# Deep Spectral Improvement for Unsupervised Image Instance Segmentation

Farnoosh Arefi, Amir M. Mansourian, Shohreh Kasaei

*Abstract*—Recently, there has been growing interest in deep spectral methods for image localization and segmentation, influenced by traditional spectral segmentation approaches. These methods reframe the image decomposition process as a graph partitioning task by extracting features using self-supervised learning and utilizing the Laplacian of the affinity matrix to obtain eigensegments. However, instance segmentation has received less attention compared to other tasks within the context of deep spectral methods. This paper addresses the fact that not all channels of the feature map extracted from a self-supervised backbone contain sufficient information for instance segmentation purposes. In fact, some channels are noisy and hinder the accuracy of the task. To overcome this issue, this paper proposes two channel reduction modules: Noise Channel Reduction (NCR) and Deviation-based Channel Reduction (DCR). The NCR retains channels with lower entropy, as they are less likely to be noisy, while DCR prunes channels with low standard deviation, as they lack sufficient information for effective instance segmentation. Furthermore, the paper demonstrates that the dot product, commonly used in deep spectral methods, is not su itable for instance segmentation due to its sensitivity to feature map values, potentially leading to incorrect instance segments. To address this issue, a new similarity metric called Bray-Curtis over Chebyshev (BoC) is proposed. It takes into account the distribution of features in addition to their values, providing a more robust similarity measure for instance segmentation. Quantitative and qualitative results on the Youtube-VIS2019 dataset highlight the improvements achieved by the proposed channel reduction methods and the use of BoC instead of the conventional dot product for creating the affinity matrix. These improvements are observed in terms of mean Intersection over Union (mIoU) and extracted instance segments, demonstrating enhanced instance segmentation performance. The code is available on: https://github.com/farnooshar/SpecUnIIS

*Index Terms*—Deep Spectral Methods, Image Instance Segmentation, Self-Supervised Learning, Unsupervised Learning, Transformer Models

## I. INTRODUCTION

OBject segmentation, including Foregorund-Background (Fg-Bg) segmentation and instance segmentation, is a fundamental aspect of computer vision with numerous applications in domains such as medical image analysis [1], autonomous driving [2], and robotics [3]. Despite the advancements in Deep Neural Networks (DNNs), tackling these problems still presents challenges due to the requirement of dense annotation, which can be time-consuming. Moreover, relying on predefined class labels may limit the applicability of these methods, especially in domains such as medical image processing where expert annotation is necessary.

The authors are with the Department of Computer Engineering, Sharif University of Technology, Tehran 11155, Iran (e-mail: far.arefi@sharif.edu; amir.mansurian@sharif.edu; kasaei@sharif.edu).
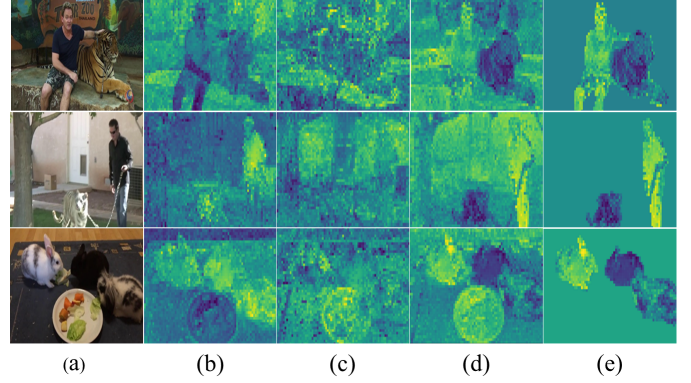


Fig. 1. Features extracted from the self-supervised backbone. (a) Input image. (b) A channel suitable for foreground-background segmentation. (c) A random channel which contains no valuable information. (d) A channel with the potential for instance segmentation. (e) A channel proper for instance segmentation, multiplied by the foreground mask. As can be seen, instances in the image are distinguishable by their pixel value.

In order to overcome the limitations of supervised learning methods, alternative approaches with reduced reliance on full supervision have emerged; such as semi-supervised learning [4], weakly-supervised learning [5], and the utilization of scribbles or clicks [6]–[8]. Nevertheless, the aforementioned methods still face challenges as they necessitate some form of annotations or prior knowledge regarding the image. On the other hand, unsupervised methods leverage self-supervised learning methods, often based on transformers, as backbones to generate attention maps that correspond to semantic segments within the input image. The extracted features from self-supervised models exhibit significant potential in aiding visual tasks; like object localization [9], object segmentation [10], semantic segmentation [11], and instance segmentation [12].

Recent advancements in this field have shown encouraging outcomes, particularly through the utilization of unsupervised learning approaches that integrate deep features with classical graph theory for object localization and segmentation tasks. Specifically, these methods employ features extracted from a self-supervised transformer-based backbone and leverage the Laplacian of the affinity matrix and patch-wise similarities for tasks such as object localization and semantic segmentation [13], [14]. However, instance segmentation has received relatively less focus. This task poses challenges as it entails recognizing and segmenting each individual object within an image, whereas semantic segmentation treats multiple objects of the same category as a single entity.

In this paper, we begin with the baseline method, Deep

Spectral Methods (DSM) [13], and conduct an experiment on the features extracted from the Dino model [15]. This analysis is motivated by the fact that some channels may contain noise and cannot be directly utilized for Fg-Bg segmentation, particularly for instance segmentation. Figure 1 illustrates that there is a specific channel within the extracted features from the self-supervised backbone that exhibits reasonable accuracy in instance segmentation, as instances can be distinguished based on their pixel values.

Based on this discovery, we hypothesize that identifying this specific channel, where instances are discernible through their pixel values, will lead to improved accuracy in instance segmentation. As depicted in Figure 2, employing this particular channel for instance segmentation yields satisfactory results across various numbers of instances. Therefore, this paper focuses on reducing the number of irrelevant channels. To achieve this, the Noise Channel Reduction (NCR) method is proposed, which reduces the number of channels based on their entropy. This channel reduction process aims to eliminate noisy channels, thereby simplifying the creation of an affinity matrix and enhancing the results of Fg-Bg segmentation.

Furthermore, an additional channel reduction method called Deviation-based Channel Reduction (DCR) is introduced. It further eliminates certain channels based on their standard deviation. The underlying idea of this module is that specific channels in the last step are useful for Fg-Bg segmentation, as they can discriminate between foreground and background regions. However, in instance segmentation, the challenge lies in distinguishing individual object instances, noting that channels with low standard deviation do not provide sufficient information for this task.

Finally, the two-step reduced channels are employed to construct an affinity matrix. At this stage, it is found that the commonly used dot product is not an optimal choice for instance segmentation, as it heavily emphasizes large or low values, which are often considered as noise. Therefore, a new metric that takes into account the distribution of features rather than their raw values is proposed. This addresses the limitations of the conventional dot product and leads to improved results in instance segmentation.

In summary, the main contributions of this work are as follows:

- Proposing Noise Channel Reduction (NCR) method for eliminating noisy channels and achieving better Fg-Bg segmentation results.
- Deviation-based Channel Reduction (DCR) method for further reducing the data dimension while preserving the channels that are valuable for instance segmentation.
- Analyzing the limitations of the dot product for instance segmentation and proposing a new metric based on the distribution of features to construct an affinity matrix suitable for the instance segmentation task.

