

# Self-Adaptively Weighted Co-Saliency Detection via Rank Constraint

Xiaochun Cao, *Senior Member, IEEE*, Zhiqiang Tao, Bao Zhang, Huazhu Fu, and Wei Feng *Member, IEEE*

**Abstract**—Co-saliency detection aims at discovering the common salient objects existing in multiple images. Most existing methods combine multiple saliency cues based on fixed weights, and ignore the intrinsic relationship of these cues. In this paper, we provide a general saliency map fusion framework, which exploits the relationship of multiple saliency cues and obtains the self-adaptive weight to generate the final saliency/co-saliency map. Given a group of images with similar objects, our method first utilizes several saliency detection algorithms to generate a group of saliency maps for all the images. The feature representation of the co-salient regions should be both similar and consistent. Therefore, the matrix jointing these feature histograms appears low rank. We formalize this general consistency criterion as the rank constraint, and propose two consistency energy to describe it, which are based on low rank matrix approximation and low rank matrix recovery, respectively. By calculating the self-adaptive weight based on the consistency energy, we highlight the common salient regions. Our method is valid for more than two input images and also works well for single image saliency detection. Experimental results on a variety of benchmark data sets demonstrate that the proposed method outperforms the state-of-the-art methods.

**Index Terms**—Saliency detection, co-saliency detection, rank constraint, low-rank.

## I. INTRODUCTION

SALIENCY detection method simulates the behavior of our early primate visual system, which captures the most salient region of a scene. It could be considered as a preferential allocation of computational resources [1]–[5]. Since proposed in [1], saliency has been applied in a broad range of

Manuscript received September 19, 2013; revised January 22, 2014; accepted June 11, 2014. Date of publication June 23, 2014; date of current version August 21, 2014. This work was supported in part by the National Natural Science Foundation of China under Grant 61332012 and Grant 61100121, in part by the National Basic Research Program of China under Grant 2013CB329305, in part by the National High-Tech Research and Development Program of China under Grant 2014BAK11B03, and in part by the 100 Talents Programme of The Chinese Academy of Sciences. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Sylvain Paris.

X. Cao is with the School of Computer Science and Technology, Tianjin University, Tianjin 300072, China, and also with the State Key Laboratory of Information Security, Institute of Information Engineering, Chinese Academy of Sciences, Beijing 100093, China (e-mail: caoxiaochun@iie.ac.cn).

Z. Tao, B. Zhang, and W. Feng are with the School of Computer Science and Technology, Tianjin University, Tianjin 300072, China (e-mail: zqtao@tju.edu.cn; zhangbao@tju.edu.cn; wfeng@tju.edu.cn).

H. Fu is with the School of Computer Engineering, Nanyang Technological University, Singapore 639798 (e-mail: hzfu@ntu.edu.sg).

This paper has supplementary downloadable material available at <http://ieeexplore.ieee.org>, provided by the author. The supplementary file contains additional material not included in the paper. The total size is 22.8 MB. Contact zqtaomail@gmail.com for further questions about this work.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TIP.2014.2332399

applications, such as object recognition [6], object segmentation [7], image compression [8] and image retrieval [9]. Most existing saliency models [3], [10]–[15] focus on detecting the salient object from an individual image, which achieve encouraging results on the public benchmarks. In recent years, co-saliency has shown its value in many image processing and multimedia tasks. Jacobs *et al.* [16] firstly defined co-saliency detection as discovering the unique object among the images with similar backgrounds. However, the requirement of the similar backgrounds limits its applicability. Thus, we prefer an alternative concept, which aims at discovering the common salient object regions existing in multiple images. In detail, we follow the definition of co-saliency in [17]–[20], which simultaneously exhibits two properties, *i.e.*, 1) the co-salient region of each image should have strong local saliency with respect to its surroundings, and 2) all these co-salient regions should be similar in appearance. The later definition is more useful in various applications, including co-segmentation [21], [22], common pattern discovery [23], [24], co-recognition [25], [26], and image retrieval [27]. More application possibilities can be found in the work of [20]. Compared with saliency detection method on the single image, which only considers the contrast, uniqueness or priors in an individual image, co-saliency takes account of the relevance information of the similar salient objects in a group of images.

