

# Probabilistic Diffusion for Interactive Image Segmentation

Tao Wang, Jian Yang, *Member, IEEE*, Zexuan Ji, *Member, IEEE*, and Quansen Sun

**Abstract**—This paper presents an interactive image segmentation approach in which we formulate segmentation as a probabilistic estimation problem based on the prior user intention. Instead of directly measuring the relationship between pixels and labels, we first estimate the distances between pixel pairs and label pairs using a probabilistic framework. Then, binary probabilities with label pairs are naturally converted to unary probabilities with labels. The higher-order relationship helps improve the robustness to user inputs. To improve segmentation accuracy, a likelihood learning framework is proposed to fuse the region and the boundary information of the image by imposing a smoothing constraint on the unary potentials. Furthermore, we establish an equivalence relationship between likelihood learning and likelihood diffusion and propose an iterative diffusion-based optimization strategy to maintain computational efficiency. Experiments on the Berkeley segmentation data-set and Microsoft GrabCut database demonstrate that the proposed method can obtain better performance than state-of-the-art methods.

**Index Terms**—Interactive image segmentation, paired distance measurement, likelihood learning, probabilistic estimation, unary potentials.

## I. INTRODUCTION

Image segmentation can be described as the partitioning of an image into several connected homogeneous regions based on similarity criteria using low-level visual features and extracting one or more objects that are of interest to the user from the background environment. Segmented semantic regions or contours associated with real-world entities or scenes are the basis for further advanced image processing. Therefore, image segmentation is a key step from image processing to image analysis, which is a fundamental problem in computer vision.

Many image segmentation methods have been proposed in the literature. Segmentation schemes can be classified into unsupervised, semi-supervised and fully supervised approaches. Unsupervised schemes can automatically segment images based on feature clustering. Due to the lack of prior knowledge

of each class, such approaches lack universality and thus are often used for specific tasks or preprocessing steps of segmentation, such as the generation of superpixels. Fully supervised schemes, such as convolutional networks [1-2], utilize a training set of images for semantic segmentation. The segmentation results are associated with the training samples of objects. However, for the same image, different users may not be interested in the same target, which causes these approaches to lack flexibility. Semi-supervised schemes, such as the graph cut approach [3], allow the user to provide simple interactions to represent label information during segmentation. Compared with the other two segmentation schemes, semi-supervised schemes can add the users' intentions to obtain results meeting their demands. Furthermore, the prior label information provided by the user helps improve segmentation performance.

This paper considers semi-supervised schemes (also called interactive approaches) for foreground-background segmentation. Given input image  $I$ , we aim to classify its pixels as one of two mutually exclusive classes,  $F$  and  $B$ , corresponding to foreground and background objects, respectively. During the last decades, many interactive segmentation approaches have been proposed, such as graph cut [3] and random walk [4]. In these approaches, unary and pairwise potentials that correspond to region and boundary information, respectively, are generally constructed for segmentation. A unary potential measures the similarity of a pixel to the labels  $F$  and  $B$ , while a pairwise potential quantifies the similarity between pairs of pixels. Unary and pairwise relationships can be represented via graphs and graph theory-based optimization algorithms can be used to produce segmentations [5]. Since the prior information of each class is provided by the user, unary potential is usually quantified as the distance between unseeded pixels and seeded pixels via some clustering algorithm, such as Gaussian mixture model (GMM). If enough seeds are given, GMM can accurately estimate the potential distribution of each label. However, due to the defects of pixel-level features, it is hard to capture the accurate label information when the user's interaction is limited (see Fig. 1(b)–(c)). In this case, the user has to work harder to obtain satisfactory results. Thus, effectively computing unary potential based on seeds is a key problem of the interactive segmentation method.

As an effective strategy to solve the above problems, methods based on perceptual grouping laws have been proposed [6]. Superpixels produced with unsupervised segmentation algorithms [7-8] are used to capture long-range grouping cues. These methods were inspired by constraints, whereby pixels constituting a particular superpixel should have the same label. They benefit from using more informative

This work was supported by the Natural Science Foundation of Jiangsu Province, China under Grant BK20180458, the National Science Foundation of China under Grant 61802188, the Postdoctoral Innovative Talent Support Program of China under Grant BX201700121, the China Postdoctoral Science Foundation under Grant 2017M621750, the National Science Foundation of China under Grants 61673220, U1713208 and 61472187, the 973 Program No.2014CB349303 and Program for Changjiang Scholars. (Corresponding author: Jian Yang and Quansen Sun.)

Tao Wang, Jian Yang, Zexuan Ji and Quansen Sun are with the School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing 210094, China (e-mail: [wangtaoatnjust@163.com](mailto:wangtaoatnjust@163.com); [csjyang@njust.edu.cn](mailto:csjyang@njust.edu.cn); [jizexuan@njust.edu.cn](mailto:jizexuan@njust.edu.cn); [sunquansen@njust.edu.cn](mailto:sunquansen@njust.edu.cn)).

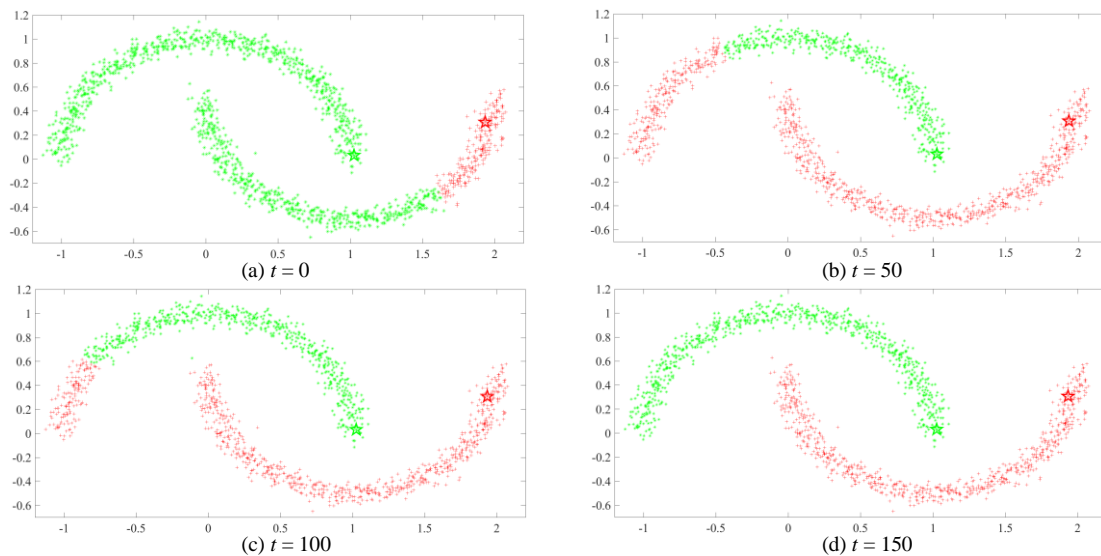


Fig. 2. Illustration of the effectiveness of probabilistic diffusion on toy data, where two seed points are defined and highlighted as stars, and all other elements are labeled according to their affinity to the seed points. (a)–(d) clustering results w.r.t. the number of iterations, where  $t$  represents the diffusion number. As shown in (a), the initial probability (without diffusion) based on the relationship between unseeded elements and seeded elements is not sufficient to capture the intrinsic structure of the data manifold when seed information is limited. Probabilistic diffusion can capture intrinsic global affinities, and thus significantly improve clustering performance.

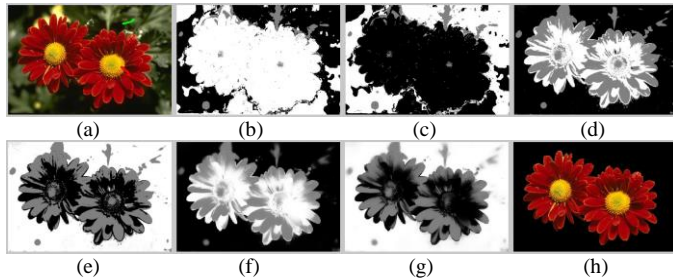


Fig. 1. Estimations of unary potentials from user inputs: (a) test image with limited seeds; (b)–(c) probability maps estimated by GMM with labels  $F$  and  $B$ , respectively; (d)–(e) probability maps captured by pixel-pair-based measurement; (f)–(g) results via pairwise likelihood learning; and (h) final segmentation result.

features extracted from the pixels within the regions. However, superpixels tend to not emphasize proximity sufficiently and thereby generate isolated regions in the segmentation results. To obtain reliable results, the relationships among pixels and superpixels are fused to enforce proximity and continuity. Both geometrical adjacency and long-range cues are critical for segmentation. As suggested in [9], superpixels are not always consistent with boundaries in natural images. To overcome the influence of inaccurate superpixels, many methods combining multiple superpixels of the same image have been proposed [9–11]. Multiple superpixels can be produced based on one unsupervised segmentation algorithm with different controlling parameters or different unsupervised segmentation algorithms. The propagation of superpixel-based grouping cues can help improve robustness to user inputs and obtain more accurate segmentation results. However, more superpixel variables need to be defined, and more relationships between pixels and superpixels need to be computed in these approaches, which increases algorithm complexity. Furthermore, several parameters are used to control the influence of pixel-level and superpixel-level relationships in these models, and the segmentation results are generally sensitive to these controlling parameters. To simplify the connections among multiple

superpixels, superpixels are divided into small, medium and large sets in [12], where small and large superpixels are used to encode local smoothness, and medium superpixels are used to propagate sparse long-range grouping cues through  $l_0$  sparsity. However, the generation of multiple superpixels also has considerable costs.