The remaining sections of this paper are structured as follows. Some related work relevant to the proposed method are reviewed in Section II. This is followed by a detailed explanation of the proposed method in Section III. Lastly, extensive experiments and ablation studies are discussed in Section IV. Finally, Section V concludes the paper.
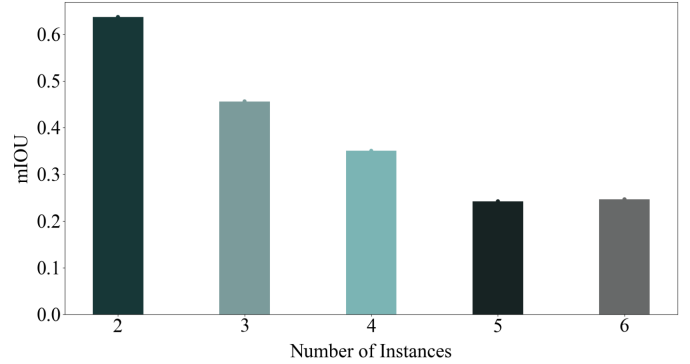


Fig. 2. mIoU results for instance segmentation, using the specific channel with significant potential for this task, are presented for various numbers of instances.

## II. RELATED WORK

In the following, the literature most relevant to this work is reviewed. This includes state-of-the-art research surrounding self-supervised learning and deep spectral methods.

### A. Self-Supervised Learning

In recent years, there have been significant advancements in self-supervised learning for visual recognition tasks. Initial approaches in this field involve training a self-supervised backbone using pretext tasks such as colorization, inpainting, or rotation [16]–[20]. The trained backbone is then utilized for the target task. More recently, the contrastive learning has been employed to generate similar embeddings for different augmentations of the image while pushing the embeddings of different images apart [21]–[24]. However, some approaches have aimed to avoid reliance on negative samples [25], [26] or have utilized clustering techniques [27], [28].

Inspired by the Bidirectional Encoder Representations from Transformers (BERT) [29], various methods have been proposed that train models by reconstructing masked tokens [30]–[32]. For instance, Masked Autoencoder (MAE) [31] takes input images with a high ratio of masked pixels and employs an encoder-decoder architecture to reconstruct the missing pixels.

In the context of self-supervised learning with vision transformers, there has been a surge of attention. The Momentum Contrast (MoCo-v3) [33], for instance, has achieved impressive results by employing contrastive learning on vision transformers. The self-Distillation with NO labels (DINO) [15], on the other hand, has introduced a self-distillation approach for training vision transformers, creating features explicitly beneficial for image segmentation.

This study employs DINO as a self-supervised backbone, leveraging its feature maps that have demonstrated efficacy in image segmentation. By incorporating deep spectral methods, the study accomplishes foreground-background separation and instance segmentation.

### B. Spectral Methods

The concept of spectral graph theory was initially introduced in [34]. Subsequent research focused on the discrete
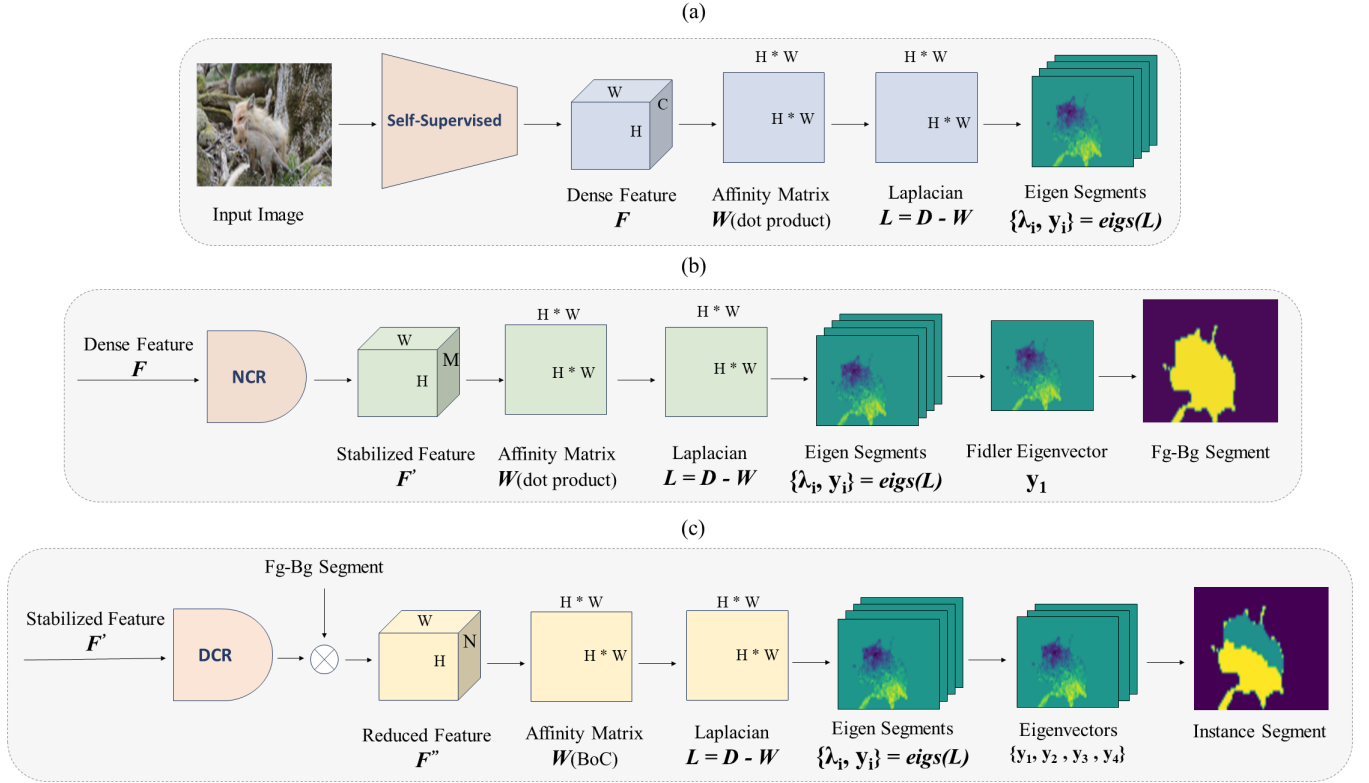
(a)



(b)



(c)



Fig. 3. Overall pipeline of proposed improved deep spectral method. (a) Workflow of [13]. An affinity matrix is created using dot product with features from a self-supervised backbone. Eigenvectors of the Laplacian matrix derived from the affinity matrix are utilized for segmentation tasks. (b) Application of the proposed CNR module on the features from the self-supervised backbone to remove noisy channels. Fiedler eigenvector is then employed for foreground-background segmentation. (c) Pipeline for instance segmentation. Stable feature map channels are further reduced based on their standard deviation to enhance richness. Feature map is multiplied by the foreground mask, and using BoC metric the affinity matrix is created. Finally, a clustering method employing eigenvectors of the Laplacian matrix is applied for instance segmentation.

formulation of graphs, establishing a connection between global graph features and the eigenvalues/eigenvectors of their Laplacian matrix [35], [36]. The work in [35] proposed that the second smallest eigenvalue of a graph, known as the Fidler eigenvalue, serves as a measure of graph connectivity. Additionally, [36] demonstrated that the eigenvectors of graph Laplacians can be utilized to achieve graph partitions with minimum energy. However, with the advancements in machine learning and computer vision, [37] and [38] introduced the pioneering methods for image segmentation through spectral clustering.