Existing co-saliency detection methods [18], [20] generate the final co-saliency maps by directly combining various multi-image saliency maps and single-image saliency maps. The multi-image saliency map is defined as the inter-image correspondence among multiple images, which is often obtained by feature matching [18] or clustering [20]. The single-image saliency map is used to describe the intra-image saliency in an individual image, and is obtained by performing any single image saliency method. We observe that the fusion step has a direct and significant influence on the co-saliency detection result. However, the confusion of similar backgrounds and the different advantages of various methods may lead these saliency maps to disagree with each other. As a result, one key in the fusion process is how to weight the maps to highlight the co-salient regions. Fig. 1 illustrates two co-saliency examples, where (b–c) are two multi-image saliency maps provided by [18], and (d–e) are two single-image saliency maps computed by [3] and [15], respectively. The multi-image saliency maps perform well in the first image pair for predicting the major salient regions. In contrast, confused by the complex background of the second image (the second row), the single-image saliency maps detect the entire object at the cost of obtaining many non-salient regions. However, for the second image pair,

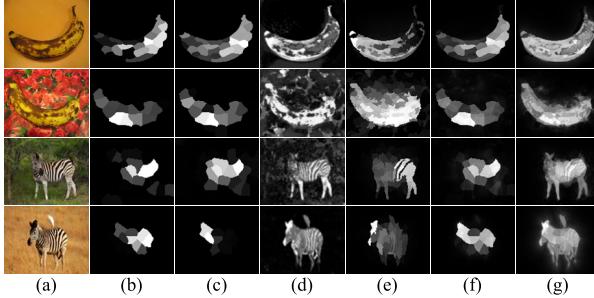


Fig. 1. Visual comparison of various saliency maps. (a) Image pairs: banana (upper) and zebra (bottom). (b)-(c) Multi-image saliency maps proposed by [18]. (d) Single-image saliency map generated by [3]. (e) Single-image saliency map generated by [15]. (f) Co-saliency map with the fixed weight in [18]. (g) Our co-saliency map under the rank constraint.

the multi-image saliency maps lose their effectiveness and the single-image saliency maps perform better. Overall, the performance of each saliency map varies with input images, and the fixed weight fusion is not the best choice for generating the final co-saliency map as shown in Fig. 1(f).

In this paper, we focus on the multiple saliency cues fusion. Different methods utilize various saliency cues to compute saliency maps. We effectively integrate these multiple saliency cues by weighting the maps self-adaptively. The feature histograms are employed to describe all the salient regions and stacked as a feature matrix. Ideally, the features of co-salient regions should be both similar and consistent, so the matrix rank appears low. We formalize this consistency property as the rank constraint. Two kinds of consistency energy are proposed to describe it, which are based on low rank matrix approximation and low rank matrix recovery, respectively. We call the first rank-one energy and the other low-rank energy for short. In fact, low-rank is a supplementary version to explain the rank constraint, which relieves the rank-one condition to low rank. The self-adaptive weights of the fusion maps are calculated by the consistency energy. As shown in Fig. 1, we improve the performance of co-saliency detection by combining multi-image saliency maps and single-image saliency maps self-adaptively. Our method performs well for both co-saliency and saliency detection, and is applicable for any number of images.

The rest of our paper is organized as follows. After a brief introduction of the related work about saliency and co-saliency in Section I-A, more details of our method are given in Section II, including rank constraint, consistency energy calculation, and co-saliency (saliency) fusion. Experimental results on the benchmark datasets are shown in Section III, followed by the conclusion in Section IV.

