignment for each patch and than unflatten the result to get an segmentation map.

## 2.1. Two-stage segmentation

The clustering functionals in this paper exhibit are biased towards larger clusters (e.g background-foreground), resulting in a tendency to underperform on finer details by merging them together, or in some cases, failing and introducing a significant amount of noise to the segmentation process. To address this issue, a solution is proposed by utilizing a two-stage segmentation approach, where the background and foreground segments are separately applied to avoid the aforementioned biases and limitations. Example can be seen at Fig. 12.

## 2.2. Training

To optimize object localization and object segmentation task, we perform individual optimization for each image for a duration of 10 epochs, with model weights being reset between images. For part semantic segmentation, we carry out separate optimization for each image over a span of 100 epochs, without resetting the model weights between them.

## 2.3. Performance

All of the results presented in the paper demonstrate DeepCut without the utilization of any post-processing (e.g. bilateral solver). All experiments were conducted using the same hardware: Tesla V100 GPU and an Intel Xeon 32 core CPU.

| Model | DUTS [mIoU] | ECSSD [mIoU] | Throughput [img/sec] |
|---|---|---|---|
| TokenCut + Bilateral Solver | 62.4 | 77.2 | 0.5 |
| TokenCut w/o Bilateral Sol. | 57.6 | 71.2 | 1 |
| Ours | 59.5 | 74.6 | 5 |

Figure 10: **Method example:** Object localization using DeepCut(NCut).



Figure 11: **Method example:** Random foreground-background segmentation samples using DeepCut(NCut) on VOC07.

Figure 12: **Method example:** Two-stage segmentation using DeepCut(NCut).



Figure 13: **Method example:** Random foreground-background segmentation samples using DeepCut(NCut/CC) on CUB-200. DeepCut segments the birds accurately without including other objects such as branches and leaves (which is a common failure point of previous methods).