Following DSM [13], numerous works have been presented that primarily focus on object localization [14], [39], [40], semantic segmentation [11], [41], [42], and video object segmentation [10], [43], [44]. The work in [45] adopts a novel approach by utilizing DINO, along with image and flow features, as inputs. They construct a fully-connected graph, based on image patches, and employ graph-cut techniques to generate binary segmentation. The generated masks are then utilized as pseudo-labels to train a segmentation network. Additionally, [43] employs the Ncut algorithm [37] to perform segmentation using the same similarity matrix derived from image patches. They further employ the Fidler eigenvector to bi-partition the graph, followed by a refinement step for video segmentation.

However, deep spectral methods have generally received less attention in the context of instance segmentation compared to object localization and semantic segmentation. In an effort to address this gap, [40] introduces the MaskCut, a method capable of extracting masks for multiple objects within an image without relying on supervision. Similar to [43], Mask-Cut constructs a similarity matrix on a patch-wise basis using a self-supervised backbone. The Normalized Cuts algorithm is then applied to this matrix, yielding a single foreground mask for the image. Subsequently, the affinity matrix values corresponding to the foreground mask are masked out, and the process is repeated to discover masks for other objects. Another recent method, proposed by [12], investigates the features extracted from various self-supervised transformer-based backbones. That approach employs two different feature extractors. Initially, a new feature extractor is utilized, followed by spectral clustering with varying numbers of clusters. However, parallel to this step, DINO is employed along with spectral clustering using two clusters to generate a foreground mask. Finally, candidate masks from the previous step that exhibit significant intersection with the foreground mask are selected as the final masks.

In this work, starting from [13], an improvement over the deep spectral methods for instance segmentation is proposed. Through the analysis conducted in this study, it is demon-

strated that not all channels of the feature maps obtained from DINO are useful for effective instance segmentation. As a result, two steps of channel reduction are proposed to enhance the overall segmentation performance. Furthermore, the study reveals that the dot product is not a suitable option for creating the affinity matrix in the context of instance segmentation. To address this limitation, a new metric is proposed that is specifically tailored for generating the affinity matrix for instance segmentation tasks.

## III. PROPOSED METHOD

### A. Background

Consider a weighted undirected graph, denoted by $G = (V, E)$, with its corresponding adjacency matrix $W = \{w(u, v) : (u, v) \in E\}$. The Laplacian matrix $L$ of this graph can be defined as $L = D - W$. Here, $D$ represents a diagonal matrix with entries being the row-wise sums of $W$. In the realm of spectral graph theory, the eigenvectors and eigenvalues of $L$ hold significant data. The eigenvectors, denoted as $y_i$, correspond to the eigenvalues $\lambda_i$ and form an orthogonal basis that allows for the smoothest representation of functions on graph $G$. Hence, it is natural to express functions defined on $G$ by using the eigenvectors of the graph Laplacian.

In the realm of classical image segmentation, the graph $G$ can be associated with the pixels of an image $I \in \mathbb{R}^{HW}$, where $H$ and $W$ represent the dimensions of the image. The edge weights $W \in \mathbb{R}^{HW \times HW}$ correspond to the affinities or similarities between pairs of pixels. The eigenvectors $y \in \mathbb{R}^{HW}$ can be interpreted as representing soft image segments, providing a way to group pixels into coherent regions.

Graph partitions that divide the image into distinct segments are often referred to as graph cuts. In this context, the value of a cut between two partitions reflects the total weight of edges that are removed by the partitioning process. A normalized version of graph cuts, as described in [37], emerges naturally from the eigenvectors of the normalized Laplacian. In this case, achieving an optimal bi-partitioning becomes equivalent to identifying the appropriate eigenvectors of the Laplacian.

Recently, [13] proposed a method that combines the explained background with deep learning. Figure 3,(a) illustrates the pipeline of [13] approach. It utilizes feature maps extracted from a self-supervised backbone, denoted as $F$, to create a patch-wise affinity matrix $W$. From $W$, it extracts the eigenvectors of its Laplacian, $L = D - W$, which enables the decomposition of an image into soft segments: $\{y_0, ..., y_{n-1}\} = eigs(L)$. These eigensegments are then utilized for both Fg-Bg segmentation and semantic segmentation. For more details regarding this method, please refer to [13].

### B. Noise Channel Reduction (NCR)

As mentioned earlier, not all feature maps extracted from a self-supervised backbone are suitable for segmentation. Some of them may contain noise, which can negatively impact the creation of the affinity matrix. Additionally, a large number of channels can increase the computational burden when constructing the affinity matrix. To address these challenges,
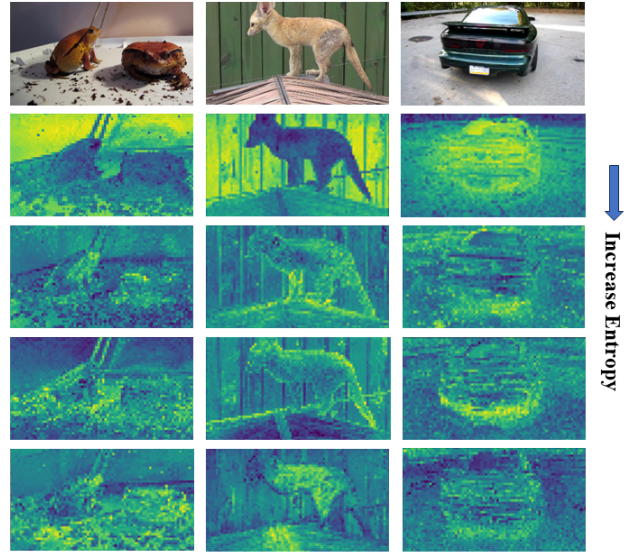


Fig. 4. Visualization of some channels from the self-supervised backbone. Last rows represent channels with higher entropy values, indicating a higher likelihood of being noise channels.

a stabilization method is proposed for reducing the number of channels based on their entropy. In this context, instability refers to the presence of unrelated or irrelevant information in the channels for the target task, compared to more stable features that contain rich information with few irrelevant components. Finding the most stable feature map is challenging as it heavily depends on the specific task or dataset. However, in this study, the aim is to move towards a more stable feature map. To achieve this, the entropy of the channels is utilized to eliminate channels with less informative content. In fact, the entropy of a channel quantifies the level of disorder or randomness within the channel. Given a feature map, the probability distribution function (or the histogram) is first calculated for each channel as:

$$PDF(c) = \frac{Hist(c)}{H \times W}; \forall c \in C, \qquad (1)$$

where $Hist(c)$ denotes the histogram of the $c$-th channel, and $H$, $W$, and $C$ represent the height, width, and the number of channels of a feature map, respectively. The entropy of a channel is then defined as follows:

$$Entropy(c) = -\sum_{b=1}^{B} PDF^b(c).log_2(PDF^b(c)), \qquad (2)$$

where $PDF^b(c)$ represents the probability of the $b$-th bin in the histogram of the channel $c$ and $B$ is equal to the total number of bins, which is considered equal to 30 in this research. The channels of the feature map $F$ are then sorted based on their entropy. The first $M$ channels with the lowest entropy are retained, where $M$ is a hyper parameter. This selected subset is denoted as $F' \in \mathbb{R}^{H \times W \times M}$, which represents a more stable version of $F$.

