

Revisiting Graph Construction for Fast Image Segmentation

Zizhao Zhang¹, Fuyong Xing², Hanzi Wang⁴,
Yan Yan⁴, Ying Huang⁴, Xiaoshuang Shi³, Lin Yang^{1,2,3,*}

¹*Dept. of Computer and Information Science and Engineering, University of Florida, FL 32611, USA*
²*Dept. of Biostatistics and Informatics, Colorado School of Public Health, University of Colorado Denver,
Denver, CO 80045 USA*

³*J. Crayton Pruitt Family Dept. of Biomedical Engineering, University of Florida, FL 32611, USA*

⁴*Fujian Key Laboratory of Sensing and Computing for Smart City, School of Information Science and
Engineering, Xiamen University, Fujian 361005, China*

Abstract

In this paper, we propose a simple but effective method for fast image segmentation. We re-examine the locality-preserving character of spectral clustering by constructing a graph over image regions with both global and local connections. Our novel approach to build graph connections relies on two key observations: 1) local region pairs that co-occur frequently will have a high probability to reside on a common object; 2) spatially distant regions in a common object often exhibit similar visual saliency, which implies their neighborhood in a manifold. We present a novel energy function to efficiently conduct graph partitioning. Based on multiple high quality partitions, we show that the generated eigenvector histogram based representation can automatically drive effective unary potentials for a hierarchical random field model to produce multi-class segmentation. Sufficient experiments, on the BSDS500 benchmark, large-scale PASCAL VOC and COCO datasets, demonstrate the competitive segmentation accuracy and significantly improved efficiency of our proposed method compared with other state of the arts.

Keywords: Image segmentation, Graph partition, Manifold

*Corresponding author

Email address: zizhaozhang@ufl.edu, fuyong.xing@ucdenver.edu,
hanzi.wang@xmu.edu.cn, yanyan@xmu.edu.cn, yinghuang@stu.xmu.edu.cn
xssh2013@gmail.com, lin.yang@bme.ufl.edu (Zizhao Zhang¹, Fuyong Xing², Hanzi
Wang⁴,
Yan Yan⁴, Ying Huang⁴, Xiaoshuang Shi³, Lin Yang^{1,2,3,*})

1. Introduction

Image segmentation is a challenging and critical computer vision task. Graph-based algorithms have been shown as an effective approach for image segmentation [1, 2, 3]. Among various graph based approaches, spectral clustering becomes a major trend [4, 5].

Recent methods attempt to solve several primary issues of spectral clustering (referring to normalized cuts (NCut) [6]) based image segmentation to segment image into meaningful partitions. First, NCut based methods tend to segment image into spatially connected components [6, 7]. Multiscaling processing [8, 9] is a common way to address this problem by building the affinity for distant pixel affinities [10, 9]. However, the usage of these methods for real large-scale datasets is not clear. Most current cutting-edge methods do not follow this direction. Instead, recent methods, like gPb [11] and MCG [12, 7] based methods [13] use the boundary-preserving property of NCut to trace boundary orientation information rather than direct segmentation. Building effective affinity matrices [7, 4, 12] usually uses sophisticated low-level features [11]. These features can effectively measure the local changes but are not effective in capturing high-level knowledge for segmentation. They are not good options for fast segmentation either due to high computational cost [11]. Different from previous approaches, our method re-examines spectral clustering from a manifold learning perspective to construct a graph to model the high-level image knowledge (i.e., pixel pair co-occurrence and saliency relationship) for unsupervised image segmentation. More importantly, our method provides the possibility of enabling graph partitioning to directly segment challenging natural images rather than just boundary tracing.

To better illustrate the motivation, we first explain the latent relation of NCut to manifold learning. Both NCut and Laplacian eigenmaps [14] take advantage of the locality-preserving character [15] of graph Laplacian to conduct clustering and dimensionality reduction. In fact, locality-based dimensionality reduction methods are implicitly tied to clustering [14, 16]. Preserving locality is the key factor that drives effective clustering. Let's assume pixels of an image lie on a certain manifold where pixels belonging to a common object are adjacent (within a small range), but far away in the

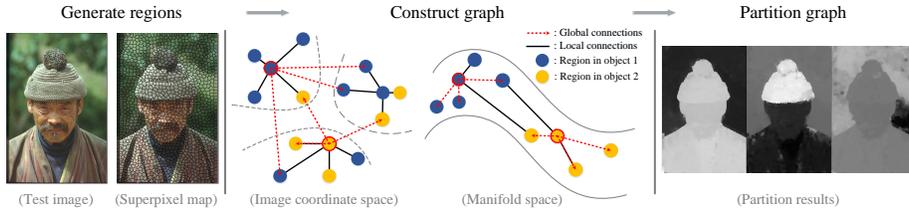


Figure 1: Illustration of the graph construction and partition procedures. The test image is first divided into over-segmented regions (left). Each region is treated as a graph node (middle). Local connections (black solid lines) link each region to its spatially adjacent regions in the image plane. When viewing those regions in a certain manifold space, spatially distant regions in the same object will be adjacent; directed global connections (red dotted arrows) link each region to its neighbors. We cut the graph to obtain multiple partitions (right).

spatial image plane. These pixels supposed to have strong connections to be grouped together, but these connections are not encoded in the sparse affinity matrix of NCut due to their Euclidean distances. Although multi-scale affinity matrices [8, 10] can alleviate this issue, increasing the range of an affinity matrix and connecting all pixel pairs in the range also introduce unavoidable noises. The method to construct affinities between spatially distant and adjacent pixels should be considered respectively in order to better capture their respective characteristics in image statistics.

In this paper, we propose a novel approach to construct an image region graph to address the aforementioned problems. The overall idea is illustrated in Figure 1. The graph nodes are connected among both spatially adjacent and distant regions through different and independent cues. We build local connections between spatially adjacent regions with an affinity matrix. The estimation of the similarity between two regions is based on an observation that adjacent region pairs co-occurring frequently often reside on a common object. Oppositely, global connections are built among adjacent regions in the manifold which might be spatially distant, with an objective to preserve their relationships and encourage them to be clustered together. We introduce a simple cue to discover the similar saliency of those regions as the global connection measurement. We present a new energy function to partition the constructed graph, which formulates the minimization problem as a single and efficiently solvable eigenvector system.

Based on the generated high quality graph partitions, we present a simple eigenvector histogram based representation to represent image regions and automatically drive effective unary potentials for the hierarchical random field of the Pylon model [17], yielding high-quality multi-class segmentation.

In brief, the contributions of this paper are:

- We propose a sophisticatedly connected graph to build the connection of image regions yet with very efficient graph partitioning capability.
- We exploit various simple and efficient cues to capture the high-level image information in order to segment objects with complex inner-variances and background.
- We present a multi-class segmentation strategy by utilizing graph partitions to generate clear and smooth segmentation.
- Extensive experiments and comprehensive analysis are conducted, on BSDS500 [11], large-scale PASCAL VOC [18] and COCO [19] datasets, to validate the effectiveness of our proposed method, its generalization ability to different datasets with diverse scenes, and the high efficiency compared with other state of the arts.

The rest of the paper are structured as follow: Section 2 discusses the related work. Section 3 introduces the graph construction and partition of our method. Section 4 introduces the proposed multi-class segmentation by utilizing graph partitions. Finally, Section 5 conducts experiments and detailed analysis. Section 6 concludes the paper.

2. Related Works

Image segmentation has been studied in the computer vision community for decades. Shi *et al.* [6] propose normalized cuts (NCut), which advanced spectral clustering based image region segmentation. [20] enables its multi-class segmentation. Among the region based segmentation, diffusion based approaches [21, 22], GraphCut [23], GrabCut [23], etc [24, 25, 26, 27], have been explored to partition images. Building successful affinity matrices is critical [28]. Many subsequent approaches have computed more effective affinity matrices using elaborately designed low-level features and

metrics [4, 10, 11, 29]. To solve the limitation of NCut to capture affinities of distant pixels, several methods [8, 9, 29, 30] have been proposed base on multi-scaling affinity strategies. However, dense affinity suffers from optimization bottleneck, although approximation algorithms are explored [10, 9, 12]. Our method is able to capture both local and global affinities as well keeps the sparsity of the affinity matrix.