To extract accurate object details, interactive segmentation methods should satisfy the robustness to seeds while maintaining a low running time. As described above, pixel-level-based approaches are sensitive to user inputs and perceptual grouping approaches are limited by high algorithm complexity. To address these problems, in this work, we propose an interactive image segmentation method based on probabilistic diffusion. Fig. 2 shows a toy example where two seed points (shown as stars) are defined and all other elements are labeled according to their affinity to the seed points. As shown in Fig. 2(a), the initial probability (without diffusion) based on the relationship between unseeded elements and seeded elements is not sufficient to capture the intrinsic structure of the data manifold when seed information is limited. In comparison, after diffusing the probabilities through the manifold and capturing the intrinsic global manifold structure, we obtain significantly improved clustering results. The contributions of the proposed approach are concluded as follows:

**First**, a probabilistic framework is proposed to estimate the distances between pixel pairs and label pairs. Instead of the original relationship measurement between pixels and two labels,  $F$  and  $B$ , the higher-order measurement between pixel pairs and four label pairs,  $(F, F)$ ,  $(F, B)$ ,  $(B, F)$  and  $(B, B)$ , helps produce more accurate relationships between unseeded pixels and seeded pixels. **Second**, the binary probabilities with label pairs can be naturally converted to unary probabilities with labels (see Fig. 1(d)–(e)). **Third**, to further improve the segmentation accuracy, a likelihood learning framework is proposed to impose a smoothing constraint on unary potentials based on the pairwise similarities of pixels (see Fig. 1 (f)–(g)).

**Fourth**, an equivalence relationship between likelihood learning and likelihood diffusion is established, and an iterative diffusion-based optimization strategy is proposed to improve computational efficiency. Partial pilot data presented in our previous work proposed a pairwise likelihood learning method for interactive image segmentation [13].

The rest of this paper is organized as follows: In Section II, we will review related work to provide insight into interactive graph-theory-based approaches. In Section III, we will introduce the proposed method in detail. In Section IV, we will show a series of experimental results on different datasets. Section V will present the paper's conclusions.

## II. RELATED WORK

**Graph cut:** Boykov and Jolly [3] first proposed the interactive graph cut method to segment grayscale medical images. Lazy snapping [14] constructed a graph based on superpixels instead of pixels to improve efficiency. A coarse-to-fine user interface was also designed to provide instant visual feedback. GrabCut [15] extended the graph cut approach to color images by using GMMs to model the  $F$  and  $B$  regions. Incomplete trimaps were also provided to simplify user interaction through an iterative optimization process. Deep GrabCut [16] combined a convolutional encoder-decoder network trained end-to-end to overcome issues with the size of the interactive bounding box. In these works [17-18], color and texture information were efficiently combined to overcome difficulties handling images containing textures. For objects with specific shapes, shape priors were introduced into the graph cut framework [19] to restrict the segmentation results to a particular class of shapes, which helps improve the accuracy of objects lacking salient edges. ACP-cut [20] used semi-supervised kernel matrix learning to preserve the details around object boundaries. Moreover, seed information was propagated to achieve discriminative structure learning and reduce computational complexity. In addition to the graph cut for the optimization of binary cost functions, in this work [21], segmentation is formulated as an inference problem based on unary and pairwise assignment probabilities via a probabilistic graph matching scheme.

**Random walk:** Grady *et al.* [22] first proposed interactive random walk for medical image segmentation and extended it for general image segmentation in [4]. To overcome weak boundary and texture problems, random walk with restart [23] constructed a generative segmentation model by using the steady-state probability to reduce dependence on seeds. Lazy random walk [24] with self-loops considered the global relationships between all pixels and seeds to solve the superpixel segmentation problem in the weak boundary and texture regions. Sub-Markov random walk [25] introduced the label prior by adding auxiliary nodes to further improve segmentation accuracy—especially in thin and elongated regions. Constrained random walk [26] advocated the use of multiple intuitive user inputs to better reflect a user's intentions. In contrast to the methods mentioned above that formally minimize the “distance” between pairwise pixels, Laplacian coordinates [27] minimize average distances while better controlling anisotropic propagation of labels, which ensures a

better fit on image boundaries. Normalized random walk [28] incorporated a degree-aware term into the original model to account for the node centrality of every neighboring node and weigh the contribution of every neighbor to the underlying diffusion process.

**Perceptual grouping approaches:** To improve the robustness to user inputs, many perceptual grouping methods have been proposed that use superpixels to capture long-range grouping cues [9-11]. The robust  $P^n$  model [10] constructed higher-order parametric potentials based on multiple superpixels with conventional unary and pairwise constraints by using higher-order condition random fields in a principled manner. Instead of estimating the number of pixels that do not belong to the dominant label, in this work [29], the sum of weights for pixels in the superpixel not taking the dominant label is measured, which helps produce finer higher-order potentials. The nonparametric higher-order model [9] used the pairwise relationship between pixels and their corresponding superpixels using a multi-layer graph. A higher-order cost function of pixel likelihoods is designed to enforce label consistency in superpixels. To further improve segmentation performance, this work [11] introduced prior label information to the nonparametric higher-order model, and the multi-layer constraints among pixels, superpixels and labels are fused for segmentation. As suggested in [9], these superpixel-based methods are less sensitive to user inputs and produce high-quality segmentation results.

The local relationship used in the graph cut and random walk approaches makes them sensitive to seed quantity and position; thus, it is hard for them to keep global coherence with limited user interactions. Perceptual grouping approaches extend the local relationship of pixels to the long-range regional connectivity of superpixels. However, the generation of multiple superpixels and the computation of higher-order relationships lead to high algorithm complexity in the corresponding algorithms. Compared with these works, the advantages of the proposed method are summarized as follows. First, a higher-order measurement between pixel pairs and label pairs is used to obtain a more accurate unary potential under limited seeds. Second, a likelihood learning method is proposed to diffuse local probabilities, which can help capture the intrinsic global manifold structure, making the diffusion process simpler and more efficient than the perceptual grouping between pixels and superpixels.

## III. IMAGE SEGMENTATION BY PAIRWISE LIKELIHOOD LEARNING

The segmentation model is formulated as a probabilistic estimation problem. Fig. 3 illustrates the framework of the proposed algorithm. The initial probability  $p^0$  is first estimated based on the seed information. A paired distance estimation method is proposed to measure the relationships between pixel pairs and label pairs (shown in III-A), which can help improve the accuracy of unary potential with limited seeds. After the binary probability transformation,  $p^0$  is updated as the unary assignment probability  $p^1$  (shown in III-B). A likelihood learning method is proposed to fuse the region and boundary information of the image to extend  $p^1$  to the final

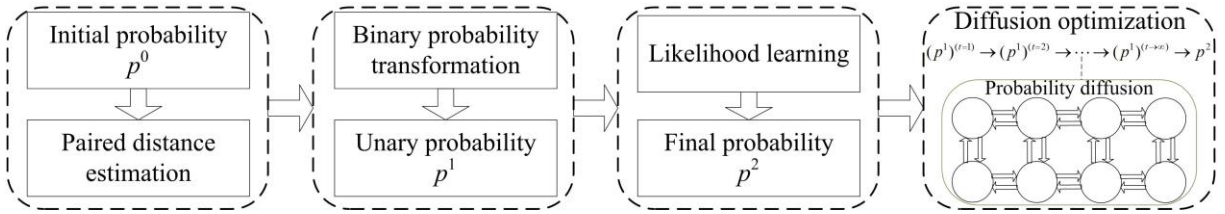


Fig. 3. Overview of the segmentation framework, where  $p^0$  represents the initial probability estimated from the seeds,  $p^1$  represents the updated probability after the paired distance estimation,  $p^2$  represents the final probability by the likelihood learning and diffusion optimization and  $t$  represents the number of iterations.

probability  $p^2$  (shown in III-C). After equivalence analysis with likelihood diffusion (shown in III-D), the likelihood learning process can capture the intrinsic global manifold structure by diffusing the local probabilities. A diffusion-based optimization technique is proposed to improve optimization efficiency.

#### A. Paired Distance Estimation

Input image  $I$  can be represented by graph  $G = (X, W)$ , where  $X = \{x_i\}_{i=1}^N$  is a collection of pixels,  $W = [W_{ij}]_{N \times N}$  is the relationships of pairwise pixels and  $N$  is the number of pixels. The task of segmentation is to classify each image element  $x_i$  as  $f_{x_i} \in \{F, B\}$ . The core of the proposed approach is the estimation of marginal assignment probabilities based on a paired relationship metric and pairwise likelihood learning.

The similarity  $W_{ij}$  between pairwise pixels  $(x_i, x_j)$  is defined as a typical Gaussian function:

$$W_{ij} = \begin{cases} \exp(-\beta^w \|c_i - c_j\|_2^2) & \text{if } (x_i, x_j) \in \mathbb{N} \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

where  $c_i$  and  $c_j$  denote the intensity feature at pixels  $x_i$  and  $x_j$ , respectively,  $\mathbb{N}$  is the set of all neighboring pixel pairs in the image and  $\beta^w > 0$  is a constant that controls the strength of the weight, automatically selected as:

$$\beta^w = \frac{1}{2E[\|c_i - c_j\|_2^2]} \quad (2)$$

where  $E[\cdot]$  represents the expectation over all pixel pairs in  $\mathbb{N}$ . It can be noticed that if two neighboring pixels have similar features, their weight is large, and vice versa.

User inputs represent the prior label information, and the conventional methods measure the distances between pixels and seeds to estimate the unary assignment probabilities (shown in Fig. 4(a)). These methods are generally sensitive to user inputs, and it is hard to obtain accurate results when the number of seeds is limited. To improve robustness to user inputs, this paper considers paired probability estimation by measuring the similarities between pixel pairs and label pairs:  $(F, F)$ ,  $(F, B)$ ,  $(B, F)$  and  $(B, B)$  (shown in Fig. 4(b)).