Figure 4 illustrates various channels of an input image, each exhibiting different levels of entropy. It demonstrates

that channels with higher entropy tend to contain more noise, while channels with lower entropy exhibit semantically richer information for distinguishing between Fg-Bg regions. The proposed feature stabilizer module not only reduces the number of channels, thereby reducing the computational burden for creating the affinity matrix, but also results in improved feature maps for Fg-Bg segmentation.

Finally, the Fg-Bg segmentation task is performed, as depicted in Figure 3,(b). Initially, the features extracted from the self-supervised backbone, denoted by $F$, pass through the NCR module, resulting in more stabilized feature maps, $F'$. Subsequently, the affinity matrix is constructed by taking the dot product between $F'$ and its transpose. The Laplacian matrix is then computed, and the second smallest eigen vector (commonly referred to as the Fiedler eigen vector), is utilized for Fg-Bg segmentation. The impact of the NCR module on segmentation will be further validated in the subsequent section.
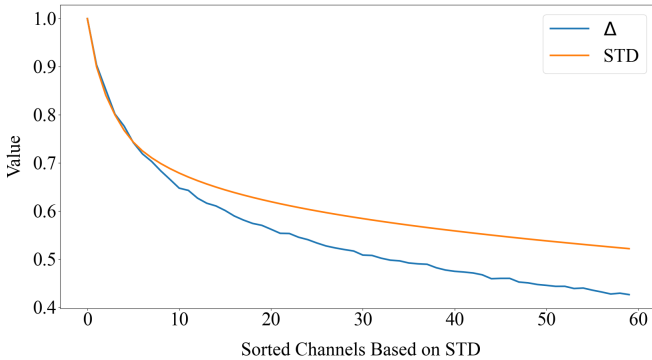


Fig. 5. Influence of DCR on distinction between instances in YouTube-VIS2019 dataset.

### C. Deviation-based Channel Reduction (DCR)

As illustrated in Figure 1, certain channels exhibit more distinguishable instances, indicating their potential usefulness in instance segmentation tasks. However, reducing the number of channels in the feature map is also desirable. To address this, the channels are sorted based on their standard deviation (STD). A channel with a lower STD tends to have less variation between instances, making it less suitable for instance segmentation. It is important to note that a higher STD does not necessarily guarantee better performance for instance segmentation, as noise can also contribute to increased STD. Nonetheless, in general, channels with higher STD are more likely to be beneficial. Furthermore, as the NCR module prunes some channels, channels that are more likely to contain noise are removed. Therefore, selecting channels based on their STD is a reasonable choice. To accomplish this, the STD of each channel is calculated using the following formula:

$$STD(c) = \sqrt{\frac{1}{H \times W} \sum_{x=1}^{H \times W} (x - \bar{x})^2}; \forall c \in F', \qquad (3)$$

where $\bar{x}$ is the mean of the values in a specific channel. Then, all channels are sorted based on their STD value, and the first $N$ channels with the highest STD are retained, where the $N$ is a hyper parameter. This process results in a final feature map denoted by $F'' \in \mathbb{R}^{H \times W \times N}$. Figure 5 demonstrates the impact of DCR with N=60. As indicated by the orange curve, the standard deviation of the channels decreases, and simultaneously, $\Delta$, representing the average difference between instances, also decreases. Therefore, there is a strong likelihood that channels with a higher standard deviation will exhibit a greater average difference between instances.
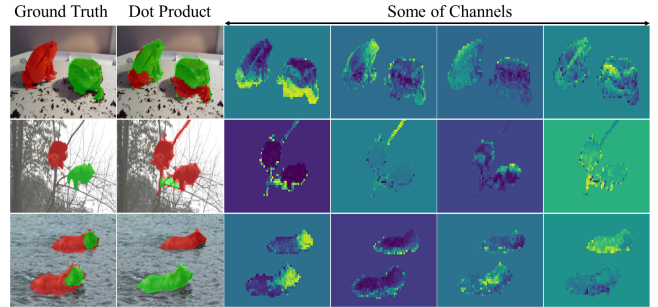


Fig. 6. Qualitative results for instance segmentation, when using dot product to create the affinity matrix. As illustrated, in some channels from the feature map, there are pixels with very high or low values. Using the dot product, which is sensitive to these values, can lead to incorrect instance segmentation outputs.

### D. Deep Spectral Instance Segmentation

*1) Dot Product for Instance Segmentation:* As shown in Figure 3,(a), [13] utilizes the dot product to create an affinity matrix, which is subsequently employed in downstream tasks such as Fg-Bg segmentation or semantic segmentation. However, in this section, we challenge the suitability of the dot product as a proper method for generating the affinity matrix specifically for the task of instance segmentation. This claim is supported by two reasons:

- **Sensitivity to the values of feature vector**: The dot product is highly sensitive to the values in the feature maps. If there are extreme values, either very high or very low (referred to as irregular values), within a feature vector, they can significantly influence the affinity matrix even after normalization. Figure 6 illustrates an instance where one of the channels in the feature maps contains irregular values. Consequently, the final affinity matrix is heavily impacted by these values, resulting in their segmentation as a distinct region.

  While the dot product can be useful for Fg-Bg segmentation, as it takes into account the significant differences in pixel values, especially along edges, it may not be as effective for instance segmentation. In instance segmentation, the focus is on capturing more than just the values of the pixels; additional details are required to accurately segment different instances.

- **Neglecting the distribution of features**: In the context of image instance segmentation, it is crucial to take into

account the distribution of features, rather than solely relying on their exact values. Pixels with similar feature distributions should be segmented as same instances. As depicted in Figure 6, there are irregular values present in some of the channels. However, it is important to note that these seemingly irregular patterns are actually part of the same instance. The problem with using the dot product is that it treats these regions as separate instances, disregarding the underlying feature distribution. As a result, the dot product is not suitable for instance segmentation, as it fails to consider the distribution of features.

An ideal affinity matrix would possess the property that pixels belonging to the same instance exhibit similar feature distributions while having maximum dissimilarity with pixels from other instances. To achieve this, a metric that takes into account the feature distribution is necessary. Additionally, for effective comparison, normalization should be carried out using a similarity metric that is sensitive to the values present in the feature maps.

*2) Capturing Feature Distribution:* In order to effectively capture information about the distribution of features, we propose the use of the Bray-Curtis similarity metric. Unlike metrics that rely solely on exact values, the Bray-Curtis metric places emphasis on the distribution aspect. Originally developed for ecological or community data analysis [46], this metric has also found valuable applications in machine learning tasks. The Bray-Curtis similarity and dissimilarity between two feature vectors, $U$ and $T$, are defined as follows:

$$BC_{diss} = \frac{\sum |u_i - t_i|}{\sum |u_i + t_i|}, \qquad (4)$$

$$BC_{sim} = \frac{1}{1 + BC_{diss}}. \qquad (5)$$

The Bray-Curtis similarity metric is capable of effectively capturing instances with similar patterns but different values in the affinity matrix. This can be better understood by examining Figure 7, which highlights the difference between the Bray-Curtis similarity ($BC_{sim}$) and the dot product. In this example, three pixels belonging to the same instance, characterized by nearly identical feature distributions, are depicted. Some channels contain randomly added or subtracted high values to the original values. When employing the dot product, the similarity matrix, which reflects the resemblance between these three pixels, contains distinct values. On the contrary, when utilizing the Bray-Curtis similarity metric, the three pixels exhibit a significant degree of closeness. This served as an abstract example for comparing the Bray-Curtis and Dot product metrics. In the subsequent section, a quantitative comparison among various metrics will be conducted to measure the ratio of variance within an instance to that outside the instance using the tested dataset.