Contour driven image region segmentation is widely studied. Arbelaez *et al.* [11] propose the globalized probability of boundary (gPb), which utilizes the boundary-preserving characteristic of NCut with sophisticatedly designed features to detect object boundaries and incorporate it into the oriented watershed transform and ultrametric contour map (OWT-UCM) to conduct image segmentation. This approach becomes the main support of many subsequent segmentation approaches [31, 4, 32, 33, 12, 13]. Kim *et al.* [34] formulate a hypergraph-based model and perform correlation clustering for image segmentation. Recently, Yu *et al.* [7] minimize an ℓ_1 -normed energy function of NCut to obtain piecewise smooth embeddings for gPb-owt-ucm [11], which obtains state-of-the-art image segmentation performance. However, all these methods suffer from expensive computations for feature extraction or optimization. Speed issues are considered in several following work. Multiscale combinatorial grouping (MCG) [12] segment images with multi-scale UCMs and it uses more advanced edge detection methods [35] to largely reduce the computation bottleneck of NCut used by gPb-owt-ucm. Chen [13] provides a solution to the scale-alignment in MCG. However, all these methods suffer from expensive computations for feature extraction or optimization. Sometimes several minutes are required to process a single 321×481 image, which significantly limits their practical usages. On the other hand, Pont-Tuset *et al.* [36] propose a downsampled approximating algorithm to accelerate the graph partitioning and use richer information in multiscale UCMs. Taylor *et al.* [37] and many others [4, 34] reduce the size of the affinity matrix using superpixel techniques. In this paper, we develop a method that is much faster than the aforementioned methods with competitive accuracy.

Edge detection plays an extreme importantly role in region based image segmentation [38, 39, 40, 41, 42]. For example, Convolutional Oriented Boundaries (COB) proposes an accurate boundary detection method using convolutional neural networks

(CNNs) and combines with [36] to perform image and object segmentation. Another popular image segmentation direction is semantic segmentation. Current methods use CNNs [43, 44, 45] to predict the semantic label of each pixel. These methods rely on large-scale training data. In contrast, our method aims at partition images into regions that can accurately segment objects from an image by observing its internal statistics in an unsupervised manner.

Designing feature to build the affinity between pixels/regions is important. Several studies have explored different cues, such as sophisticated combination of mixed image features [11], texture information [46], or saliency [47]. Different from these low-level image features, we argue that high-level cues are equally important and sometimes even more effective. For example, co-occurrence statistics have been used to capture the semantic object context knowledge based on training data to help the inference in, for example, condition random field (CRF) [48]. Different from this direction of research, our approach models region-wise co-occurrence probability based on pointwise mutual information [49] to build local connections of our proposed graph learned from the image itself.

Laplacian eigenmaps [50] computes a low-dimensional embedding to preserve the pairwise affinity of data points in the manifold. Local linear embedding (LLE) [51], alternatively, preserves the linear structure among the local neighboring points. The locality-preserving character of these two methods implicitly encourages the clustering of data. However, Isomap [52], which preserves global data geodesic distances, does not possess the nature of clustering. Our method shows a distinct point of view on the side of manifold learning to enhance spectral clustering for image segmentation.

3. Global-local Connected Graph Partitioning

In this section, we present the approach to build the local and global connections of the graph. Then we introduce the proposed energy function to partition the constructed graph.

3.1. Local connection with co-occurrence cues

Our proposed method begins with an over-segmentation with a set of regions, defined as $\mathcal{S} = \{S_1, \dots, S_N\}$. The over-segmentation is favorable considering its local spatial consistency and computational efficiency. Denote the graph by $\mathcal{G}_{local} = \{\mathcal{S}, W\}$, where $W \in \mathbb{R}^{N \times N}$ is the affinity matrix with each entry W_{ij} representing the affinity between regions S_i and S_j . W is sparse such that only spatially adjacent region pairs within a small range have nonzero values. Given a test image with large appearance variations inside the object (see Figure 2), a desirable affinity matrix should be able to discover the strong affinity between two visually different neighboring regions belonging to the common object. However, it is difficult for low-level features to achieve this goal because of their limitations in learning high-level knowledge.

One type of high-level knowledge comes from the fact that a neighboring region pair residing on an object is more likely to co-occur (i.e., have a high joint probability) due to the color patterns inside the object [31], such as the strip pattern on the clothes of the images in Figure 2. If we treat regions $\{S_i | i = 1, \dots, N\}$ as random variables, we can define the co-occurrence of two regions as

$$CO(S_i, S_j) = \log \frac{1}{A} P(S_i, S_j), \quad (1)$$

where $P(S_i, S_j)$ is the joint probability over S_i and S_j . Let $A = P(S_i)P(S_j)$ represent a normalization term, which is crucial to penalize the biased-high $P(S_i, S_j)$ of background region pairs against foreground object region pairs, because the background area usually has larger proportion than foreground objects. This normalization term will eliminate this unbalance accordingly. In addition, CO also contains information about object boundaries, because a region pair across the object boundary is a small-probability event [31].

We estimate $P(S_i, S_j)$ and marginal distribution $P(S_i)$ by using a nonparametric kernel density estimator [53] following [31]. But differently, we densely sample region pairs of each region and its adjacent regions within a certain (denoted as e_1) distance apart without repetition (which means $P(S_1, S_2) = P(S_2, S_1)$). Basically, we place estimator kernels on all regions $\{S_i\}$, and compute the image feature (gray values) co-

occurrence probability over all region pairs. So for each feature value pair, we have a co-occurrence frequency. Then we can simply normalize them and obtain $P(S_i)$ and the final co-occurrence cue $CO(S_i, S_j)$.

Our approach shares some similarities with [31] (denoted as PMI) for using point-wise mutual information, but is different from PMI in several perspectives. PMI interests in low pixel-wise joint probability to discover the rare boundaries, but we are interested in high region-wise probabilities and simultaneously maintain the boundary detection ability of PMI. PMI relies on raw image pixels, the probability in Eq. (1) is estimated over limited number of randomly sampled pixels. We rely on coherent regions to estimate this probability over most of adjacent region pairs, which yields probability distribution estimation closer to the actual distribution for the regions, and the estimation process is much less computational expensive.

Energy function: The first term E_{local} in our proposed energy function will encourage frequently co-occurring region pairs to be clustered into a group, and vice versa. Minimizing E_{local} is defined as the following:

$$\min_{\mathbf{y}} \sum_{i=1}^N \sum_{j=1}^N \|y_i - y_j\|^2 W_{ij}, \quad s.t. \mathbf{y}^T D \mathbf{y} = 1, \quad (2)$$

where D is a diagonal matrix and its i -th diagonal element is $d_{ii} = \sum_j W_{ij}$. The constraint is the key to normalize the cut of the graph. Minimizing E_{local} enforces y_i and y_j to take a similar value when W_{ij} is large. $\mathbf{y} = [y_1, \dots, y_N]^T$ is a real-valued vector, which is interpreted as a binary graph partition in Ncut or an one dimensional embedding in Laplacian eigenmaps. W_{ij} is defined as

$$W_{ij} = \exp \left(\sum_o CO(S_i^{f_o}, S_j^{f_o}) \right), \quad (3)$$

where the superscript f_o specifies a feature representation of the corresponding region. For each region, we calculate the pixel mean of Lab color space and the diagonal values of the RGB color covariance matrix in a 3×3 window. W_{ij} is computed between S_i and S_j within a certain distance apart, denoted as e_2 ($e_2 > e_1$).

The affinity matrix W is designed to measure the similarity between spatially adjacent region pairs based on their latent co-occurrence statistics. In order to preserve the

ignored relationships among spatially distant regions in the common object, we propose an additional energy term by building the global connections of the graph in the following section.

3.2. Global connection with saliency cues

The graph associated with global connections is denoted by $\mathcal{G}_{global} = \{\mathcal{S}, K\}$. Our approach strengthens the locality-preserving character by discovering the underlying linear structures among spatially distant regions (i.e., each region can be linearly represented by several neighboring regions so that the global connections are directed) belonging to a common object, while these regions are adjacent on a certain manifold. This goal is achieved by minimizing the second energy term E_{global} :

$$\min_{\mathbf{y}} \sum_{i=1}^N \|y_i - \sum_{j \neq i} K_{ij} y_j\|^2, \quad s.t. \mathbf{y}^T \mathbf{y} = 1, \quad (4)$$

where $K \in \mathbb{R}^{N \times N}$ is the coefficient matrix with R non-zero entries in each row to specify the linear combination coefficients of the representing neighbors. The constraint avoids degenerated solutions. \mathbf{y} is interpreted as an embedding in the original locally linear embedding (LLE) [51] method. Note that both \mathbf{y} and K are unknown; minimizing this energy function consists of three steps: 1) finding R neighbors for each region, 2) computing coefficient matrix K , and 3) computing \mathbf{y} .