For any pixel pair  $(x_i, x_j) \in \mathbb{N}$ , the distances  $d_{ij}^{FF}$ ,  $d_{ij}^{FB}$ ,  $d_{ij}^{BF}$  and  $d_{ij}^{BB}$  with label pairs  $(F, F)$ ,  $(F, B)$ ,  $(B, F)$  and  $(B, B)$  are defined as:

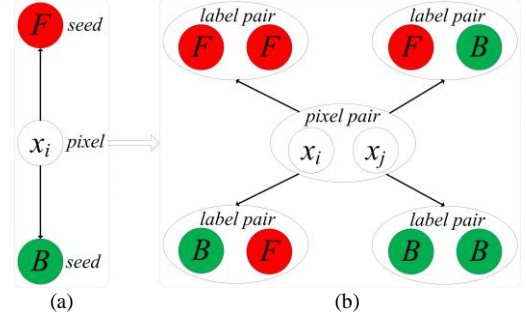


Fig. 4. Sketch map of prior information estimation from user inputs: (a) the distance measurement between pixels and seeds by conventional methods, (b) the paired relationship measurement between pixel pairs and label pairs by the proposed method.

$$d_{ij}^{FF} = p(x_i \in F, x_j \in F) = \frac{1}{2} [p(x_i \in F | x_j \in F) p^0(x_j \in F) + p(x_j \in F | x_i \in F) p^0(x_i \in F)] \quad (3)$$

$$d_{ij}^{FB} = p(x_i \in F, x_j \in B) = \frac{1}{2} [p(x_i \in F | x_j \in B) p^0(x_j \in B) + p(x_j \in B | x_i \in F) p^0(x_i \in F)] \quad (4)$$

$$d_{ij}^{BF} = p(x_i \in B, x_j \in F) = \frac{1}{2} [p(x_i \in B | x_j \in F) p^0(x_j \in F) + p(x_j \in F | x_i \in B) p^0(x_i \in B)] \quad (5)$$

$$d_{ij}^{BB} = p(x_i \in B, x_j \in B) = \frac{1}{2} [p(x_i \in B | x_j \in B) p^0(x_j \in B) + p(x_j \in B | x_i \in B) p^0(x_i \in B)] \quad (6)$$

The initial probabilities of pixels with labels can be estimated based on prior user inputs. For simplicity, in this paper, the k-means algorithm is used to cluster both the foreground and background seeds. Cluster centers  $\{c_k^L\}_{k=1}^K$  are then produced, where  $L = F/B$  and  $K$  is the number of clusters. For each pixel  $x_i \in X$ , the initial probability  $p^0(x_i \in L)$  is defined as:

$$p^0(x_i \in L) = \max_{k \in \{1, \dots, K\}} \exp(-\|c_i - c_k^L\|_2) \quad (7)$$

The value of  $p^0(x_i \in L)$  is normalized under the constraint  $p^0(x_i \in F) + p^0(x_i \in B) = 1$ .

The conditional probabilities can be associated with the relationships between the two pixels. For example, if the similarity weight between  $x_i$  and  $x_j$  is large, they are likely to belong to the same label. Therefore, under the condition  $f_{x_j} \in \{F, B\}$ , the probability that pixel  $x_i$  belongs to the same label ( $f_{x_i} = f_{x_j}$ ) is high, and the probability that pixel  $x_i$  belongs to the different label ( $f_{x_i} \neq f_{x_j}$ ) is low. Otherwise, if the



similarity weight between  $x_i$  and  $x_j$  is small, they are likely to belong to different labels. Therefore, the probability that pixel  $x_i$  belongs to the different label ( $f_{x_i} \neq f_{x_j}$ ) is high, and the probability that pixel  $x_i$  belongs to the same label ( $f_{x_i} = f_{x_j}$ ) is low. Based on the above observation, in this paper, we simply define the conditional probabilities as follows:

$$p(f_{x_i} | f_{x_j}) = \begin{cases} \hat{W}_{ij} & \text{if } f_{x_i} = f_{x_j} \\ 1 - \hat{W}_{ij} & \text{if } f_{x_i} \neq f_{x_j} \end{cases} \quad (8)$$

where  $f_{x_i} \in \{F, B\}$  and  $f_{x_j} \in \{F, B\}$  represent the labels of pixels  $x_i$  and  $x_j$ , respectively.  $\hat{W}_{ij}$  is an extended relationship weight of pixels  $x_i$  and  $x_j$ , which both considers the similarity in the intensity feature and the prior probability. In this way, it is better to measure the relationship between different components in the same object. The value of  $\hat{W}_{ij}$  is defined as:

$$\hat{W}_{ij} = \frac{1}{2}(W_{ij} + DF_{ij}) \quad (9)$$

where  $W_{ij}$  represents the similarity of pixel pair  $(x_i, x_j)$  in the intensity feature and  $DF_{ij}$  represents the similarity in probabilities with the foreground label, which is defined as:

$$DF_{ij} = \exp(-\beta^{DF} \|p^0(x_i \in F) - p^0(x_j \in F)\|_2^2) \quad (10)$$

where  $\beta^{DF} > 0$  is a constant that controls the strength of  $DF_{ij}$ , automatically selected as:

$$\beta^{DF} = \frac{1}{2E[\|p^0(x_i \in F) - p^0(x_j \in F)\|_2^2]} \quad (11)$$

### B. Prior Assignment Probabilities

Because  $W_{ij} = W_{ji}$  and  $DF_{ij} = DF_{ji}$ , the values of  $d_{ij}^{FF}$ ,  $d_{ij}^{FB}$ ,  $d_{ij}^{BF}$  and  $d_{ij}^{BB}$  can be simplified as:

$$d_{ij}^{FF} = \frac{1}{2}\hat{W}_{ij}[p^0(x_i \in F) + p^0(x_j \in F)] \quad (12)$$

$$d_{ij}^{FB} = \frac{1}{2}(1 - \hat{W}_{ij})[p^0(x_i \in F) + p^0(x_j \in B)] \quad (13)$$

$$d_{ij}^{BF} = \frac{1}{2}(1 - \hat{W}_{ij})[p^0(x_i \in B) + p^0(x_j \in F)] \quad (14)$$

$$d_{ij}^{BB} = \frac{1}{2}\hat{W}_{ij}[p^0(x_i \in B) + p^0(x_j \in B)] \quad (15)$$

Since  $p^0(x_i \in F) + p^0(x_i \in B) = 1$  and  $p^0(x_j \in F) + p^0(x_j \in B) = 1$ , it can be seen that:

$$\begin{aligned} d_{ij}^{FF} + d_{ij}^{FB} + d_{ij}^{BF} + d_{ij}^{BB} \\ = \frac{1}{2}[p^0(x_i \in F) + p^0(x_i \in B) + p^0(x_j \in F) + p^0(x_j \in B)] = 1 \end{aligned} \quad (16)$$

The binary assignment probabilities of pixels  $x_i$  and  $x_j$  can be converted to unary assignment probabilities. The initial probabilities  $p^0(x_i \in L)$  and  $p^0(x_j \in L)$  are updated to  $p^1(x_i \in L)$  and  $p^1(x_j \in L)$ , respectively. For any pixel  $x_i \in X$ :

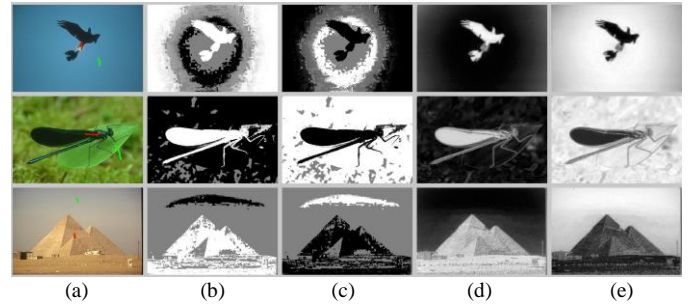


Fig. 5. Comparison of unary potentials: (a) the test image with seeds; (b)–(c) the probability maps of  $F$  and  $B$ , respectively, based on the pixel-level measurement by GMM; (d)–(e) the probability maps of  $F$  and  $B$ , respectively, based on the proposed pixel-pair-level measurement.

$$\begin{aligned} p^1(x_i \in F) &= \sum_{x_j \in \mathbb{V}_i} d_{ij}^{FF} + d_{ij}^{FB} \\ &= \sum_{x_j \in \mathbb{V}_i} \frac{1}{2} [p^0(x_i \in F) + \hat{W}_{ij} p^0(x_j \in F) + (1 - \hat{W}_{ij}) p^0(x_j \in B)] \end{aligned} \quad (17)$$

$$\begin{aligned} p^1(x_i \in B) &= \sum_{x_j \in \mathbb{V}_i} d_{ij}^{BB} + d_{ij}^{BF} \\ &= \sum_{x_j \in \mathbb{V}_i} \frac{1}{2} [p^0(x_i \in B) + \hat{W}_{ij} p^0(x_j \in B) + (1 - \hat{W}_{ij}) p^0(x_j \in F)] \end{aligned} \quad (18)$$

where  $\mathbb{V}_i$  represents the neighborhood of  $x_i$ . After the binary assignment probabilities of all pixel pairs in set  $\mathbb{S}$  are transformed into unary probabilities, the value of  $p^1(x_i \in L)$  for each pixel  $x_i \in X$  is normalized under the constraint  $p^1(x_i \in F) + p^1(x_i \in B) = 1$ .