*3) Affinity Matrix Creation:* As mentioned previously, when constructing the affinity matrix, it is crucial to consider a criterion for the similarity of pixels, taking into account the distribution of features. For a more precise and accurate comparison, normalization to the similarity criterion with higher sensitivity values is necessary. We have previously
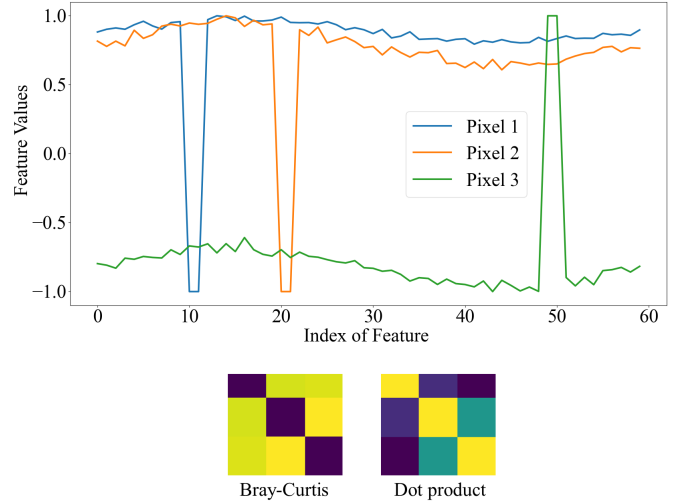


Fig. 7. Comparison of dot product and Bray-Curtis metrics in creating the affinity matrix. Three pixels, belonging to the same instance, are selected from the feature maps, and random noise is added to some of their channels. The affinity matrix created by the dot product is not suitable for our purpose, while the matrix created by the Bray-Curtis matrix correctly demonstrates the correlation between the three pixels.

observed that specific feature channels provide a highly accurate representation of the instance segmentation task. In these channels, the feature values of different instances significantly differ from each other, contributing to an enhancement in segmentation when appropriately considered. When examining the distribution of features, less emphasis is placed on the values. To prioritize the feature values, the similarity of two feature vectors can be equated to the inverse of their maximum distance, known as the Chebyshev distance. The Chebyshev dissimilarity ($CH_{diss}$) and the Chebyshev similarity criterion ($CH_{sim}$) between two feature vectors, $U$ and $T$ are then defined as follows:

$$CH_{diss} = max_i(|u_i - t_i|), \qquad (6)$$

$$CH_{sim} = \frac{1}{1 + CH_{diss}}. \qquad (7)$$

With respect to the two proposed criteria, the Braycurtis over Chebyshev ($BoC$) is defined for constructing the affinity matrix, derived from the ratio of Bray Curtis and Chebyshev criteria:

$$BoC = \frac{BC_{sim}}{CH_{sim}}. \qquad (8)$$

Utilizing the Chebyshev similarity alone may result in the inclusion of non-significant regions. This is because giving maximum attention to the difference between two feature vectors increases the probability of capturing irrelevant information. However, when considering the feature distribution, we can apply a penalty value proportional to the maximum difference between the two vectors by utilizing the Chebyshev distance. If there is a small Chebyshev distance between two vectors, the similarity between the vertices reflects the similarity of the feature distribution. However, the larger

this distance, the greater the penalty applied to this feature distribution.

*4) Instance Segmentation Paradigm:* The final framework for instance segmentation is illustrated in Figure 3,(c). Stabilized features, denoted as $F'$, are fed into the proposed DCR module, where channels with higher standard deviation are retained, $F''$. These selected channels are then multiplied by the foreground mask, and an affinity matrix is created using the proposed BoC similarity metric. In the last step, the first four small eigensegments, excluding $y_0$ (which is equivalent to noise), are utilized for instance segmentation. To achieve this, a clustering algorithm with an appropriate number of classes is applied to the eigensegments, resulting in the extraction of instance masks.

## IV. EXPERIMENTS

This section begins by introducing the dataset and evaluation metrics, and implementation details. Then, experimental results and ablation studies are discussed to validate the proposed method.

### A. Dataset and Evaluation Metrics

For the Fg-Bg segmentation task, two datasets are utilized: YouTube-VIS2019 (train) [47] and PascalVOC 2012 [48] (train/validation), comprising 61,845 and 2,913 images, respectively.

Since the focus of this work is instance segmentation, images containing more than one instance are specifically chosen for analysis. Additionally, due to the small spatial size of the feature maps generated by the self-supervised backbone, images are excluded if the size of their objects is smaller than 0.07 times the image size or if the ratio of the smallest instance to the largest instance is less than 0.3. Consequently, the final dataset consists of 10,285 images that meet these criteria and are used as the test set.

For evaluating the Fg-Bg segmentation task, the F-score metric is employed. The F-score combines precision and recall to provide a balanced measure of the model's performance. In binary segmentation, the objective is to classify each pixel or region in an image as either foreground or background. Precision measures the proportion of correctly classified foreground pixels out of all the pixels predicted as foreground, while recall quantifies the proportion of correctly classified foreground pixels out of all the true foreground pixels in the image.

For the instance segmentation task, a linear assignment is performed using the Hungarian algorithm to determine the correspondence between predictions and instances. The average mIoU of all instances is then reported. The mIoU measures the intersection of the prediction and ground truth masks divided by their union, providing an overall evaluation of the instance segmentation performance.

### B. Implementation Details

After extracting the mask using the proposed method for Fg-Bg segmentation, a post-processing step is performed.

To eliminate small regions, a median kernel with a size of 5x5 is applied to the mask. Additionally, in cases where the foreground and background are mistakenly predicted, the foreground and background labels are swapped while taking into account the boundaries according to DINO [15], similar to [13]. The backbone we employed in the upcoming tests is ViT-s16 [15].
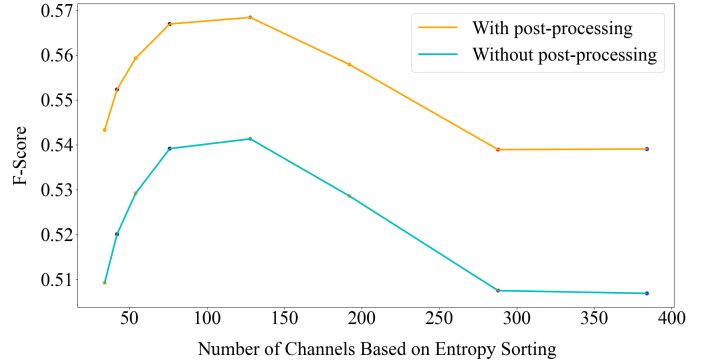


Fig. 8. Results of Fg-Bg segmentation, for various values of "M" on Youtube-VIS2019 dataset. Channels are sorted in ascending order based on their entropy
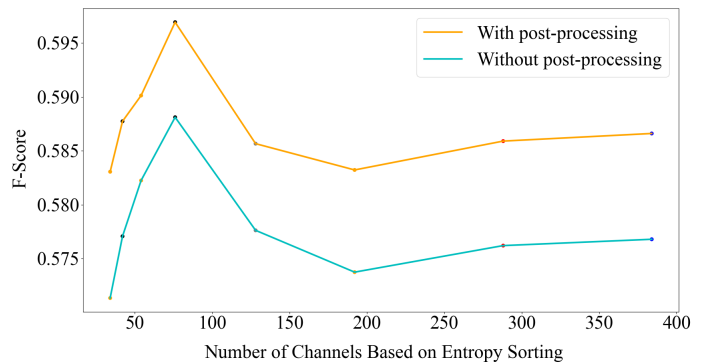


Fig. 9. Results of Fg-Bg segmentation, for different values of "M" on PascalVOC 2012 dataset. Channels are arranged in ascending order based on their entropy.