Geodesic distance based neighbors: For each region, we consider its candidate neighbors from all regions within a large range of the defined geodesic distance, such that the distance (s_{ij}) between regions S_i and S_j is defined as follows:

$$s_{ij} = \left(\min(|w_{S_i} - w_{S_j}|, I_w - |w_{S_i} - w_{S_j}|)^2 + \min(|h_{S_i} - h_{S_j}|, I_h - |h_{S_i} - h_{S_j}|)^2 \right)^{\frac{1}{2}}, \quad (5)$$

where w_{S_i} and h_{S_i} denote the spatial x - and y -coordinates of region S_i , respectively. I_w and I_h denote the width and height of the image, respectively. Intuitively, this metric treats the image as if it was wrapped along its four corners into a sphere and describes the geodesic distance along this resulting surface. The measurement can trace the connections of the regions belonging to foreground objects or background with an arbitrary shape and range.

For a region S_i , we select R nearest regions with each region represented as a feature vector calculated by a saliency cue mapping σ . Then we find its coefficients K_i by

$$\min_{K_i} \|\sigma(S_i) - \sum_{j \neq i} K_{ij} \sigma(S_j)\|^2 + \alpha \text{Tr}(K_i^T K_i), \quad s.t. \quad \sum_j K_{ij} = 1. \quad (6)$$

The regularization term is necessary to prevent ill-conditioned solutions when neighboring regions have similar feature values (i.e., making the Gram matrix singular). The regularization parameter is chosen as $\alpha = 1e-10$. The constraint ensures the translation invariance.

Saliency cue complying linearity: Spatially distant regions inside the same object may have large appearance variances, for example, the face, hairs, and clothes of a person exhibit totally different appearances (see Figure 2). Therefore, it is difficult to measure their latent similarity with traditional cues. However, those visually different regions usually exhibit similar saliency degree in the human visual system [54]. This characteristic remedies the ‘‘imperfection’’ of pairwise co-occurrence affinity and satisfies the requirement to build global connections. We take advantage of the empirical knowledge that salient objects in images have distinctive colors from the background under a certain linear combination of mixed color spaces [55]. To this end, we choose RGB, Lab, and hue and saturation channels of HSV (8 channels) and their nonlinear transformations with gamma correction (with three gamma values, [0.5, 1.5, 2.0]) to consider the human vision’s nonlinear responses, thereby yielding a 24-dimensional feature vector for each region, $\sigma : S \mapsto \mathbb{R}^{24}$.

Our method to incorporate saliency in the graph connection is elaborate. Unlike saliency detector [55], we do not compute the coefficient explicitly based on any supervised information. Since the correlation of each region feature vector $\sigma(S)$ is consistent under arbitrary linear transformation, its saliency characteristic between regions will be implicitly expressed in Eq. (4).

3.3. Proposed energy function to partition graph

Overall, the proposed full graph is defined as $\mathcal{G} = \{S, (W, K)\}$, where W specifies the local undirected connections and K specifies the global directed connections. Our



Figure 2: The first row shows two sample images exhibiting large object internal appearance variances. The second and third row shows the graph partitioning results of NCut and ours, respectively. As can be observed, NCut preserves the edge information but segments images into the connected components (middle). Our approach is able to separate the entire object from its background (bottom).

goal is to pursue global partitioning of the graph \mathcal{G} , i.e., minimizing the two energy terms simultaneously. Therefore, the energy function E can be defined and derived as follows (detailed derivations are skipped):

$$\begin{aligned}
 E &= E_{local} + \mu E_{global} \\
 &= \sum_{i=1}^N \sum_{j=1}^N \|y_i - y_j\|^2 W_{ij} + \mu \sum_{i=1}^N \|y_i - \sum_{i \neq j} K_{ij} y_j\|^2 \\
 &= \mathbf{y}^T (D - W + \mu M) \mathbf{y},
 \end{aligned} \tag{7}$$

where $(D - W) \in \mathcal{R}^{N \times N}$ is the Laplacian matrix and $M = (I - K)^T (I - K) \in \mathcal{R}^{N \times N}$. μ is a regularization parameter to balance E_{global} and E_{local} . How to select the optimal value of μ is discussed in the experimental section.

It is straightforward to see that minimizing E is to solve a generalized eigenvector

system:

$$(D - W + \mu M)\mathbf{y} = \lambda D\mathbf{y}, \quad (8)$$

which produces a set of eigenvectors \hat{Y} , where each column is an eigenvector representing a binary partition of the graph. In practice, the number of segments of an arbitrary test image is unknown and the expected partition is not guaranteed to be the eigenvector associated with the second smallest eigenvalue [6]. In Section 4, we present a novel approach to address this issue for multi-class segmentation.

Leverage E_{local} and E_{global} : The two energy terms are designed for different purposes. E_{local} preserves the pairwise similarity of spatially adjacent region pairs, while E_{global} preserves the linear structure of spatially distant regions in the common object. Both emphasize the locality-preserving for the purpose of clustering (or graph partitioning). Compared with the hard constraint of E_{local} , E_{global} encourages soft (i.e., likelihood) clustering of the regions [14]. In Figure 2, we visualize several graph partition results of the proposed approach and compare it to NCut. As can be observed, the significantly improved graph partitioning quality demonstrates the effectiveness of the global connections introduced in E_{global} .

4. Multi-class Segmentation

In this section, we introduce the approach to use graph partitions for multi-class segmentation.

4.1. EigenHistogram

We have computed a set of eigenvectors (i.e., image partitions) $\hat{Y} = [\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_d] \in \mathbb{R}^{N \times d}$ corresponding to the first d smallest eigenvalues (excluding the zero eigenvalue) using Eq. (8). The i -th region can now be represented as a d -dimensional vector $S_i^{\hat{Y}}$. The k -means algorithm is applied to group regions into L segments, $\mathcal{Z} = \{\mathcal{Z}_k\}_{k=1}^L$, to produce a hard partition [4, 6]. To obtain more reliable multi-class segmentation that can be generalized to arbitrary images with different number of classes, we treat it as a prior segmentation to provide the class likelihood for the multi-class segmentation.

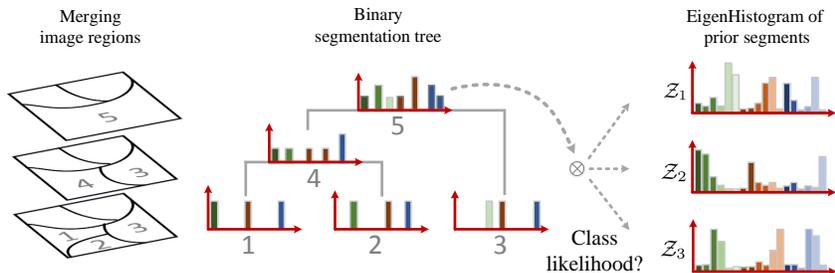


Figure 3: In the merged binary segmentation tree, each leaf node representing an initial image region has an EigenHistogram while the EigenHistograms of the inner nodes are accumulated and normalized from that of their descendant nodes. Each node computes its class likelihood based on the EigenHistograms of the prior segmentation \mathcal{Z} .

Note that our method can deal with the number of segments regardless of pre-defined L . We will discuss this in the experiments.

For each dimension of $S_i^{\hat{Y}}$, we compute a histogram with L bins uniformly spaced between $[0, 1]$ based on the corresponding normalized eigenvector. As a consequence, a region will be represented as a $(d \times L)$ -dimensional concatenated histogram (we set $d=6$ empirically and we will discuss the parameter L in the experimental section). For each segment \mathcal{Z}_k , we accumulate and normalize the histograms of all regions belonging to this segment. We term this region representation as EigenHistogram (see Figure 3).

4.2. Multi-class segmentation

Considering the image regions as a random field, we are interested in incorporating the unary potential (the class likelihood) of each region based on the prior segmentation \mathcal{Z} and pairwise potentials between neighboring regions into a unified energy function, to achieve a holistic multi-class segmentation. Numerous literatures have investigated to learn effective unary potentials for random field based algorithms via structured support vector machine [56, 17] or convolution neural network [57, 58] to perform semantic segmentation. In contrast, EigenHistogram can be treated as a high-level representation which possesses spatial consistency, thereby intrinsically scalable to image segments of arbitrary size. Furthermore, it is easy and fast to compute without any supervision as other methods [56, 17, 57] conduct.

Following the Pylon model [17], we can configure the regions into a hierarchical binary segmentation tree. Different from the traditional “flat” random field models [59, 2], each node in our tree structure stands for a region nested from bottom to top, which enables the features to be extracted at different levels of the hierarchy to enrich the feature representation of the segments. In total, the constructed tree has $2N+1$ regions (the root node is the whole image), $\mathcal{S}^+ = \{S_i | i = 1, \dots, 2N+1\}$, which define a hierarchical random field.