Examples of prior probability estimation are shown in Fig. 5, where (a) shows the test image with seeds, (b)–(c) show the results obtained by GMM based on the pixel-level distance and (d)–(e) show the results produced by the proposed method based on the pixel-pair-level relationship measurement. It can be seen that the pixel-level distances between pixels and seeds are not enough to discriminate the label information when the user inputs are limited. The pixel-pair-level relationships between pixel pairs and labels pairs can be regarded as the higher-order information of the pixel-level distances and can produce more accurate estimation of the prior probability under the same seed information. Furthermore, instead of the two labels  $F$  and  $B$ , the distances to the four label pairs,  $(F, F)$ ,  $(F, B)$ ,  $(B, F)$  and  $(B, B)$ , make the relationships more discriminating.

The relationships between doublets can be naturally extended to higher-order relationships. More interactions can be considered in higher-order measurement. For example, eight label groups are involved when measuring the distances between triplets; however, the number of label groups increases exponentially with an increase in order, which leads to higher algorithm complexity. To maintain computational efficiency, in this paper, only pixel-pair-level relationships are considered.

### C. Pairwise Likelihood Learning

The above probability estimation only considers the region information regardless of the boundary information of the image. To improve segmentation accuracy, we impose a

smoothing constraint on the unary prior potentials by a pairwise likelihood learning strategy.

Let  $\bar{P}^F = [p^1(x_i \in F)]_{N \times 1}$  and  $\bar{P}^B = [p^1(x_i \in B)]_{N \times 1}$  denote the probability vectors (abbreviate  $p^1(x_i \in F)$  as  $\bar{p}_i^F$  and  $p^1(x_i \in B)$  as  $\bar{p}_i^B$ ). Let  $\hat{P}^F = [p^2(x_i \in F)]_{N \times 1}$  and  $\hat{P}^B = [p^2(x_i \in B)]_{N \times 1}$  denote the novel probability vectors after the likelihood learning (abbreviate  $p^2(x_i \in F)$  as  $\hat{p}_i^F$  and  $p^2(x_i \in B)$  as  $\hat{p}_i^B$ ), where  $p^2(x_i \in F)$  and  $p^2(x_i \in B)$  represent the novel probabilities of pixel  $x_i$  belonging to the labels  $F$  and  $B$ , respectively.

The likelihood learning process of  $\hat{P}^F$  and  $\hat{P}^B$  is defined as minimizing the following two cost functions:

$$E(\hat{P}^F) = E_{\text{boundary}}(\hat{P}^F) + E_{\text{region}}(\hat{P}^F) \\ = \sum_{(x_i, x_j) \in N} W_{ij} \cdot (\hat{p}_i^F - \hat{p}_j^F)^2 + \lambda \sum_{x_i \in X} d_i \cdot \left[ \bar{p}_i^F \cdot (\hat{p}_i^F - \frac{1}{d_i})^2 + \bar{p}_i^B \cdot (\hat{p}_i^F)^2 \right] \quad (19)$$

$$E(\hat{P}^B) = E_{\text{boundary}}(\hat{P}^B) + E_{\text{region}}(\hat{P}^B) \\ = \sum_{(x_i, x_j) \in N} W_{ij} \cdot (\hat{p}_i^B - \hat{p}_j^B)^2 + \lambda \sum_{x_i \in X} d_i \cdot \left[ \bar{p}_i^B \cdot (\hat{p}_i^B - \frac{1}{d_i})^2 + \bar{p}_i^F \cdot (\hat{p}_i^B)^2 \right] \quad (20)$$

where  $E_{\text{boundary}}$  represents the boundary energy term,  $E_{\text{region}}$  represents the region energy term,  $d_i = \sum_{j=1}^N W_{ij}$  and the parameter  $\lambda = \alpha / (1 - \alpha)$  ( $0 < \alpha < 1$ ) is used to balance these two energy terms.  $E_{\text{boundary}}$  constrains neighboring pixels with high similarities to have similar likelihood probabilities.  $E_{\text{region}}$  constrains the likelihood estimation to maintain consistency with the prior probabilities. Pixel  $x_i$  should be assigned high  $\hat{p}_i^F$  and low  $\hat{p}_i^B$  if its prior probability  $\bar{p}_i^F$  is high, and vice versa. Both region and boundary information are considered in the likelihood learning process, which maintains regional connectivity and piecewise smooth in the segmentation results.

Reformulate the cost functions in Eqs. (19, 20) in matrix form:

$$E(\hat{P}^F) = (\hat{P}^F)^T (D - W) \hat{P}^F + \lambda D \left[ \left( \hat{P}^F - \frac{O}{D} \right)^T \Omega^F \left( \hat{P}^F - \frac{O}{D} \right) + (\hat{P}^F)^T \Omega^B \hat{P}^F \right] \quad (21)$$

$$E(\hat{P}^B) = (\hat{P}^B)^T (D - W) \hat{P}^B + \lambda D \left[ \left( \hat{P}^B - \frac{O}{D} \right)^T \Omega^B \left( \hat{P}^B - \frac{O}{D} \right) + (\hat{P}^B)^T \Omega^F \hat{P}^B \right] \quad (22)$$

where  $D = \text{diag}([d_1, \dots, d_N])$ ,  $O = [1]_{N \times 1}$ ,  $\Omega^F = \text{diag}(\bar{P}^F)$  and  $\Omega^B = \text{diag}(\bar{P}^B)$ . Differentiating  $E(\hat{P}^F)$  and  $E(\hat{P}^B)$  with respect to  $\hat{P}^F$  and  $\hat{P}^B$ , respectively, and setting to zero, we have:

$$\frac{\partial E(\hat{P}^F)}{\partial \hat{P}^F} = (D - W) \hat{P}^F + \lambda D (\Omega^F + \Omega^B) \hat{P}^F - \lambda \Omega^F O = 0 \quad (23)$$

$$\frac{\partial E(\hat{P}^B)}{\partial \hat{P}^B} = (D - W) \hat{P}^B + \lambda D (\Omega^F + \Omega^B) \hat{P}^B - \lambda \Omega^B O = 0 \quad (24)$$

Since  $(\Omega^F + \Omega^B)$  is an identity matrix,  $\Omega^F O = \bar{P}^F$ ,  $\Omega^B O = \bar{P}^B$  and  $\lambda = \alpha / (1 - \alpha)$ :

$$\hat{P}^F = T \bar{P}^F \quad (25)$$

$$\hat{P}^B = T \bar{P}^B \quad (26)$$

where  $T = \alpha(D - (1 - \alpha)W)^{-1}$ . It can be seen that the likelihood probabilities  $\hat{P}^F$  and  $\hat{P}^B$  can be obtained by multiplying matrix  $T$  and prior probability vectors  $\bar{P}^F$  and  $\bar{P}^B$ .

#### D. Iterative Diffusion-based Optimization

Iterative diffusion-based optimization involves solving the inverse of a matrix of size  $N \times N$  to compute matrix  $T$ ; time complexity is  $O(N^3)$ . Although there are some approximate strategies to improve computational efficiency, such as the MATLAB division operator '\', the algorithm burden is heavy when the image size is large.

To overcome the above limitation, likelihood learning is converted to likelihood diffusion. Instead of computing matrix  $T$ , we diffuse the probabilities iteratively until convergence; time complexity is  $O(N^2)$ . Following [30], the diffusion strategy is described as:

$$(\hat{P}^L)^{(t)} = \alpha Q (\hat{P}^L)^{(t-1)} + (1 - \alpha) \bar{P}^L \quad (27)$$

where  $Q$  is the row-normalized matrix of  $W$ :  $Q = D^{-1} \times W$  and  $t$  is the iterative step.

**Convergence Analysis:** Referring to [30, 31], the closed form of the diffusion matrix  $\hat{P}^L$  at step  $t$  can be written as:

$$(\hat{P}^L)^{(t)} = (1 - \alpha)^{t-1} Q^{t-1} \bar{P}^L + \alpha \sum_{i=0}^{t-2} (1 - \alpha)^i Q^i \bar{P}^L \quad (28)$$

Because  $0 < \alpha < 1$ , we can derive:

$$\lim_{t \rightarrow \infty} (1 - \alpha)^{t-1} Q^{t-1} \bar{P}^L = 0$$

and

$$\lim_{t \rightarrow \infty} \sum_{i=0}^{t-2} (1 - \alpha)^i Q^i = (I - (1 - \alpha)Q)^{-1} \quad (29)$$

where  $I$  is the identity matrix.