### C. Experimental Results

*1) Fg-Bg Segmentation:* To validate the effectiveness of the proposed NCR module, experiments were conducted to evaluate the Fg-Bg segmentation results for different values of "M" explained in the last section. Figure 8 illustrates the impact of stabilizing features with varying values of "M", both with and without post-processing, on the Youtube-VIS2019 dataset. The results indicate that preserving 1/3 of the channels (M=1/3 C) with the lowest entropy yields the best F-score for this task. Figure 9 present similar results obtained from the PascalVOC 2012 dataset, where M=1/5C leads to the best performance. Table I provides quantitative numbers for different values of "M" for both datasets. It demonstrates that, for Youtube-VIS2019 dataset, retaining nearly 1/3 of the channels with the lowest entropy can improve segmentation results by 2% with post-processing and 3% without post-processing.

Notably, even retaining only 20% of the channels can yield better results than using all the channels, validating the fact that self-supervised learning can introduce noisy channels and not all channels contribute to effective segmentation.

TABLE I
F-SCORE RESULTS FOR FG-BG SEGMENTATION, CONSIDERING DIFFERENT VALUES OF M, WITH AND WITHOUT POST-PROCESSING.

| Youtube-VIS2019 | | |
|---|---|---|
| **Value of M** | **w post-processing** | **w/o post-processing** |
| M = C | 53.9 | 50.68 |
| M = 3C/4 | 53.89 | 50.74 |
| M = C/2 | 55.79 | 52.85 |
| M = C/3 | **56.84** | **54.13** |
| M = C/5 | 56.69 | 53.91 |
| M = C/7 | 55.93 | 52.91 |
| M = C/9 | 55.23 | 52 |
| M = C/11 | 54.33 | 50.92 |
| PascalVoc 2012 | | |
| **Value of M** | **w post-processing** | **w/o post-processing** |
| M = C | 58.66 | 57.68 |
| M = 3C/4 | 58.59 | 57.62 |
| M = C/2 | 58.32 | 57.37 |
| M = C/3 | 58.56 | 57.76 |
| M = C/5 | **59.69** | **58.81** |
| M = C/7 | 59.01 | 58.22 |
| M = C/9 | 58.77 | 57.70 |
| M = C/11 | 58.30 | 57.13 |

TABLE II
A COMPARISON OF INSTANCE SEGMENTATION RESULTS BETWEEN DIFFERENT METRICS FOR CREATING THE AFFINITY MATRIX, WITH THE PROPOSED METRIC, IN TERMS OF MIOU ON YOUTUBE-VIS2019 DATASET.

| **Metric** | **mIoU (%)** |
|---|---|
| Mahalanobis | 25.27 |
| L1 | 31.53 |
| Dot product | 32.71 |
| L2 | 32.77 |
| Chebyshev | 33.09 |
| Cosine | 33.56 |
| Correlation | 34.08 |
| Braycurtis | 34.14 |
| **BoC** | **34.41** |

TABLE III
QUANTITATIVE RESULTS FOR DIFFERENT METRICS UNDER VARYING LEVELS OF OCCLUSION IN TERMS OF MIOU.

| **Metric** | **MBOR 0.01-0.14** | **MBOR 0.14-0.34** | **MBOR ≥ 0.34** |
|---|---|---|---|
| Mahalanobis | 24.46 | 27.20 | 26.86 |
| L1 | 31.92 | 34.18 | 30.25 |
| Dot product | 33.50 | 35.25 | 30.69 |
| L2 | 33.42 | 35.26 | 30.99 |
| Chebyshev | 33.65 | 35.91 | 31.31 |
| Cosine | 34.03 | 36.26 | 31.38 |
| Correlation | 34.76 | 36.67 | 31.64 |
| Braycurtis | 34.67 | 37.06 | 31.88 |
| **BoC** | **34.94** | **37.38** | **32.33** |

The findings suggest that reducing channels based on their entropy results in more stable feature maps, which enhances segmentation performance while reducing the memory and time required for calculating the affinity matrix. However, reducing a significant number of channels also leads to a decrease in results, indicating that useful channels from the backbone network have been discarded. Figure 10 provides

TABLE IV
QUANTITATIVE RESULTS FOR INSTANCE SEGMENTATION CONSIDERING DIFFERENT VALUES OF THE RATIO, IN TERMS OF MIOU.

| **Metric** | **Ratio 1.26-1.64** | **Ratio ≥1.64** |
|---|---|---|
| Mahalanobis | 26.57 | 23.98 |
| L1 | 32.10 | 30.97 |
| Dot product | 34.23 | 31.20 |
| L2 | 33.54 | 31.99 |
| Chebyshev | 34.08 | 32.10 |
| Cosine | 34.23 | 32.90 |
| Correlation | 34.98 | 33.18 |
| Braycurtis | 34.91 | 33.38 |
| **BoC** | **35.10** | **33.71** |

visual examples illustrating the impact of stabilizing the feature map on segmentation results. It demonstrates that a more stable feature map leads to improved segmentation quality, as unrelated objects belonging to the background are correctly identified and not considered as foreground.

*2) Instance Segmentation:* Table II presents the mIoU results for instance segmentation using various metrics for creating the affinity matrix. Cosine, correlation, L1, L2, and Mahalanobis metrics are reversed, akin to the relationships in equations 5 and 7, to demonstrate similarity. NCR and DCR with M=128 and N=60 are utilized in all tables II,III and IV. The results indicate that the proposed metric, BoC, achieves the best performance, with approximately 2% higher mIoU compared to the dot product metric used in [13]. The proposed metric demonstrates robustness to the values of feature vectors and performs well in tasks where the value alone is insufficient for discriminating the data, such as instance segmentation. In contrast, metrics like cosine similarity, correlation, L1 and L2 distances are more sensitive to the values of feature vectors.

Another advantage of the proposed metric is its performance in occlusion scenarios. It effectively highlights differences and is less susceptible to the effects of occlusion. Table III showcases the instance segmentation results under various levels of occlusion. The Mean Boundary Overlap Ratio (MBOR) introduced by [49] indicates the degree of occlusion, with a higher MBOR value indicating more severe occlusion. It can be observed that the proposed metric performs better in scenarios with heavy occlusion. Figure 10 provides a qualitative comparison of different metrics for varying MBOR values.

Furthermore, when two instances in an image have different sizes, it is common for the smaller object to be considered as part of the larger object. However, the proposed metric handles this situation more effectively. Table IV presents the results, where the ratio represents the sum of the ratios of each object's size to the size of the largest object.