Our goal is to assign labels $\mathbf{p} = [p_1, \dots, p_{(2 \cdot N+1)}]^T$ to all regions in \mathcal{S}^+ . Therefore, we minimize the following object function:

$$\Theta(\mathbf{p}) = \sum_{i=1}^{2 \cdot N+1} U(p_i) + \sum_{(i,j) \in \mathcal{N}(\mathcal{S})} B(S_i, S_j), \quad p_i \in \{0, 1, \dots, L\}, \quad (9)$$

where $U(p_i)$ is the unary potential of the region $S_i \in \mathcal{S}^+$ to specify the cost of assigning label p_i to S_i , and $B(S_i, S_j)$ is the pairwise potential to specify the boundary cost (exponentiated boundary strength [56]) between any two neighboring regions $(i, j) \in \mathcal{N}(\mathcal{S})$ in the child nodes, which is used to encourage the spatial smoothness. Note that p_i is allowed to take a zero label such that it satisfies the non-overlapping requirement [17] by using the constraint:

$$\forall i \neq j, S_i, S_j \in \mathcal{S}^+, \text{ if } S_i \cap S_j \neq \emptyset, \text{ then } p_i \cdot p_j = 0, \quad (10)$$

which ensures that any subtree can have only one single non-zero label.

Since we have clustered image regions into L segments, the unary potential of region S_i assigned to the k -th segment has the cost:

$$U_{p_i=k} = -\beta \cdot \langle \mathcal{H}(S_i), \mathcal{H}(\mathcal{Z}_k) \rangle, \quad S_i \in \mathcal{S}^+, \mathcal{Z}_k \in \mathcal{Z}, \quad (11)$$

where U_{p_i} is the unary potential of the region $S_i \in \mathcal{S}^+$ to specify the cost of assigning label p_i to S_i . β determines the weight of the unary potential against the pairwise potential. \mathcal{H} transforms a region into the EigenHistogram representation, where the class likelihood is calculated for each region in the tree. Following [17], we compute the pairwise potentials as the exponentiated boundary strength.

EigenHistograms of the internal nodes of the binary segmentation tree are accumulated and normalized from that of the corresponding descendant nodes (see Figure 3).

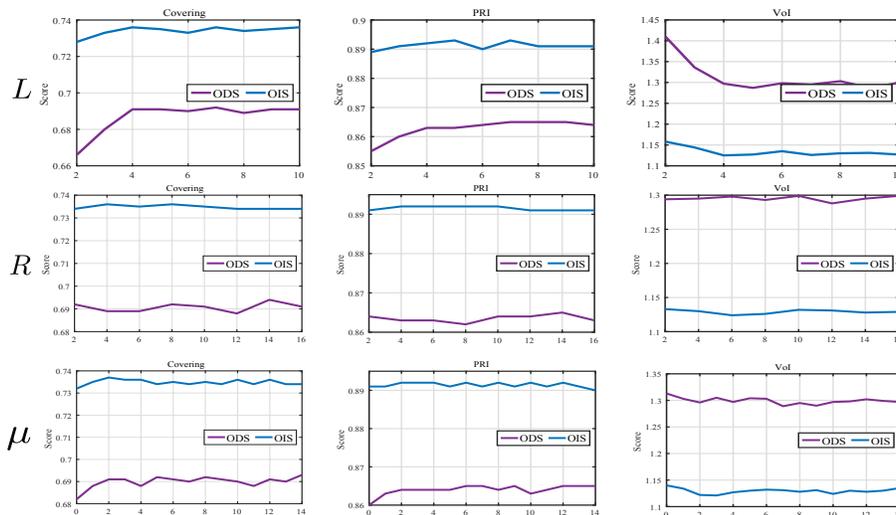


Figure 4: Segmentation evaluation with different parameters on the training dataset. The top, middle, and bottom rows show the segmentation results by varying L , R , and μ , respectively, as indicated in x -axis. It can be observed that the proposed method is insensitive to variations of the parameters.

Therefore, the partial non-smoothness effects of the eigenvectors (i.e., isolated regions as visualized in the right panel of Figure 2) reflected in the EigenHistograms of top-level nodes will be suppressed. Finally, we can leave the rest computation to the whole inference procedure to produce a holistic multi-class segmentation as the final output, by using the alpha-expansion based graph cut [2].

5. Experimental Results

This section evaluates the segmentation performance of the proposed method. We first analyze the parameter settings. Then we evaluate and demonstrate the segmentation results, and compare to several state-of-the-art methods.

We mainly evaluate the proposed segmentation approach using the challenging Berkeley Segmentation Dataset (BSDS500) [11]. BSDS500 is widely used as the benchmark for image segmentation and boundary detection, which contains 200 training, 100 validation, and 200 test images. We use several standard evaluation criteria [11] to conduct quantitative analysis: Segmentation Covering, Probability Rand Index

Table 1: The comparison of segmentation results and runtime on the BSDS500 dataset.

Method	Covering		PRI		VoI		Time(s)
	ODS	OIS	ODS	OIS	ODS	OIS	
NCut [6]	.45	.53	.78	.80	2.23	1.89	-
Felz-Hutt [1]	.52	.57	.80	.82	2.21	1.87	-
Mean Shift [60]	.54	.58	.79	.81	1.85	1.64	-
ISCRA [61]	.59	.66	.82	.85	1.60	1.42	30
gPb-owt-ucm [11]	.59	.65	.83	.85	1.69	1.48	240
cPb-owt-ucm [4]	.59	.65	.83	.86	1.65	1.45	>240
red-spectral [37]	.56	.62	.83	.85	1.78	1.56	~12
DC-Seg [33]	.58	.63	.82	.85	1.75	1.59	6
DC-Seg _{full} [33]	.59	.64	.82	.85	1.68	1.54	144
PMI _{low} [31]	.61	.66	.83	.86	1.58	1.42	30*
MCG [12]	.61	.66	.83	.86	1.57	1.39	18
PFE+ucm [7]	.61	.66	.83	.86	1.64	1.46	>900·b [†]
PFE+MCG [7]	.62	.68	.84	.87	1.56	1.36	>900·b [†]
Ours	.62	.66	.83	.86	1.59	1.43	9

*Time is tested on half-sized images. [†] $b=\{4, 8, 16\}$ is the number of the embedding needs to compute.

(PRI), and Variation of Information (VoI), which measure per-pixel segment overlapping, pairwise pixel matching, and segmentation-wise entropy, respectively. For each measurement, we report the values with the optimal dataset scale (ODS) and optimal image scale (OIS). We further evaluate our method on large-scale PASCAL VOC and COCO datasets to show the generalization ability of our method for object segmentation and compare to two state-of-the-art methods.



Figure 5: Segmentation results on BSDS500. The first and second columns show the test images and the ground truth, respectively. The third to the sixth columns show the results obtained by gPb-owt-ucm [11], DC-Seg_{full} [33], MCG [12] and our method, respectively. All results are visualized with the optimal scale (ODS) of the corresponding methods used for quantitative evaluation. Figure 6 presents one graph partition result of the proposed method.

5.1. Implementation details

We investigate the parameter sensitivity of the proposed method and select the optimal values based on the training set. Then we apply these values to the independent

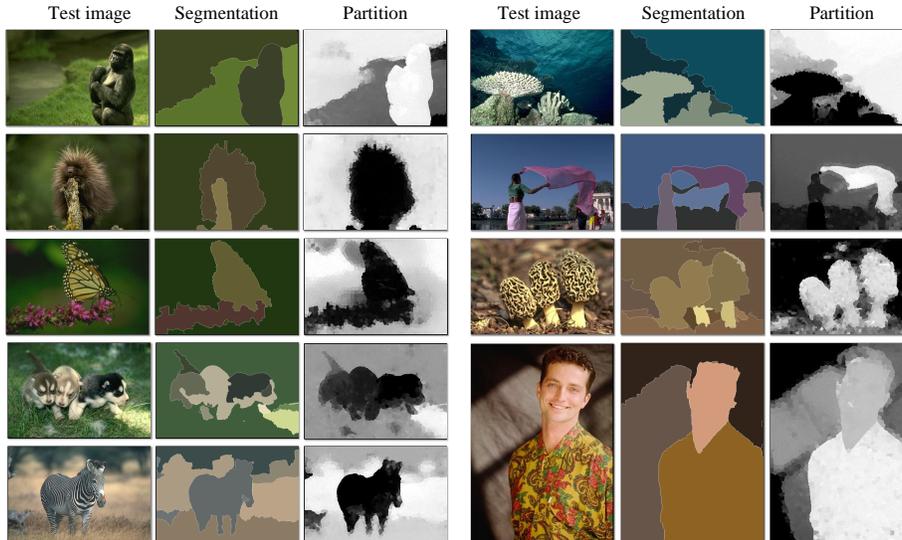


Figure 6: The graph partition results of the proposed method on the test images in Figure 5. As can be observed, our graph partition is able to segment the objects from the background on the challenging images, where the objects have similar color or texture with that of background.

test set.