### A. Related Work

Generally, there are two major research directions in visual saliency modeling: human fixation prediction [28]–[30] and salient object region detection [3], [10]–[15]. The human fixation prediction model is used to predict the locations where the observer may fixate. It produces a saliency distribution map based on the eye fixation points. In contrast, the salient object detection model prefers to detect the salient object, which is a meaningful foreground. One of the important works in human

fixation prediction is proposed by Judd *et al.* [30]. They train a saliency model with a set of low, mid and high-level image features. Different from their work, our algorithm focuses on detecting salient objects, and fuses known multiple saliency cues without a learning process.

Salient object region detection in its essence is a segmentation problem [31]. However, compared to the traditional image foreground segmentation, which generally partitions an image into several regions with certain labels, a saliency detection method is based on human visual attention system and predicts the salient foreground localization automatically. The output of a saliency detection method is usually a saliency distribution map, which infers the probability of each pixel belonging salient in the image [31]. In this paper, we focus on detecting the common salient object regions in a group of images. The output of our co-saliency model is also a group of saliency maps, where each map offers a “soft” co-salient object segmentation.

A similar research work to co-saliency is co-segmentation [21], [22], [32], [33], which aims to segment out the common objects from multiple images. However, compared with the co-segmentation, co-saliency provides the saliency distribution maps to predict the common salient objects, based on the visually salient stimuli [20]. In practice, co-saliency could be used as an effective preprocessing step for co-segmentation [19]. Here, we briefly introduce the related work on single image saliency detection and co-saliency detection.

There are a lot of saliency detection methods for a single image. Various local and global contrast-based algorithms have been proposed. Cheng *et al.* [3] compute saliency values of image pixels using color statistics histogram contrast (HC), and evaluate the global contrast differences with the spatial coherence. To obtain a uniform high-response map for small-scale high-contrast images, Yan *et al.* [14] provide a hierarchical saliency (HS) detection method with three image layers in different scales, where the saliency cue in each layer includes local contrast. Besides the contrast-based ones, some methods utilize the frequency-domain features. Hou and Zhang [28] calculate the residual of the image log-spectrum, compared with an average log-spectrum from a set of natural images, to obtain saliency map (SR). By subtracting the average color from the low-pass filtered input, Achanta *et al.* [11] provide a frequency tuned algorithm (FT) to directly assign pixel saliency. Numerous saliency models are graph-based. Harel *et al.* [29] employ random walks on their designed graph to compute visual saliency (GB). Yang *et al.* [15] generate saliency map (MR) by performing manifold ranking on a close-loop graph with superpixels as nodes, which are ranked based on the similarity to background and foreground queries. Recently, low rank matrix recovery is employed for saliency detection in [13] and [34]. Shen *et al.* [13] provide a unified approach (LR) which incorporates low-level visual features with higher-level priors. In their model, an image is decomposed via low rank matrix recovery, and salient regions are explained as sparse noises. However, all these methods concentrate on single images without considering the consistency property among the multiple images.