For further validation of the proposed BoC metric, an experiment was conducted. Firstly, 10 pixels were randomly sampled from each ground-truth mask. Subsequently, using various similarity metrics, the intra-instance distances between these pixels were calculated. It was expected that these intra-instance distances, being from the same instance, would exhibit low variance and similarity. The mean of intra-instance variance, $mVar_{intra}$, was then computed. Additionally, the
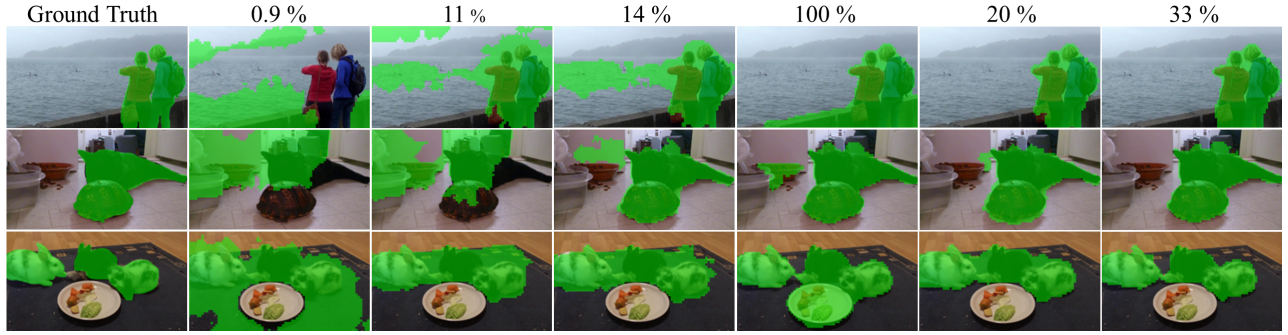
Fig. 10. Qualitative outcomes of Fg-Bg segmentation on Youtube-VIS2019 dataset. Percentages indicate the proportion of channels that are preserved after being sorted according to their entropy. Precise number of channels to be retained varies for each dataset. It is determined during the generation of final masks.
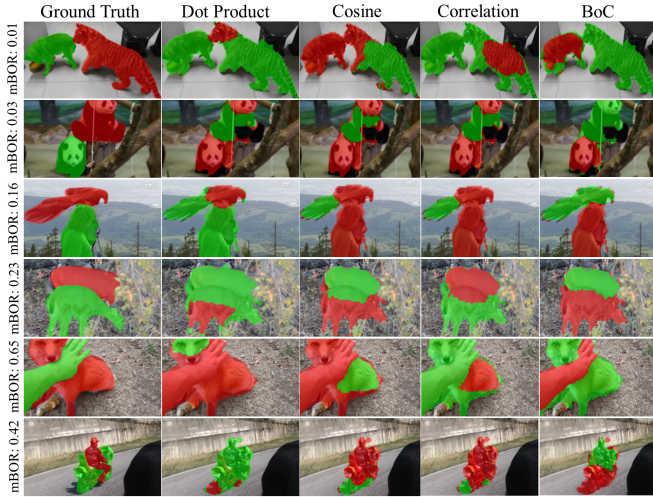


Fig. 11. Instance segmentation results under varying levels of occlusion, represented by the MBOR value, while utilizing different metrics for creating the affinity matrix. Proposed metric, $BoC$, outperforms other metrics and produces more accurate masks, even in scenarios with heavy occlusions.

mean of inter-instance distance variance, $mVar_{inter}$, was calculated. The ratio of $mVar_{intra}$ to $mVar_{inter}$, denoted as $mR$, was compared for different similarity metrics. As shown in Figure 13, the proposed BoC metric demonstrated a lower mR value compared to other metrics. This result validates that the BoC metric better accounts for both intra-instance and inter-instance similarities compared to other metrics, resulting in improved instance segmentation.

Table V includes the ablation analysis for the proposed components. The results demonstrate that all components significantly contribute to improving the performance of instance segmentation. Starting from a baseline mIoU of 31.75%, the addition of NCR leads to a significant improvement to 32.92%, showcasing its individual effectiveness. Interestingly, introducing DCR alongside NCR results in a marginal decrease to 32.70%, suggesting a nuanced interaction between the two components that warrants further exploration.

Isolating BoC alone yields a lower mIoU of 30.50%, indicating that BoC's standalone contribution may not be as impactful in enhancing instance segmentation. However, combining NCR with BoC produces a synergistic effect, boosting

mIoU to 33.62%. The most compelling outcome arises when all three components—NCR, DCR, and BoC—are integrated, achieving the highest mIoU of 34.41%. This holistic approach underscores the complementary nature of the components, overcoming individual limitations and emphasizing their collective significance in advancing instance segmentation accuracy.
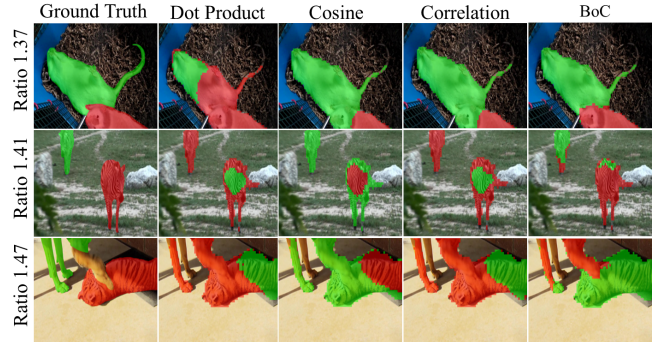


Fig. 12. Extracted masks for different metrics on Youtube-VIS2019 dataset. A lower value of ratio indicates greater variation in object sizes.
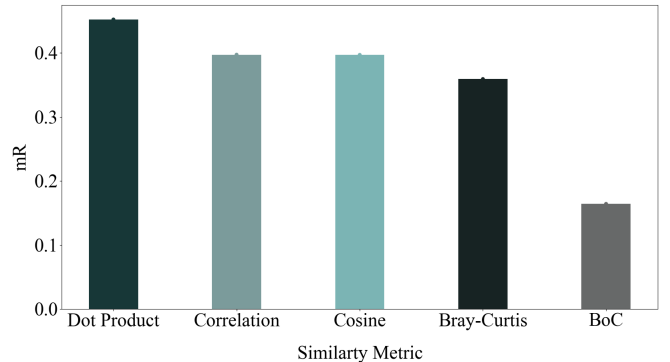


Fig. 13. Ratio of variance of intra-instance similarity to inter-class similarity for different similarity metrics.

## V. CONCLUSION

This paper focused on enhancing the performance of deep spectral methods specifically for instance segmentation purposes. To achieve this, two modules were proposed to retain

TABLE V
AN ABLATION STUDY TO ANALYZE THE IMPACT OF PROPOSED
COMPONENTS IN TERMS OF mIoU.

| NCR | DCR | BoC | mIoU (%) |
|:---:|:---:|:---:|:---:|
|  |  |  | 31.75 |
| ✓ |  |  | 32.92 |
| ✓ | ✓ |  | 32.70 |
|  |  | ✓ | 30.50 |
| ✓ |  | ✓ | 33.62 |
| ✓ | ✓ | ✓ | **34.41** |

channels with the most informative content from the feature maps obtained from a self-supervised backbone. Moreover, it was shown that the conventional use of dot product for creating the affinity matrix has limitations when it comes to instance segmentation. To address this issue, a novel similarity metric was introduced, which aims to improve the affinity matrix for instance segmentation tasks. The proposed components were designed to enhance the quality of eigensegments extracted through deep spectral methods for instance segmentation. These components can easily be integrated with any deep spectral method that aims to solve instance segmentation problems. It should be noted that in the context of supervised learning, it would also be possible to train networks, specifically for both channel reduction modules, enabling the effective retention of the most informative channels.

## REFERENCES

[1] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18.* Springer, 2015, pp. 234–241.

[2] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2017.

[3] J. Dai, K. He, and J. Sun, "Instance-aware semantic segmentation via multi-task network cascades," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3150–3158.