Figure 4 shows the performance of the proposed method with respect to L for clustering, R for selecting nearest regions in Eq. (6), and μ for graph partitioning. The selection of the optimal value of the number of segments L is dependent on the test set, but we do not select the best L based on the test set, by which we aim to demonstrate the strong generalization ability of the proposed method. For the parameter μ , compared with $\mu = 0$, which means that only E_{local} is considered in the energy function, E_{global} with $\mu = 8$ improves the accuracy by $>.10DS$ (Covering). Section 5.4.2 further validates the effectiveness of E_{global} . As can be observed, the proposed method is insensitive to the three parameters. As a result, we set $L = 6$, $\mu = 8$, and $R = 14$ throughout the following experiments.

We empirically set $e_1 = 20$ for kernel density estimation in Eq. (1) and $e_2 = 40$ for computing W in Eq. (3). Since the test set has approximately equal image sizes, we can assume that these two values can be generalized to all test images. We empirically

Table 2: The running time (in second) of each phase of the proposed method.

Phase	Min	Max	Mean	Var.
1: Region structure generation	2.60	3.73	3.08	0.05
2: Graph construct. and partition	3.30	6.89	4.51	0.56
3: Multi-class segmentation	1.39	2.26	1.71	0.03
Total	7.54	12.6	9.30	1.00

found that the parameter β in Eq. (11) varies from images to images. In practice, we run the inference procedure to obtain multiple segmentations by varying the β value between $[200, 300, \dots, 800]$ with an interval equal to 100, and take the average of all these outputs and the superpixel map as the final segmentation.

We use the toolbox provided by Dollár *et al.* [35] to generate the superpixel map (i.e., the structured edge (SE) detector followed by UCM) with roughly uniform region sizes. Our implementation is based on Matlab running on a standard Intel i7 desktop.

5.2. Segmentation result comparison

We evaluate the performance and efficiency of the proposed method, and compare it to several state-of-the-art methods.

In Table 1, we compare the proposed approach to several state-of-the-art methods in terms of segmentation accuracy and running time on the BSDS500 test set. As one can see, the proposed method significantly outperforms most of the comparative methods. PMI_{low} [31] is a boundary detection method, which embeds the edge map into OWT-UCM [11] to obtain accurate segmentation. We report its the best accuracy, which is achieved on low resolution images. The recently proposed multiscale combinatorial grouping (MCG) [12] and piecewise flat embedding (PFE) [7] obtain significant improvement compared with the early method, such as red-spectral [37] and DC-Seg [33] (see Table 1). MCG uses hierarchical UCMs to boost the segmentation performance. PFE integrates its computed graph partitions into the gPb-owt-ucm [11] and MCG, which achieves good segmentation performance. However, PFE suffers from the computationally expensive optimization. The proposed method outper-

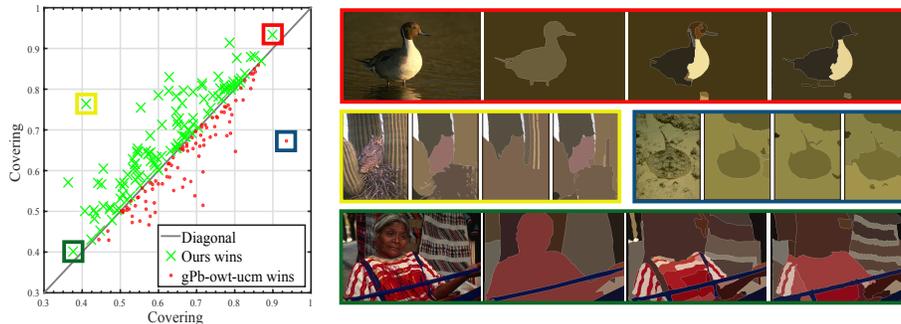


Figure 7: Pairwise segmentation comparison between gPb-owt-ucm [11] and our method on the 200 test images. In the left figure, \times above the diagonal indicates the image our method wins, and \bullet below the diagonal indicates gPb-owt-ucm wins (best viewed in electronic form). We also highlight several extreme cases corresponding to two with the largest discrepancy and the worst and the best cases. Each is marked with an unique colored rectangular. In each rectangular of the right slide containing four images, from left to right, shows the test image, ground truth, the gPb-owt-ucm result, and our result, respectively.

forms PEF+owt-ucm and it achieves close segmentation performance compared with PFE+MCG. More importantly, the proposed method is hundreds of times faster than the PFE based methods. DC [33] and red-spectral [37] also emphasize on fast image segmentation, but their segmentation accuracy is not as accurate as ours.

Figure 5 shows the qualitative segmentation results. Figure 6 presents the graph partitioning result obtained by the proposed method, which provides good initial segmentation proposals. Compared with other methods, the proposed method is able to resist the object internal variances to avoid small segments, so that the segments are much more spatially consistent. In addition, the proposed method can implicitly figure out the best number of segments regardless of the pre-defined L value. It is because that Eigen-Histogram can penalize over-segmentation since homogeneous segments have similar EigenHistogram and thus proximate unary potentials, encouraging them to be merged. The first three rows of Figure 5 particularly highlight the above-mentioned capability. To provide more detailed comparison, in Figure 7, we show the pairwise segmentation results obtained by our proposed method compared to the classical gPb-owt-ucm [11] method. As can be observed, the proposed method shows obvious improvement on a large number test images.

Time Efficiency: The comparison of running time is shown in the rightmost column of Table 1. The test image size is 321×481 . The proposed method is much faster than other competing methods because of several important aspects:

1. The proposed method does not need complex feature computation, which is superior than gPb based methods [7, 11].
2. We construct the graph model based on superpixels rather than raw pixels. Although we incorporate multiple cues into a graph with complicated constraints, the graph partitioning is a single eigenvector system. While in the PFE method [7], performing graph partitioning is particularly computationally expensive.
3. EigenHistogram is efficient to compute and very scalable to regions with arbitrary size for hierarchical multi-class segmentation.

The proposed method is executed in three phases: 1) generating the superpixel map and constructing the hierarchical segmentation tree, 2) constructing and partitioning the graph, and 3) conducting multi-class segmentation. Given an $H \times W$ resolution image, phase 1 takes low logarithmic time of random forest tree depth to predict edge map with a random forest, $O(HW + N)$ to compute superpixels with N regions, and $(\log N)$ to construct the hierarchical binary tree. Phase 2 takes up to a factor of $O(HW + N)$ to compute all image features with respect to pixels and regions and approximately $O(fN^2)$ to compute the affinity matrix, where f is the feature dimension. Since $(D - W + \mu M)$ in Eq. (8) is sparse, solving the eigen decomposition problem with a $N \times N$ affinity matrix takes $O(N(\tilde{R} + R))$ using a Lanczos algorithm according to [6], where $\tilde{R} + R \ll N$ is the adjacencies from both local (\tilde{R}) and global connections (R) in the graph. Phase 3 optimizes Eq. (9) with L classes with approximately $O(N^2L)$ using graph cut with alpha-expansion [59, 17]. Therefore, the overall method has approximate time complexity $O(HW + fN^2)$, bounded by phase 2.

We also evaluate the detailed running time of the proposed method on the 200 BSDS500 test images. Table 2 shows the detailed time cost of each phase. Compared with most comparative methods, the proposed method is more scalable for practical usages.

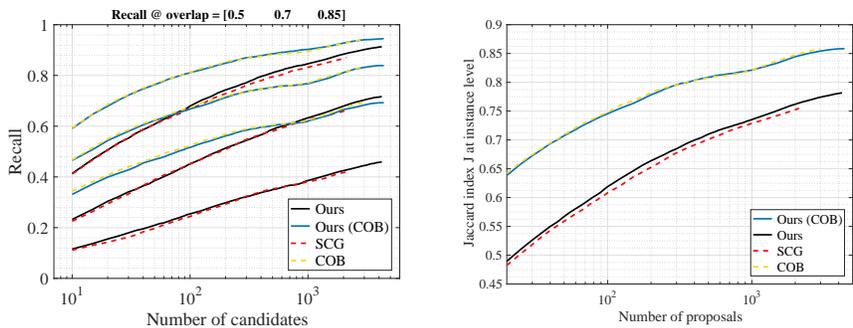


Figure 8: Object segmentation evaluation on PASCAL VOC 2012. The different lines in the same color indicates results under different Jaccard overlapping thresholds. [36] introduces detailed evaluation metrics.