Hence, after self-normalization, the diffusion matrix converges to:

$$\lim_{t \rightarrow \infty} (\hat{P}^L)^{(t)} = \alpha(D - (1 - \alpha)W)^{-1} \bar{P}^L \quad (30)$$

By comparing the Eqs. (25, 26) and the Eq. (30), it can be concluded that the likelihood learning is equivalent to the likelihood diffusion. Therefore, we can design the algorithm as:

1. Initialize parameters  $K$  and  $\alpha$
2. Estimate the initial probabilities  $p^0$  with Eq. (7)
3. Compute values of  $W$ ,  $DF$  and  $\hat{W}$  with Eqs. (1, 10, 9)
4. Estimate the unary probabilities  $\bar{P}^L$  with Eqs. (17-18)
5. Set  $t = 1$ ,  $(\hat{P}^L)^{(t)} = \bar{P}^L$  and  $f^t = [0]_{N \times 1}$
6. Update the value of  $(\hat{P}^L)^{(t+1)}$  with Eq. (27)
7. Compute labels:  $f^{t+1} = [f_i^{t+1}]_{N \times 1}$   
where  $f_i^{t+1} = 1$  if  $(\hat{p}_i^F)^{(t+1)} > (\hat{p}_i^B)^{(t+1)}$ ; otherwise  $f_i^{t+1} = 0$
8. Check the termination condition:  
If  $f^{t+1}$  equals  $f^t$ , stop iteration; otherwise  $t = t + 1$ , go to 6

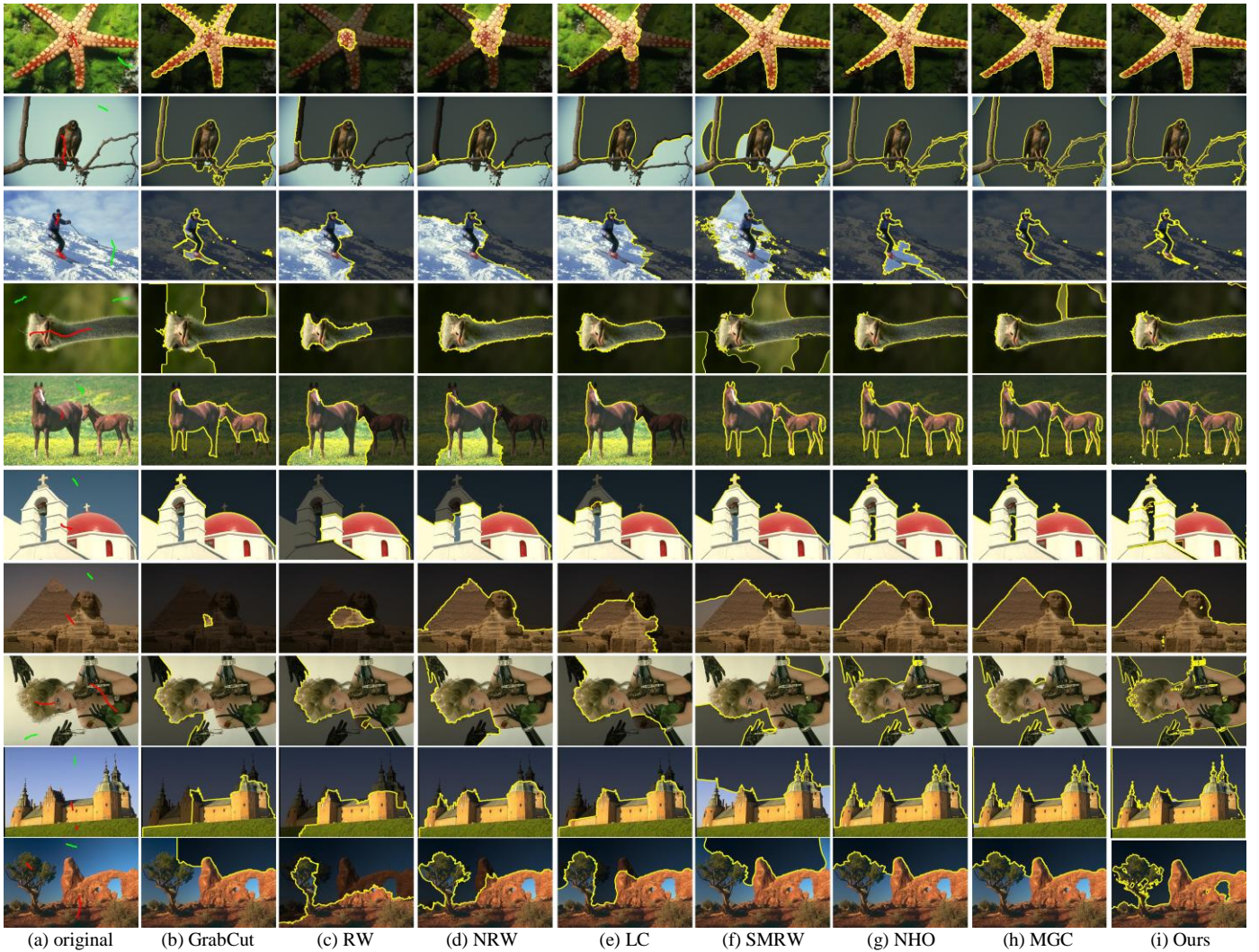


Fig. 6. Comparison to state-of-the-art approaches with a few scribbles. (a) Test image from the Berkeley segmentation dataset with scribbles (red: foreground, green: background); (b)–(i) segmentation results of GrabCut [15], RW [4], NRW [28], LC [27], SMRW [25], NHO [9], MGC [11] and the proposed method.

#### IV. EXPERIMENTAL RESULTS

The proposed method was experimentally verified by comparing it with state-of-the-art approaches—GrabCut [15], random walk (RW) [4], normalized random walk (NRW) [28], Laplacian coordinates (LC) [27], sub-Markov random walk (SMRW) [25], nonparametric higher-order method (NHO) [9], multi-layer graph constraints method (MGC) [11] and DeepGrabCut [16]—on the Berkeley segmentation dataset<sup>1</sup> and Microsoft GrabCut database<sup>2</sup>, all of which have ground-truth annotations. There are two parameters involved in the proposed scheme, and their values are set as follows: the number of clusters  $K$  is set to 4 and the region and boundary controlling parameter  $\alpha$  is set to 0.2. The 4-neighborhood relationship is used in the proposed method. The implementation codes of the compared algorithms are offered by the original authors, and the suggested parameters in their papers are used for the comparison experiments.

<sup>1</sup> <http://www.eecs.berkeley.edu/Research/Projects/CS/vision/grouping/segbench/S>.

<sup>2</sup> <http://research.microsoft.com/enus/um/cambridge/projects/visionimagevideoediting/segmentation/grabcut.htm>

##### A. Qualitative Benchmark Results

Fig. 6 shows the comparison with state-of-the-art interactive segmentation methods. Fig. 6(a) shows the test images from the Berkeley segmentation dataset with a few scribbles (red: foreground, green: background). Fig. 6(b)–(i) shows the segmentation results of GrabCut [15], RW [4], NRW [28], LC [27], SMRW [25], NHO [9], MGC [11] and the proposed method. It can be seen that GrabCut and RW cannot obtain satisfactory results when the number of seeds is limited. The RW extension methods, i.e., NRW, LC and SMRW, can obtain better results than RW; however, they are also sensitive to the seeds. The pixel-level information learnt from limited seeds is not enough to discriminate the foreground and background labels. Compared with these pixel-level-based approaches, the perceptual grouping methods NHO and MGC, obtain more robust results by using multiple superpixels to propagate long-range grouping laws. However, the object details around the boundaries cannot be well preserved, especially for many thin and slender regions. It can be clearly seen that the proposed method produces the best results with accurate object details. Paired distance estimation can help obtain accurate label prior information with limited seeds, and pairwise likelihood learning can help produce smooth object boundaries.



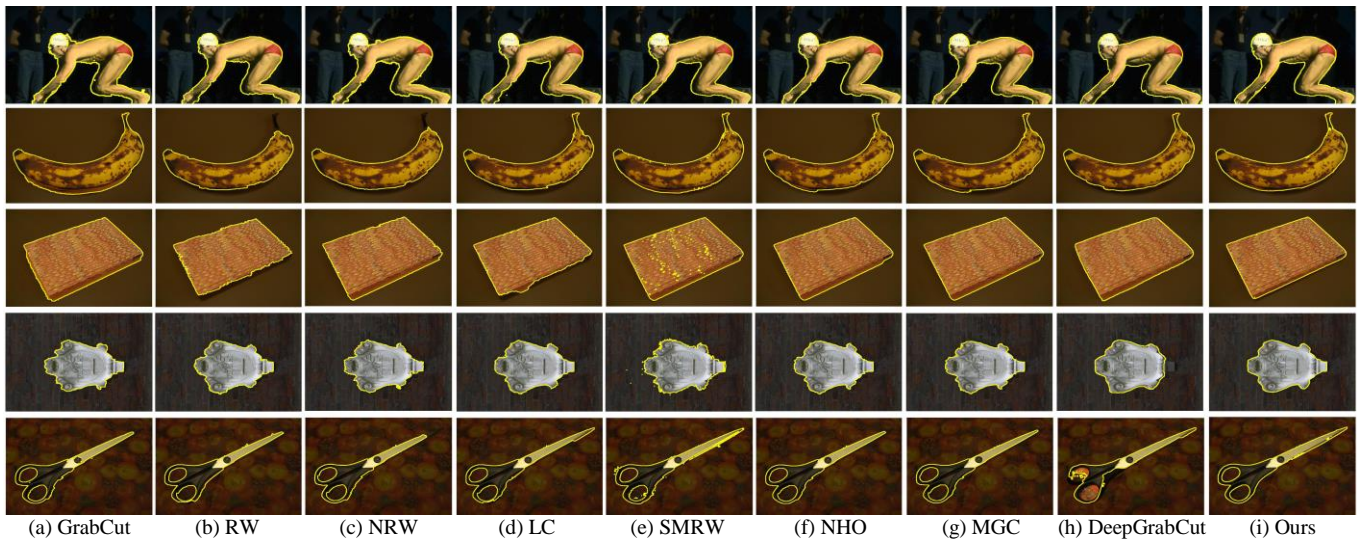


Fig. 7. Example segmentations on the Microsoft GrabCut database. (a)–(i) Segmentation results of GrabCut [15], RW [4], NRW [28], LC [27], SMRW [25], NHO [9], MGC [11], DeepGrabCut [16] and the proposed method, where DeepGrabCut is initialized based on the tight bounding boxes provided by [32] and other methods are initialized based on the public trimaps in this database.

TABLE I

Mean  $\pm$  standard deviation (Std) and the average rank (Ar) of PRI and VoI values for the compared methods on the Berkeley segmentation dataset.