[4] Y. Wei, H. Xiao, H. Shi, Z. Jie, J. Feng, and T. S. Huang, "Revisiting dilated convolution: A simple approach for weakly-and semi-supervised semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7268–7277.

[5] D. Pathak, P. Krahenbuhl, and T. Darrell, "Constrained convolutional neural networks for weakly supervised segmentation," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1796–1804.

[6] M. Tang, F. Perazzi, A. Djelouah, I. Ben Ayed, C. Schroers, and Y. Boykov, "On regularized losses for weakly-supervised cnn segmentation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 507–522.

[7] A. Khoreva, R. Benenson, J. Hosang, M. Hein, and B. Schiele, "Simple does it: Weakly supervised instance and semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 876–885.

[8] T.-W. Ke, J.-J. Hwang, and S. X. Yu, "Universal weakly supervised segmentation by pixel-to-segment contrastive learning," *arXiv preprint arXiv:2105.00957*, 2021.

[9] G. Shin, S. Albanie, and W. Xie, "Unsupervised salient object detection with spectral cluster voting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3971–3980.

[10] S. Choudhury, L. Karazija, I. Laina, A. Vedaldi, and C. Rupprecht, "Guess what moves: Unsupervised video and image segmentation by anticipating motion," *arXiv preprint arXiv:2205.07844*, 2022.

[11] K. Li, Z. Wang, Z. Cheng, R. Yu, Y. Zhao, G. Song, C. Liu, L. Yuan, and J. Chen, "Acseg: Adaptive conceptualization for unsupervised semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7162–7172.

[12] P. Engstler, L. Melas-Kyriazi, C. Rupprecht, and I. Laina, "Understanding self-supervised features for learning unsupervised instance segmentation," *arXiv preprint arXiv:2311.14665*, 2023.

[13] L. Melas-Kyriazi, C. Rupprecht, I. Laina, and A. 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*, 2022, pp. 8364–8375.

[14] Y. Wang, X. Shen, S. X. Hu, Y. Yuan, J. L. Crowley, and D. Vaufreydaz, "Self-supervised transformers for unsupervised object discovery using normalized cut," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 14 543–14 553.

[15] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 9650–9660.

[16] C. Doersch, A. Gupta, and A. A. Efros, "Unsupervised visual representation learning by context prediction," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1422–1430.

[17] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," *arXiv preprint arXiv:1803.07728*, 2018.

[18] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.

[19] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2536–2544.

[20] R. Zhang, P. Isola, and A. A. Efros, "Colorful image colorization," in *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part III 14.* Springer, 2016, pp. 649–666.

[21] S. A. Koohpayegani, A. Tejankar, and H. Pirsiavash, "Mean shift for self-supervised learning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10 326–10 335.

[22] Y. Kalantidis, M. B. Sariyildiz, N. Pion, P. Weinzaepfel, and D. Larlus, "Hard negative mixing for contrastive learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 21 798–21 809, 2020.

[23] X. Chen, H. Fan, R. Girshick, and K. He, "Improved baselines with momentum contrastive learning," *arXiv preprint arXiv:2003.04297*, 2020.

[24] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning.* PMLR, 2020, pp. 1597–1607.

[25] X. Chen and K. He, "Exploring simple siamese representation learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 15 750–15 758.

[26] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar *et al.*, "Bootstrap your own latent-a new approach to self-supervised learning," *Advances in neural information processing systems*, vol. 33, pp. 21 271–21 284, 2020.

[27] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, "Unsupervised learning of visual features by contrasting cluster assignments," *Advances in neural information processing systems*, vol. 33, pp. 9912–9924, 2020.

[28] J. Li, P. Zhou, C. Xiong, and S. C. Hoi, "Prototypical contrastive learning of unsupervised representations," *arXiv preprint arXiv:2005.04966*, 2020.

[29] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[30] H. Bao, L. Dong, and F. Wei, "Beit: Bert pre-training of image transformers. arxiv 2021," *arXiv preprint arXiv:2106.08254*.

[31] Z. Li, Z. Chen, F. Yang, W. Li, Y. Zhu, C. Zhao, R. Deng, L. Wu, R. Zhao, M. Tang *et al.*, "Mst: Masked self-supervised transformer for visual representation," *Advances in Neural Information Processing Systems*, vol. 34, pp. 13 165–13 176, 2021.

[32] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 16 000–16 009.

[33] X. Chen, S. Xie, and K. He, "An empirical study of training self-supervised vision transformers. in 2021 ieee," in *CVF International Conference on Computer Vision (ICCV)*, pp. 9620–9629.

[34] J. Cheeger, "A lower bound for the smallest eigenvalue of the laplacian, problems in analysis (papers dedicated to salomon bochner, 1969)," 1970.

[35] M. Fiedler, "Algebraic connectivity of graphs," *Czechoslovak mathematical journal*, vol. 23, no. 2, pp. 298–305, 1973.

[36] W. E. Donath and A. J. Hoffman, "Lower bounds for the partitioning of graphs," *IBM Journal of Research and Development*, vol. 17, no. 5, pp. 420–425, 1973.

[37] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 22, no. 8, pp. 888–905, 2000.

[38] A. Ng, M. Jordan, and Y. Weiss, "On spectral clustering: Analysis and an algorithm," *Advances in neural information processing systems*, vol. 14, 2001.

[39] O. Siméoni, G. Puy, H. V. Vo, S. Roburin, S. Gidaris, A. Bursuc, P. Pérez, R. Marlet, and J. Ponce, "Localizing objects with self-supervised transformers and no labels," *arXiv preprint arXiv:2109.14279*, 2021.

[40] X. Wang, R. Girdhar, S. X. Yu, and I. Misra, "Cut and learn for unsupervised object detection and instance segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 3124–3134.

[41] A. Zadaianchuk, M. Kleindessner, Y. Zhu, F. Locatello, and T. Brox, "Unsupervised semantic segmentation with self-supervised object-centric representations," *arXiv preprint arXiv:2207.05027*, 2022.

[42] A. Aflalo, S. Bagon, T. Kashti, and Y. Eldar, "Deepcut: Unsupervised segmentation using graph neural networks clustering," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 32–41.

[43] Y. Wang, X. Shen, Y. Yuan, Y. Du, M. Li, S. X. Hu, J. L. Crowley, and D. Vaufreydaz, "Tokencut: Segmenting objects in images and videos with self-supervised transformer and normalized cut," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

[44] G. Ponimatkin, N. Samet, Y. Xiao, Y. Du, R. Marlet, and V. Lepetit, "A simple and powerful global optimization for unsupervised video object segmentation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 5892–5903.

[45] S. Singh, S. Deshmukh, M. Sarkar, and B. Krishnamurthy, "Locate: Self-supervised object discovery via flow-guided graph-cut and bootstrapped self-training," *arXiv preprint arXiv:2308.11239*, 2023.

[46] J. R. Bray and J. T. Curtis, "An ordination of the upland forest communities of southern wisconsin," *Ecological monographs*, vol. 27, no. 4, pp. 326–349, 1957.

[47] L. Yang, Y. Fan, and N. Xu, "Video instance segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5188–5197.

[48] 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.

[49] J. Qi, Y. Gao, Y. Hu, X. Wang, X. Liu, X. Bai, S. Belongie, A. Yuille, P. H. Torr, and S. Bai, "Occluded video instance segmentation: A benchmark," *International Journal of Computer Vision*, vol. 130, no. 8, pp. 2022–2039, 2022.