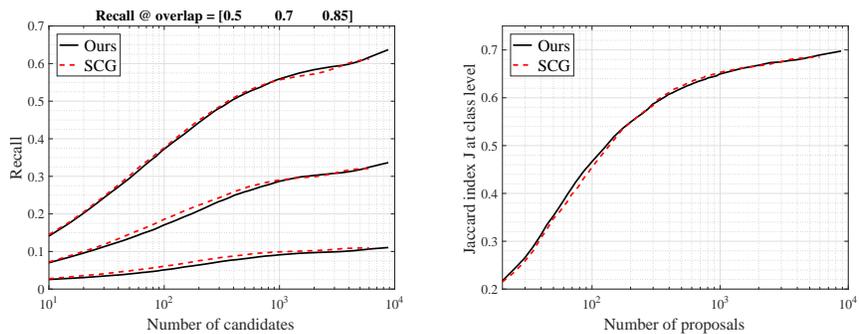


Figure 9: Object segmentation evaluation on COCO. See text for explanations.



Figure 10: Qualitative evaluation of the segmentation on COCO. Our method (right column) generates obviously clearer segmentation results with significantly reduced over-segmentation (e.g. the flower and vase in the left images and swimwear in the right images) than the comparing method (middle column) [36].

5.3. Towards large-scale object segmentation

This section further demonstrates our proposed method on the large-scale PASCAL VOC [18] and COCO [19] segmentation datasets¹. Since our method generates region segmentation composed by a set of connected regions (the same as UCMs), we can fully use our method to generate object proposals by training an object proposal grouping classifier following [36]. We closely follow its training procedure and evaluation settings. In brief, Jaccard Index J , i.e. the size of the intersection of the pixel union of two regions, is used to evaluate the accuracy of generated objects compared with groundtruth.

Figure 8 shows the comparing results for PASCAL VOC. We compare with a method proposed by [36], denoted as singlescale combinatorial grouping (SCG). As can be observed on the two evaluation metrics, our method improves the performance of SCG on the recall evaluation metrics consistently. We also compare with a recent deep learning based method COB [39] which aims at detecting accurate object boundaries. It combines with MCG [36] to perform object segmentation and achieved significant improvement. Note that region segmentation highly relies on the quality of boundary detection (it is out of the focus of this paper). As will demonstrated in Section 4.1, our method is flexible to be an extension of arbitrary baseline methods. Hence, we use the edge maps generated by COB (denoted as Ours (COB)). As can be observed, our method improves a substantial margin compared with our original method and achieves competitive results compared with COB. Figure 9 compares the results on COCO. Our method shows better results than SCG (right) at low numbers of proposals and competitive results on the recall with respect to the number of candidates.

SCG is designed for generate image object candidates, so its generated UCMs contain very fine and small region segments, which is an advantage when computing evaluation metrics for images with multiple objects. However, our method does not have

¹According to the experiment settings of [36], for PASCAL, 1,464 training images and 1,449 validation images are used. COCO totally contains 82,783 training images and 40,504 validation images in total. In our experiments, we randomly select 5,000 and 2,500 from the training and validation set, respectively for evaluation.

Table 3: The segmentation results under different configurations. Different region generation baselines which are used by our method are indicated in (\cdot). The last row is the result of our method (SE+ucm) when E_{global} is not applied. Please see text for detailed explanations.

Method	Covering		PRI		VoI	
	ODS	OIS	ODS	OIS	ODS	OIS
SemiContour+ucm	.56	.63	.82	.85	1.79	1.57
Ours (SemiContour+ucm)	.60	.64	.83	.85	1.68	1.50
MCG	.61	.66	.83	.86	1.57	1.39
Ours (MCG)	.62	.66	.84	.86	1.57	1.40
SE+ucm	.59	.64	.83	.86	1.71	1.51
Ours (SE+ucm)	.62	.66	.83	.86	1.59	1.43

designs for this goal. Compared with it, our method is significantly more proficient at segmenting the salient objects in images. We will further analyze this behavior in the next section. The PASCAL dataset is mainly collected for image and object segmentation tasks. According to our observation, PASCAL images usually contain definite and salient objects. Therefore, our method performs better and largely improves SCG. While in COCO, most of images are outdoor scenes that usually contain many small and indefinite objects. That is the reason why the improvement on COCO for our method is not as large as that in PASCAL, compared with SCG. We qualitatively compare with SCG on COCO images with relatively definite objects. As can be seen in Figure 10, our method can significantly reduce over-segmentation and give rises to clearer segmentation results. Nevertheless, the shown results on the two large-scale datasets are sufficient to demonstrate the generalization ability of our proposed method to different datasets with diverse scenes ².

²Note that we did not select the parameters of our method on the targeting training datasets following Section 5.1 but used the unique one selected using the BSDS500 training set. We believe there is still room for improvement with careful fine-tuning.

Components				Covering	
E_l	E_g	MS	kM	ODS	OIS
✓			✓	.43	.50
✓	✓		✓	.52	.52
✓		✓		.60	.65
✓	✓	✓		.62	.66

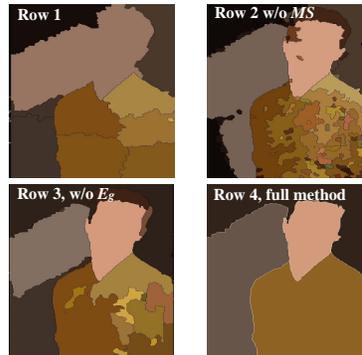


Figure 11: Ablation study to analyze the effectiveness of each component of the proposed method. **Left:** Each row shows the result of a combination of the components. E_l , E_g , MS , and kM denote E_{global} , E_{local} , multi-class segmentation, and k-means, respectively. To evaluate the effectiveness of MS , we simply use k-means to cluster eivenvectors (11 classes), as an alternative to perform MS . The 4th row is our full method. **Right:** The qualitative results corresponds to each row of the left side table. As can be observed from both quantitative and qualitative results, the proposed E_g and MS components play important roles in generating clear segmentation and better scores.

5.4. Analysis

5.4.1. Serving as an extension to improve baseline methods

We consider the cases of using different methods to generate superpixel maps as the input of the proposed method, which allow us to conduct more detailed analyses. It is necessary to notice that, although the proposed method is flexible to build upon these methods, it is not an extension of the underling methods. In contract, the proposed method is a new exploration of accurate and fast spectral clustering based image segmentation. In addition, many state-of-the-art methods use accurate supervised edge detectors and other trained classifiers [12, 13]. We are particularly interested in reducing number of training data with an aim to completely unsupervised image segmentation. Either unsupervised [62] or semi-supervised SE detector [41] can be used as the underlying edge detectors. We consider using the latter, namely Semi-Contour [41] (3 training images are used), as an alternative to the originally used SE. We compare the performance in Table 3. The obtained segmentation results consistently improve the segmentation accuracy of different baseline methods. Particularly, we observe .4ODS (Covering) improvement over SemiContour+ucm and .3ODS im-

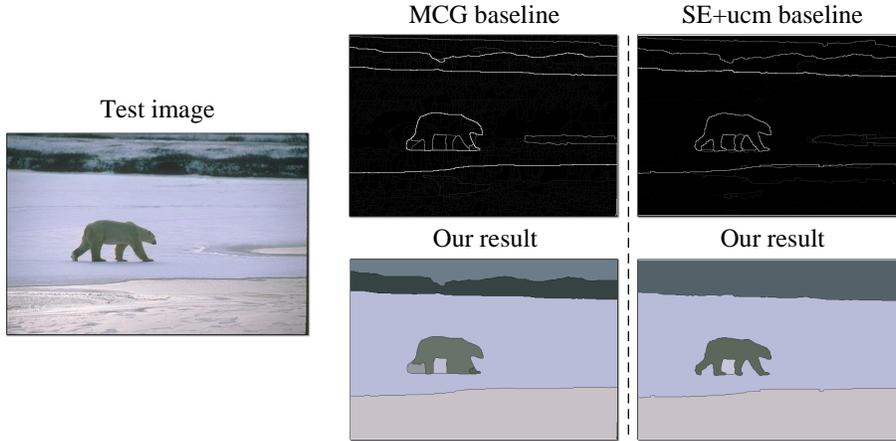


Figure 12: The segmentation results (shown in the second row of right slide) of the proposed method which using MCG and SE+ucm to generate superpixel maps as candidate regions (the first row), respectively. Note that MCG sharpens noisy edges below the bear (upper left), resulting in worse segmentation (bottom left).

provement over SE+ucm.