Method	PRI		VoI	
	Mean $\pm$ Std	Ar	Mean $\pm$ Std	Ar
GrabCut [15]	0.56 $\pm$ 0.12	2.4	2.03 $\pm$ 0.45	5.6
RW [4]	0.60 $\pm$ 0.11	3.2	2.07 $\pm$ 0.33	6.3
NRW [28]	0.67 $\pm$ 0.07	5.3	1.91 $\pm$ 0.39	4.3
LC [27]	0.68 $\pm$ 0.07	5.3	1.94 $\pm$ 0.40	4.5
SMRW [25]	0.59 $\pm$ 0.09	3.1	2.12 $\pm$ 0.50	6.0
NHO [9]	0.65 $\pm$ 0.09	4.8	1.82 $\pm$ 0.49	4.2
MGC [11]	0.66 $\pm$ 0.12	5.1	1.72 $\pm$ 0.63	2.9
Ours	<b>0.71 <math>\pm</math> 0.08</b>	<b>6.8</b>	<b>1.61 <math>\pm</math> 0.42</b>	<b>2.2</b>

Fig. 7 illustrates the example segmentations on the Microsoft GrabCut database. Fig. 7(a)–(i) show the segmentation results of GrabCut [15], RW [4], NRW [28], LC [27], SMRW [25], NHO [9], MGC [11], DeepGrabCut [16] and the proposed method, where the results of DeepGrabCut are obtained based on the tight bounding boxes provided by [32] and other methods are initialized based on the public trimaps in this database. Affected by low contrast, it can be seen that the conventional methods cannot obtain accurate object boundaries. Furthermore, they are very sensitive to thin and slender objects and cannot obtain complete contours. Comparatively, the proposed method can achieve high-quality segmentation results. For the second and third test images, the proposed method can obtain accurate boundaries of the banana and the book regardless of the influence of low contrast. For the last test image, the proposed method can detect complete contours of the scissors. These qualitative comparison results demonstrate the superior performance of the proposed method.

### B. Quantitative Benchmark Results

Probabilistic rand index (PRI) [33] and variation of information (VoI) [34] are used to quantitatively evaluate the segmentation performance in the Berkeley segmentation dataset. PRI measures the agreement between the segmented result and the manually generated ground truth. PRI ranges from 0 to 1, with a higher value representing a more accurate result. VoI measures the information content in each segmentation and how much information one segmentation

gives about the other. VoI ranges in  $[0, \infty)$ , with a smaller value representing a more accurate result. TABLE I lists the mean  $\pm$  standard deviation and the average rank from the Friedman statistical test [35] (with a significance level of 0.05) of PRI and VoI for GrabCut [15], RW [4], NRW [28], LC [27], SMRW [25], NHO [9], MGC [11] and the proposed method on the Berkeley segmentation dataset. It can be observed that the proposed method outperforms the other methods with the largest PRI value and the smallest VoI value. The Friedman test determines the chi-square ( $\chi^2$ ) value as 25.13 (24.46) and the p-value as 7.1e-04 (9.4e-04) for PRI (VoI). From the  $\chi^2$  distribution table, we find the critical value for  $(8-1)=7$  degrees of freedom with 0.05 significance level is 14.07. Since the  $\chi^2$  value is larger than the critical value,  $H_0$  is rejected and  $H_1$  is accepted, which substantiates the significant difference in behavior among the compared methods.

We then demonstrated the quality of the proposed method on the Microsoft GrabCut database. Error rate, which is defined as the ratio of the number of wrongly labeled pixels to the total number of unlabeled pixels, is used to evaluate segmentation accuracy. Fig. 8 shows the error rate curves of each test image in the Microsoft GrabCut database by applying GrabCut [15], RW [4], NRW [28], LC [27], SMRW [25], NHO [9], MGC [11], DeepGrabCut [16] and the proposed method. The images are sorted in ascending order based on the values of the proposed method. It can be clearly observed that the proposed method obtains the best performance in most cases. TABLE II summarizes the average error rates obtained by various methods. Compared with the graph cut and random walk methods [3, 4, 14, 15, 17, 23-28], the proposed method has a significant improvement in the segmentation error rate. Compared with the perceptual grouping methods [9-11], the proposed method also shows a great improvement of the error rate. Compared with the methods in [5, 16, 21, 36], it can be observed that the proposed method outperforms the deep-learning-based methods and the latest interactive methods on the Microsoft GrabCut database.



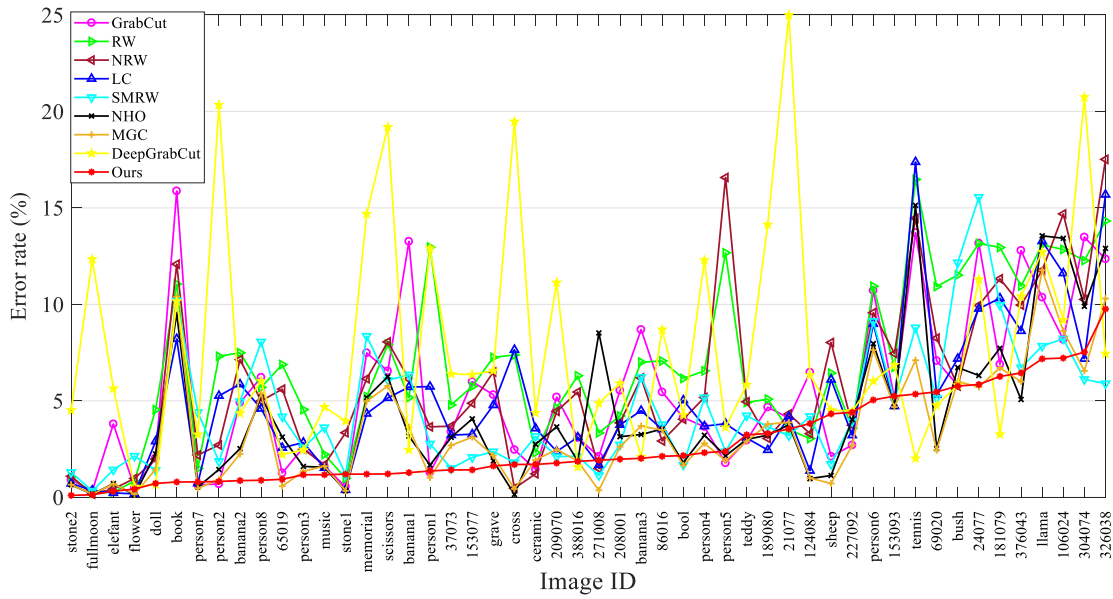


Fig. 8. Error rates of all images in the Microsoft GrabCut database by applying GrabCut [15], RW [4], NRW [28], LC [27], SMRW [25], NHO [9], MGC [11], DeepGrabCut [16] and the proposed method.

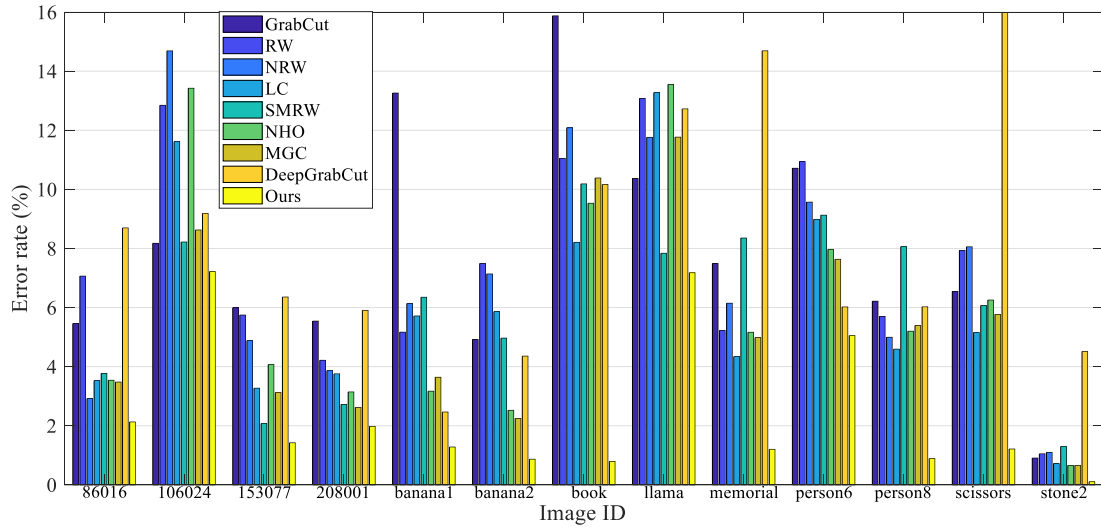


Fig. 9. Illustration of significantly improved results of the proposed method compared with GrabCut [15], RW [4], NRW [28], LC [27], SMRW [25], NHO [9], MGC [11] and DeepGrabCut [16] on the Microsoft GrabCut database.

TABLE II

Average error rates (%) of state-of-the-art approaches on the Microsoft GrabCut database.

Method	Error rate
Graph cut [3]	6.6
Lazy snapping [14]	6.7
GrabCut [15]	5.4
DeepGrabCut [16]	7.7
Texture aware method [17]	3.6
Probabilistic graph matching [21]	5.7
RW [4]	6.4
Random walk with restart [23]	5.6
Lazy random walk [24]	5.7
SMRW [25]	4.6
Constrained random walk [26]	4.1
LC [27]	5.0
NRW [28]	5.9
Robust $P^p$ model [10]	6.1
NHO [9]	4.2
MGC [11]	3.4
Robust graph model [36]	3.8
Constrained dominant sets [5]	3.8
Ours	<b>2.7</b>

Fig. 9 shows an illustration of the significantly improved results of the proposed method compared with GrabCut [15], RW [4], NRW [28], LC [27], SMRW [25], NHO [9], MGC [11] and DeepGrabCut [16] on the Microsoft GrabCut database. It can be observed that there are weak boundary problems in all the test images in Fig. 9. The compared algorithms cannot perform well on these images because the local relationships utilized in them cannot accurately describe the similarities between neighboring pixels around the weak boundary regions. The proposed algorithm outperforms other approaches and has a significant improvement in error rates. The local relationship is extended to the global affinity by the proposed probabilistic diffusion process, which helps capture the intrinsic relationship of pixels around the weak boundary regions and makes our algorithm suitable for weak boundary problems.