5.4.2. Ablation study

We analyze the effectiveness of each component of the proposed method. The proposed global connection (Section 3.2) is very effective at capturing the affinity between spatially distant regions belonging to the same objects. And the proposed multi-class segmentation is critical to generate smooth and clear segmentation map and makes our method robust to arbitrary images. Figure 11 evaluates each components both qualitatively and quantitatively. Comparing with our method without using E_{global} , we observe obvious improvement (comparing the 3rd row against 4th row and the 1st row against the 2nd row), which indicates the effectiveness of the proposed energy term E_{global} . To validate our multi-class segmentation, we conduct an experiment by simpling clustering the generated graph partitions (i.e. eigenvectors) using k-means to L classes and evaluate the performance. Simple hard clustering strategy can not adapt to arbitrary images with different number of classes and does not guarantee local smoothness, these two factors have large penalty on the evaluation metrics as shown in the first and second rows of Table 11. Therefore, we argue that our strategy to use eigenvectors

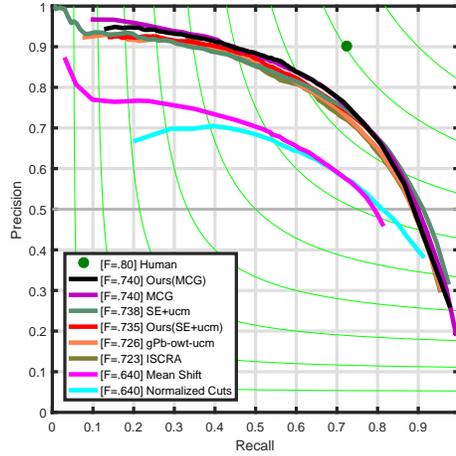


Figure 13: The Precision-Recall curves obtained by the proposed method and competing methods for boundaries on the BSDS500 dataset. (·) indicates the baseline method used by the proposed method.

for multi-class segmentation is very effective (as explained in Section 4.2).

5.4.3. Edge information

The improvement using MCG as the baseline is a small margin (i.e., .10DS) compared with cases of using the other two methods as the baselines. In fact, MCG uses SE to detect edges while it also sharpens edges. Nevertheless, we observe MCG sometimes sharpens irrelevant edges as well, such that the sharpened noisy edges will have a large penalization through pairwise potentials against unary potentials in our multi-class segmentation procedure, leading to undesirable results. Figure 12 illustrates this situation. The above results indicate that the proposed method relies less on strong edge information compared with MCG.

Additionally, since the proposed method relies less on edges, one potential weakness of the graph partitioning procedure could result in the fragmentation of homogeneous regions, which decreases the precision of the boundary detection. We compare the boundary precision-recall curve in Figure 13, from which we can see that the proposed method maintains nearly the same precision as the baseline methods, i.e., MCG and SE+ucm (though negligible 0.03 decrease for SE+ucm).

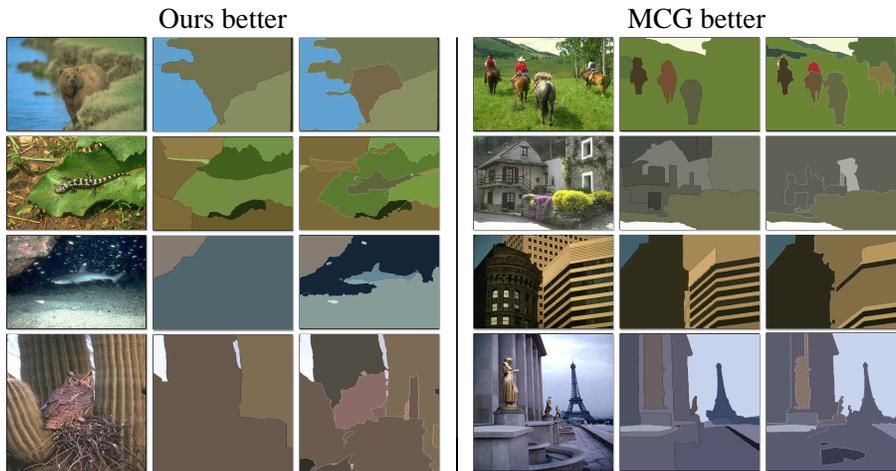


Figure 14: Each sample shows the segmentation results of our proposed method with SE+ucm baseline (right) and MCG (middle) [12]. The left side shows images that our method wins and the right side shows images that MCG wins.

5.4.4. Strengths and limitations

The proposed method is effective in discovering complex image knowledge among regions from challenging natural images and segmenting objects even when objects have weak boundaries. The proposed method is significantly better than MCG in those samples shown in Figure 14(left). However, we found that the proposed method is not that effective at images without definite objects, because our graph design emphasizes the high-level discriminative image knowledge of objects against the background. MCG outperforms ours in those samples (see Figure 14(right)).

6. Conclusions

In this paper, we present a fast yet accurate image segmentation method, which is a novel re-examination of spectral clustering based image segmentation for unsupervised image segmentation. We construct an image region graph with both local and global connections based on simple but effective high-level cues, and formulate the graph partitioning as a simple generalized eigenvector system. The high quality graph partitions are used to compute effective unary potentials of Pylon model for multi-class

image segmentation. Extensive experiments, on the BSDS500 benchmark, large-scale PASCAL VOC and COCO datasets, show that the proposed method achieves significantly faster speed and competitive performance when it is compared to state-of-the-art segmentation methods.

7. Acknowledgement

This work was partially supported by the National Natural Science Foundation of China under Grants U1605252, 61472334, and 61571379.

References

- [1] P. F. Felzenszwalb, D. P. Huttenlocher, Efficient graph-based image segmentation, *International Journal of Computer Vision* 59 (2) (2004) 167–181.
- [2] Y. Boykov, O. Veksler, R. Zabih, Fast approximate energy minimization via graph cuts, *Transactions on Pattern Analysis and Machine Intelligence* 23 (11) (2001) 1222–1239.
- [3] B. Peng, L. Zhang, D. Zhang, A survey of graph theoretical approaches to image segmentation, *Pattern Recognition* 46 (3) (2013) 1020–1038.
- [4] T. H. Kim, K. M. Lee, S. U. Lee, Learning full pairwise affinities for spectral segmentation, *Transactions on Pattern Analysis and Machine Intelligence* 35 (7) (2013) 1690–1703.
- [5] X. Shi, Z. Guo, Z. Lai, Y. Yang, Z. Bao, D. Zhang, A framework of joint graph embedding and sparse regression for dimensionality reduction, *IEEE Transactions on Image Processing* 24 (4) (2015) 1341–1355.
- [6] J. Shi, J. Malik, Normalized cuts and image segmentation, *Transactions on Pattern Analysis and Machine Intelligence* 22 (8) (2000) 888–905.
- [7] Y. Yu, C. Fang, Z. Liao, Piecewise flat embedding for image segmentation, in: *Proceedings of the International Conference on Computer Vision*, 2015, pp. 1368–1376.

- [8] S. X. Yu, Segmentation induced by scale invariance, in: Proceedings of the Conference on Computer Vision and Pattern Recognition, 2005, pp. 444–451.
- [9] M. Maire, S. X. Yu, Progressive multigrid eigensolvers for multiscale spectral segmentation, in: Proceedings of the International Conference on Computer Vision, 2013, pp. 2184–2191.
- [10] T. Cour, F. Benezit, J. Shi, Spectral segmentation with multiscale graph decomposition, in: Proceedings of the Conference on Computer Vision and Pattern Recognition, 2005, pp. 1124–1131.
- [11] P. Arbelaez, M. Maire, C. Fowlkes, J. Malik, Contour detection and hierarchical image segmentation, Transactions on Pattern Analysis and Machine Intelligence 33 (5) (2011) 898–916.
- [12] P. Arbelaez, J. Pont-Tuset, J. Barron, F. Marques, J. Malik, Multiscale combinatorial grouping, in: Proceedings of the Conference on Computer Vision and Pattern Recognition, 2014, pp. 328–335.
- [13] Y. Chen, D. Dai, J. Pont-Tuset, L. Van Gool, Scale-aware alignment of hierarchical image segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 364–372.
- [14] M. Belkin, P. Niyogi, Laplacian eigenmaps for dimensionality reduction and data representation, Neural computation 15 (6) (2003) 1373–1396.
- [15] A. Y. Ng, M. I. Jordan, Y. Weiss, et al., On spectral clustering: Analysis and an algorithm, in: Advances in Neural Information Processing Systems, Vol. 14, 2001, pp. 849–856.
- [16] X. Shi, Y. Yang, Z. Guo, Z. Lai, Face recognition by sparse discriminant analysis via joint l_2, l_1 -norm minimization, Pattern Recognition 47 (7) (2014) 2447–2453.
- [17] V. Lempitsky, A. Vedaldi, A. Zisserman, Pylon model for semantic segmentation, in: Advances in Neural Information Processing Systems, 2011, pp. 1485–1493.