### C. Sensitivity analysis

We analyzed the sensitivity of our method with respect to seed quantity and placement. The standard segmentations are

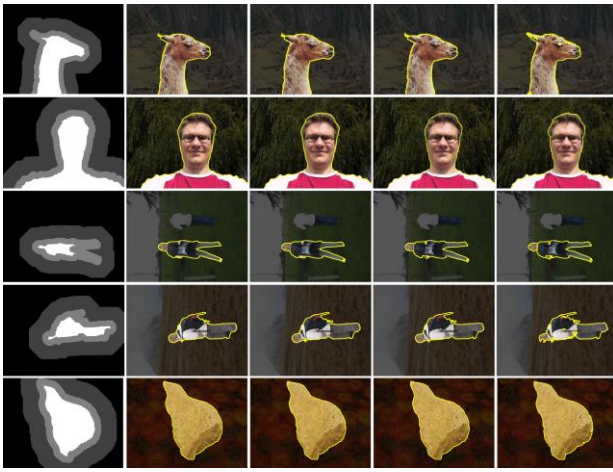


Fig. 10. Segmentation results with respect to the variation of seed quantity and placement. (a) Trimap input in the Microsoft GrabCut database; (b)–(e) perturbed segmentations with 50%, 30%, 10% and 1% percent seeds.

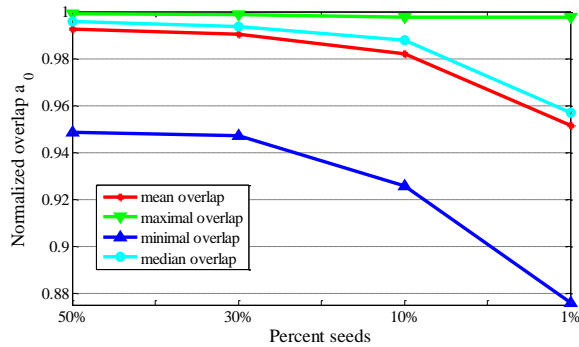


Fig. 11. Normalized overlap  $a_o$  on the Microsoft GrabCut database with 50%, 30%, 10% and 1% percent seeds.

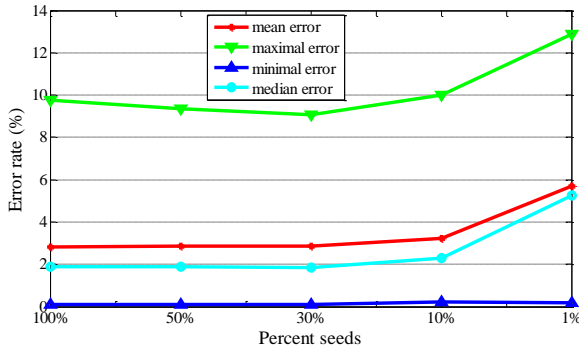


Fig. 12. Error rates (%) on the Microsoft GrabCut database with 50%, 30%, 10% and 1% percent seeds.

achieved from the initial trimaps provided by the Microsoft GrabCut database. The initial seeds are randomly taken from 50% to 1% of total seed quantity. The perturbed segmentations are recomputed from these selected seeds and compared with the standard segmentations. Let  $F_1$  and  $F_2$  denote a perturbed segmentation result and a standard segmentation result, respectively. The normalized overlap  $a_o$  is defined as  $a_o = |F_1 \cap F_2| / |F_1 \cup F_2|$ , which measures an overlap rate between a perturbed foreground and the corresponding standard foreground. Fig. 10 shows the comparison of example segmentations on the Microsoft GrabCut database with 50%, 30%, 10% and 1% percent seeds. It can be seen that almost the same results are produced, even with 1% seeds, the proposed

method can still obtain satisfactory segmentation results. Fig. 11 shows the quantitative evaluation of normalized overlap  $a_o$  on the Microsoft GrabCut database when varying the seed quantity as 50%, 30%, 10% and 1% of total seed quantity. The mean overlap rate is still over 0.95 when the percent seeds drops to 1%. Fig. 12 shows the quantitative evaluation of the error rates on the Microsoft GrabCut database with different percent seeds. From the variation of the error rates, we can find that the segmentation results are not sensitive to the variation of seed quantity and placement. These quantitative and qualitative experiments show that the proposed method has good robustness to the seeds.

#### D. Parameter settings

The parameter  $\alpha$  is used to control the influence of the region and boundary energies. Fig. 13 shows the segmentation results on the Microsoft GrabCut database by varying the value of  $\alpha$ . With a larger  $\alpha$ , region information plays a more important part and the details in objects can be preserved. However, as shown in Fig. 13(e), the boundaries are not smooth enough and it is hard to provide satisfactory segmentations. Comparatively smoother boundaries can be produced with a smaller  $\alpha$ . However, as shown in Fig. 13(a), the results may be over-smoothed and the details around boundaries cannot be preserved well. Therefore, it is important to determine an appropriate  $\alpha$  to improve boundary accuracy and reduce the over-smoothing effect. It can be seen that the best segmentation results can be obtained when  $\alpha = 0.2$ . Fig. 14 shows the quantitative evaluation of the error rates on the Microsoft GrabCut database with different values of  $\alpha$ . It can be observed that the lowest error rate is obtained when  $\alpha = 0.2$ . From the variation of the error rates, we can find that the segmentation results are somewhat sensitive to parameter  $\alpha$ .

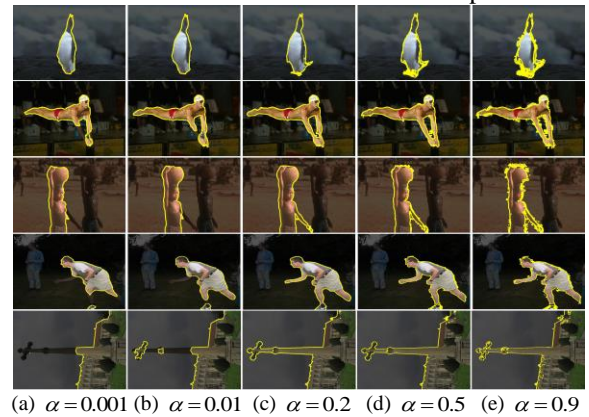


Fig. 13. Segmentation results with respect to the variation of parameter  $\alpha$ .

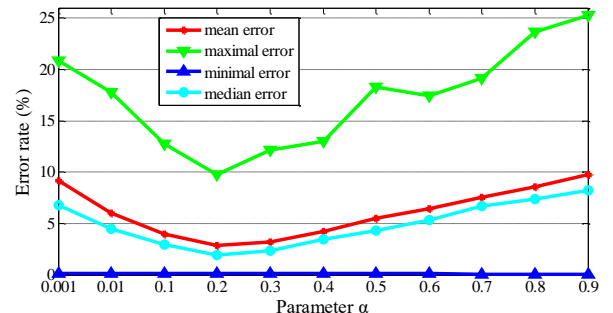
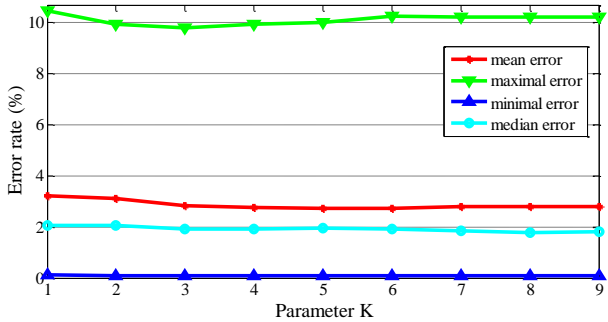


Fig. 14. Error rates (%) on the Microsoft GrabCut database by varying the values of parameter  $\alpha$ .

Fig. 15. Segmentation results with respect to the variation of parameter  $K$ .Fig. 16. Error rates (%) on the Microsoft GrabCut database by varying the value of parameter  $K$ .

The parameter  $K$  represents the number of clusters. Fig. 15 shows the segmentation results on the Microsoft GrabCut database by varying the value of  $K$ . It can be seen that almost the same results are produced with different values of  $K$ . Fig. 16 shows the quantitative evaluation of the error rates on the Microsoft GrabCut database with different values of  $K$ . It can be observed that the lowest error rate is obtained when  $K = 4$ . From the variation of error rate, we can find that the segmentation results are not sensitive to parameter  $K$ .

### E. Runtimes

TABLE III

Average running times (s) of GrabCut [15], RW [4], NRW [28], LC [27], SMRW [25], NHO [9], MGC [11] and the proposed method on all 20 images with size  $321 \times 481$  in the Microsoft GrabCut database.

	GrabCut	RW	NRW	LC	SMRW	NHO	MGC	Ours
Time	0.6	0.7	6.8	1.4	5.2	14.8	8.8	1.2

TABLE III lists the average run times of GrabCut [15], RW [4], NRW [28], LC [27], SMRW [25], NHO [9], MGC [11] and our method on all 20 test images (size:  $321 \times 481$ ) in the Microsoft GrabCut database on an Intel Core i7-7700K CPU with 16 GB memory running at 4.20 GHz in MATLAB R2017a. It can be seen that GrabCut and RW obtain the lowest run times. The run time of the proposed method is slightly higher than GrabCut and RW and is significantly lower than the perceptual grouping methods, NHO and MGC. The average number of iterations of the proposed method is 64, and the average run time is 1.2 s to segment an image with size  $321 \times 481$ . The algorithm complexity of the proposed method mainly focuses on the computation of the multiplication of matrix  $Q$  (size:

$N \times N$ ) and probability vector  $\hat{P}^L$  (size:  $N \times 1$ ). Since  $Q$  is a sparse matrix, this multiplication step is computationally efficient.