- [18] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The pascal visual object classes (voc) challenge, *International journal of computer vision* 88 (2) (2010) 303–338.
- [19] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick, Microsoft coco: Common objects in context, in: *European conference on computer vision*, Springer, 2014, pp. 740–755.
- [20] S. X. Yu, J. Shi, Multiclass spectral clustering, in: *Transactions on Pattern Analysis and Machine Intelligence*, 2003, pp. 313–319.
- [21] B. Wang, Z. Tu, Affinity learning via self-diffusion for image segmentation and clustering, in: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, 2012, pp. 2312–2319.
- [22] J. Zhang, J. Zheng, J. Cai, A diffusion approach to seeded image segmentation, in: *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, 2010, pp. 2125–2132.
- [23] C. Rother, V. Kolmogorov, A. Blake, Grabcut: Interactive foreground extraction using iterated graph cuts, in: *ACM transactions on graphics (TOG)*, Vol. 23, 2004, pp. 309–314.
- [24] L. Grady, Random walks for image segmentation, *IEEE transactions on pattern analysis and machine intelligence* 28 (11) (2006) 1768–1783.
- [25] Y. Li, X. Feng, A multiscale image segmentation method, *Pattern Recognition* 52 (2016) 332–345.
- [26] S. Yin, Y. Qian, M. Gong, Unsupervised hierarchical image segmentation through fuzzy entropy maximization, *Pattern Recognition* 68 (2017) 245–259.
- [27] K. J. F. De Souza, A. de Albuquerque Araújo, Z. K. do Patrocínio, S. J. F. Guimarães, Graph-based hierarchical video segmentation based on a simple dissimilarity measure, *Pattern Recognition Letters* 47 (2014) 85–92.

- [28] X. Zhu, C. Change Loy, S. Gong, Constructing robust affinity graphs for spectral clustering, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2014, pp. 1450–1457.
- [29] Z. Li, X.-M. Wu, S.-F. Chang, Segmentation using superpixels: A bipartite graph partitioning approach, in: Proceedings of the Conference on Computer Vision and Pattern Recognition, 2012, pp. 789–796.
- [30] X. Wang, Y. Tang, S. Masnou, L. Chen, A global/local affinity graph for image segmentation, Image Processing, IEEE Transactions on 24 (4) (2015) 1399–1411.
- [31] P. Isola, D. Zoran, D. Krishnan, E. H. Adelson, Crisp boundary detection using pointwise mutual information, in: European Conference on Computer Vision, 2014, pp. 799–814.
- [32] R. Xiao Feng, L. Bo, Discriminatively trained sparse code gradients for contour detection, in: Advances in Neural Information Processing Systems, 2012, pp. 584–592.
- [33] M. Donoser, D. Schmalstieg, Discrete-continuous gradient orientation estimation for faster image segmentation, in: Proceedings of the Conference on Computer Vision and Pattern Recognition, 2014, pp. 3158–3165.
- [34] S. Kim, C. D. Yoo, S. Nowozin, P. Kohli, Image segmentation using higher-order correlation clustering, Transactions on Pattern Analysis and Machine Intelligence 36 (9) (2014) 1761–1774.
- [35] P. Dollár, C. L. Zitnick, Fast edge detection using structured forests, Transactions on pattern analysis and machine intelligence 37 (8) (2015) 1558–1570.
- [36] J. Pont-Tuset, P. Arbelaez, J. T. Barron, F. Marques, J. Malik, Multiscale combinatorial grouping for image segmentation and object proposal generation, Transactions on pattern analysis and machine intelligence 39 (1) (2017) 128–140.
- [37] C. J. Taylor, Towards fast and accurate segmentation, in: Proceedings of the Conference on Computer Vision and Pattern Recognition, 2013, pp. 1916–1922.

- [38] W. Shen, B. Wang, Y. Jiang, Y. Wang, A. Yuille, Multi-stage multi-recursive-input fully convolutional networks for neuronal boundary detection, Proceedings of the International Conference on Computer Vision (ICCV).
- [39] K. Maninis, J. Pont-Tuset, P. Arbeláez, L. V. Gool, Convolutional oriented boundaries: From image segmentation to high-level tasks, IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI).
- [40] W. Shen, X. Wang, Y. Wang, X. Bai, Z. Zhang, Deepcontour: A deep convolutional feature learned by positive-sharing loss for contour detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 3982–3991.
- [41] Z. Zhang, F. Xing, X. Shi, L. Yang, Semicontour: A semi-supervised learning approach for contour detection, in: Proceedings of the Conference on Computer Vision and Pattern Recognition, 2016, pp. 251–259.
- [42] W. Shen, K. Zhao, Y. Jiang, Y. Wang, Z. Zhang, X. Bai, Object skeleton extraction in natural images by fusing scale-associated deep side outputs, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 222–230.
- [43] J. Long, E. Shelhamer, T. Darrell, Fully convolutional networks for semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 3431–3440.
- [44] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. L. Yuille, Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs, arXiv preprint arXiv:1606.00915.
- [45] H. Noh, S. Hong, B. Han, Learning deconvolution network for semantic segmentation, in: Proceedings of the IEEE International Conference on Computer Vision, 2015, pp. 1520–1528.
- [46] H. Zhou, J. Zheng, L. Wei, Texture aware image segmentation using graph cuts and active contours, Pattern Recognition 46 (6) (2013) 1719–1733.

- [47] Z. Li, G. Liu, D. Zhang, Y. Xu, Robust single-object image segmentation based on salient transition region, *Pattern Recognition* 52 (2016) 317 – 331.
- [48] L. Ladicky, C. Russell, P. Kohli, P. H. Torr, Graph cut based inference with co-occurrence statistics, in: *European Conference on Computer Vision*, 2010, pp. 239–253.
- [49] R. M. Fano, D. Hawkins, Transmission of information: A statistical theory of communications, *American Journal of Physics* 29 (11) (1961) 793–794.
- [50] M. Belkin, P. Niyogi, Laplacian eigenmaps and spectral techniques for embedding and clustering., in: *Advances in Neural Information Processing Systems*, 2001, pp. 585–591.
- [51] S. T. Roweis, L. K. Saul, Nonlinear dimensionality reduction by locally linear embedding, *Science* 290 (5500) (2000) 2323–2326.
- [52] J. B. Tenenbaum, V. De Silva, J. C. Langford, A global geometric framework for nonlinear dimensionality reduction, *Science* 290 (5500) (2000) 2319–2323.
- [53] E. Parzen, On estimation of a probability density function and mode, *The annals of mathematical statistics* (1962) 1065–1076.
- [54] M. Cheng, N. J. Mitra, X. Huang, P. H. Torr, S. Hu, Global contrast based salient region detection, *Transactions on Pattern Analysis and Machine Intelligence* 37 (3) (2015) 569–582.
- [55] J. Kim, D. Han, Y.-W. Tai, J. Kim, Salient region detection via high-dimensional color transform, in: *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2014, pp. 883–890.
- [56] J. Yang, Y.-H. Tsai, M.-H. Yang, Exemplar cut, in: *Proceedings of the International Conference on Computer Vision*, 2013, pp. 857–864.
- [57] C. Farabet, C. Couprie, L. Najman, Y. LeCun, Learning hierarchical features for scene labeling, *Transactions on Pattern Analysis and Machine Intelligence* 35 (8) (2013) 1915–1929.

- [58] F. Liu, G. Lin, C. Shen, Crf learning with cnn features for image segmentation, *Pattern Recognition* 48 (10) (2015) 2983–2992.
- [59] V. Kolmogorov, R. Zabini, What energy functions can be minimized via graph cuts?, *Transactions on Pattern Analysis and Machine Intelligence* 26 (2) (2004) 147–159.
- [60] D. Comaniciu, P. Meer, Mean shift: A robust approach toward feature space analysis, *Transactions on Pattern Analysis and Machine Intelligence* 24 (5) (2002) 603–619.
- [61] Z. Ren, G. Shakhnarovich, Image segmentation by cascaded region agglomeration, in: *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2013, pp. 2011–2018.
- [62] Y. Li, M. Paluri, J. M. Rehg, P. Dollár, Unsupervised learning of edges, in: *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1619–1627.