### F. Limitations

We estimate the initial foreground/background distributions from low-level features. Due to a lack of semantic information, the pixel-level features cannot be used to distinguish the foreground and background in images with similar foreground and background appearances. Fig. 17 shows an example; (a) and (c) show the same image with different scribbles (red: foreground, green: background), and (b) and (d) are the corresponding segmentation results. As seen, if we select all penguins as the foreground (shown in Fig. 17(a)), a satisfactory result can be produced with limited seeds (shown in Fig. 17(b)). If we select one penguin as the foreground (shown in Fig. 17(c)), it is difficult to remove the background of other penguins with similar appearances (shown in Fig. 17(d)). Because the deep features have a high-level understanding of objectness and semantics, combining the proposed algorithm with semantically aware information can be a good strategy to overcome the above limitations, but is out of the scope of this paper and would be considered in our future work.

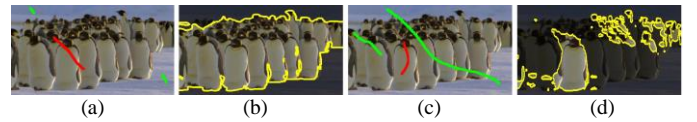


Fig. 17. Illustration of error segmentations of the proposed algorithm: (a) and (c) the same image with different scribbles (red: foreground; green: background); (b) and (d) the corresponding segmentation results.

## V. CONCLUSION

In this work, we presented an interactive approach for foreground/background image segmentation. The classification is formulated using a probabilistic framework consisting of unary potential estimation and likelihood learning. To improve the robustness to the seeds, the distances between pixel pairs and label pairs are measured to obtain prior label information. To improve segmentation accuracy, the region and boundary information are combined by a likelihood learning framework. An equivalence relation between likelihood learning and likelihood diffusion is also established, and an iterative diffusion-based optimization strategy is proposed to maintain computational efficiency. The qualitative and quantitative comparisons with state-of-the-art interactive approaches demonstrate the superior performance of the proposed method.

## REFERENCES

- [1] H. Noh, S. Hong, and B. Han, "Learning Deconvolution Network for Semantic Segmentation," in Proceedings of IEEE International Conference on Computer Vision, 2015, pp. 1520-1528.
- [2] E. Shelhamer, J. Long, and T. Darrell, "Fully convolutional networks for semantic segmentation," IEEE Transactions on Pattern Analysis & Machine Intelligence, vol.79, no.10, pp.1337-1342, 2014.
- [3] Y. Boykov, and M. Jolly, "Interactive graph cuts for optimal boundary & region segmentation of objects in ND images," in Proceedings of IEEE International Conference on Computer Vision, 2001, pp. 105-112.
- [4] L. Grady, "Random walks for image segmentation," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol.28, no.11, pp.1768-1783, 2006.



- [5] E. Zemene, and M. Pelillo, "Interactive image segmentation using constrained dominant sets," in *Proceedings of European Conference on Computer Vision*, 2016, pp. 278-294.
- [6] P. Arbeláez, and L. Cohen, "Constrained image segmentation from hierarchical boundaries," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2008, pp. 1-8.
- [7] Z. Li, and J. Chen, "Superpixel segmentation using Linear Spectral Clustering," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1356-1363.
- [8] M. V. D. Bergh, X. Boix, G. Roig, B. Capitan, L. V. Gool, "SEEDS: superpixels extracted via energy-driven sampling," in *Proceedings of European Conference on Computer Vision*, 2012, pp. 13-26.
- [9] T. Kim, K. Lee, and S. Lee, "Nonparametric higher-order learning for interactive segmentation," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2010, pp. 3201-3208.
- [10] P. Kohli, and P. Torr, "Robust higher order potentials for enforcing label consistency," *International Journal of Computer Vision*, vol.82, no.3, pp. 302-324, 2009.
- [11] T. Wang, Q. Sun, Z. Ji, Q. Chen, and P. Fu, "Multi-layer graph constraints for interactive image segmentation via game theory," *Pattern Recognition*, vol.55, pp.28-44, 2016.
- [12] X. Wang, Y. Tang, S. Masnou, and L. Chen, "A global/local affinity graph for image segmentation," *IEEE Transactions on Image Processing*, vol.24, no.4, pp.1399-1411, 2015.
- [13] T. Wang, Q. Sun, Q. Ge, Z. Ji, Q. Chen, and G. Xia, "Interactive Image Segmentation via Pairwise Likelihood Learning," in *Proceedings of International Joint Conference on Artificial Intelligence*, 2017, pp. 2957-2963.
- [14] Y. Li, J. Sun, C. Tang, and H. Shum, "Lazy snapping," *ACM Transactions on Graphics*, vol. 23, no.3, pp. 303-308, 2004.
- [15] C. Rother, V. Kolmogorov, and A. Blake, "Grabcut: Interactive foreground extraction using iterated graph cuts," in *Proceedings of the ACM SIGGRAPH Conference*, 2004, pp. 309-314.
- [16] N. Xu, B. Price, S. Cohen, J. Yang, and T. Huang, "Deep GrabCut for Object Selection," in *Proceedings of British Machine Vision Conference*, 2017.
- [17] S. Han, W. Tao, D. Wang, X. Tai, and X. Wu, "Image segmentation based on GrabCut framework integrating multiscale nonlinear structure tensor," *IEEE Transactions on Image Processing*, vol.18, no.10, pp.2289-2302, 2009.
- [18] H. Zhou, J. Zheng, and L. Wei, "Texture aware image segmentation using graph cuts and active contours," *Pattern Recognition*, vol.46, no.6, pp. 1719-1733, 2012.
- [19] D. Freedman, and T. Zhang, "Interactive Graph Cut Based Segmentation with Shape Priors," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2005, pp. 755-762.
- [20] M. Jian, and C. Jung, "Interactive Image Segmentation Using Adaptive Constraint Propagation," *IEEE transactions on image processing*, vol.25, no.3, pp.1301-1311, 2016.
- [21] A. Heimowitz, and Y. Keller, "Image Segmentation via Probabilistic Graph Matching," *IEEE Transactions on Image Processing*, vol.25, no.10, pp.4743-4752, 2016.
- [22] L. Grady, and G. Funkalea, "Multi-label Image Segmentation for Medical Applications Based on Graph-Theoretic Electrical Potentials," in *Proceedings of Computer Vision and Mathematical Methods in Medical and Biomedical Image Analysis*, 2004, pp.230-245.
- [23] T. Kim, K. Lee, and S. Lee, "Generative image segmentation using random walks with restart," in *Proceedings of European Conference on Computer Vision*, 2008, pp. 264-275.
- [24] J. Shen, Y. Du, W. Wang, and X. Li, "Lazy random walks for superpixel segmentation," *IEEE Transactions on Image Processing*, vol. 23, no. 4, pp. 1451-1462, 2014.
- [25] X. Dong, J. Shen, L. Shao, and L. Gool, "Sub-Markov random walk for image segmentation," *IEEE Transactions on Image Processing*, vol.25, no.2, pp. 516-527, 2016.
- [26] W. Yang, J. Cai, J. Zheng, and J. Luo, "User-friendly interactive image segmentation through unified combinatorial user inputs," *IEEE Transactions on Image Processing*, vol.19, no.9, pp. 2470-2479, 2010.
- [27] W. Casaca, L. G. Nonato, and G. Taubin, "Laplacian Coordinates for Seeded Image Segmentation," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 384-391.
- [28] C. G. Bampis, P. Maragos, and A. C. Bovik, "Graph-Driven Diffusion and Random Walk Schemes for Image Segmentation," *IEEE Transactions on Image Processing*, vol.26, no.1, pp. 35-50, 2016.
- [29] T. Wang, Q. Chen, Z. Ji, and Q. Sun, "Label Propagation and Higher-order Constraint based Segmentation of Fluid-associated Region in Retain SD-OCT Images," *Information Sciences*, vol.358, pp.92-111, 2016.
- [30] D. Zhou, O. Bousquet, T. N. Lal, J. Weston, and B. Scholkopf, "Learning with Local and Global Consistency," *Advances in Neural Information Processing Systems*, vol.16, no.4, pp.321-328, 2004.
- [31] T. Wang, Q. Sun, Z. Ji, Q. Chen, Q. Ge, and J. Yang, "Diffusive likelihood for interactive image segmentation," *Pattern Recognition*, vol.79, pp.440-451, 2018.
- [32] V. Lempitsky, P. Kohli, C. Rother, and T. Sharp, "Image segmentation with a bounding box prior," in *Proceedings of IEEE International Conference on Computer Vision*, 2009, pp. 277-284.
- [33] R. Unnikrishnan, C. Pantofaru, and M. Hebert, "Toward objective evaluation of image segmentation algorithms," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 29, no. 6, pp. 929-944, 2007.
- [34] M. Meilă, "Comparing clusterings: an axiomatic view," in *Proceedings of International conference on Machine learning*, 2005, pp. 577-584.
- [35] M. Friedman, "A comparison of alternative tests of significance for the problem of m rankings," *Annals of Mathematical Statistics*, vol.11, no.1, pp. 86-92, 1940.
- [36] T. Wang, Z. Ji, Q. Sun, Q. Chen, and X. Jing, "Interactive multi-label image segmentation via robust multi-layer graph constraints," *IEEE Transactions on Multimedia*, vol.18, no.12, pp.2358-2371, 2016.