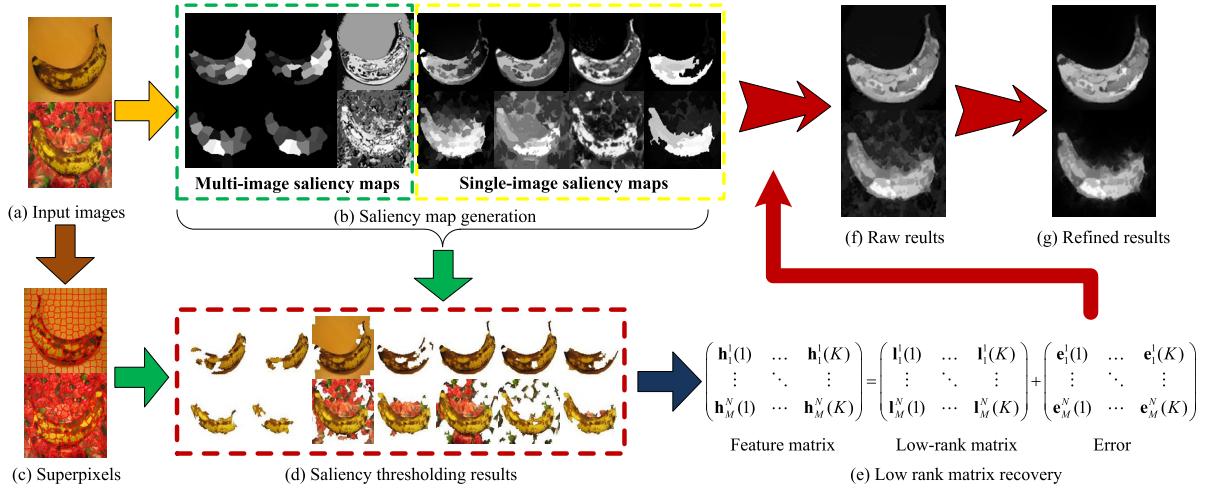


Fig. 2. The framework of our co-saliency detection method. (a) The  $N$  input images. (b) Saliency maps computed by  $M$  saliency methods. (c) The result of over-segmentation. (d) Salient regions are extracted by saliency thresholding. (e) Rank constraint with low-rank energy. (f) Our raw results. (g) Our final results after spatial refinement.

Few existing methods [17]–[20] are related to co-saliency detection. Chen [17] employs a distribution-based representation to characterize the sparse features in an image. Based on the feature distribution of a pair of images, it presents a progressive algorithm to enhance pre-attentive responses and thus can identify the common salient objects in both images. In [18]–[20], the co-saliency map is generated by combining multiple saliency cues. Li *et al.* [18] use two kinds of saliency maps (*i.e.*, single-image saliency maps and multi-image saliency maps) to define co-saliency as similarity scores, by performing single-pair SimRank algorithm on a co-multilayer graph. Three types of single-image saliency maps are adopted in their work, which are Itti’s model (IT) [1], frequency-tuned saliency (FT) [11], and spectral residual saliency (SR) [28]. The final map in [18] is modeled as a linear combination of three single-image saliency maps and two multi-image saliency maps with fixed weights. These two efficient approaches [17], [18], however, are only available for image pairs. It is hard to extend them to the case of multiple images. Chang *et al.* [19] consider the single-view saliency map and concentrate on those salient parts that frequently repeat in most images. But if the single-view saliency map is invalid, the co-saliency result will degenerate. Fu *et al.* [20] utilize clustering process to learn global correspondence information among the multiple images. Three bottom-up saliency cues (*i.e.*, contrast, spatial, and corresponding) are proposed to measure the cluster saliency, and the co-saliency map is obtained by fusing these cues with multiplication. However, these methods mostly utilize the fixed weights to combine the saliency maps/cues, by summation [18] or multiplication [19], [20]. They ignore the intrinsic relationship among the multiple saliency cues. In our method, we employ the rank constraint to integrate the superiority of each cue by weighting it self-adaptively.

## II. OUR PROPOSED APPROACH

As mentioned above, our method self-adaptively weights each saliency map that participates in fusion process under

the rank constraint. Fig. 2 shows the flowchart of our method, consisting of six steps. The first two steps (b), (c) generate saliency maps and superpixels for all the images, respectively. Various saliency detection methods are used in our framework, and their outputs (Fig. 2(b)) are employed as the elemental saliency maps. Next, we perform saliency thresholding to extract the salient regions for each image, by selecting salient superpixels based on different maps. The cut results are shown in Fig. 2(d). The histogram is utilized to represent the feature statistics of the salient region. As in Fig. 2(e), we stack all the histograms to form a feature matrix, which is used to calculate the consistency energy of each elementary map. Finally, we combine all the maps with self-adaptive weights decided by the energy values to generate the raw co-saliency map, and then refine the raw map to obtain our final one.

**Notations:** Given  $N$  input images  $\{I^i\}_{i=1}^N$ , we employ  $M$  saliency detection methods to compute  $N \times M$  saliency maps  $S_j^i$ ,  $1 \leq i \leq N$ ,  $1 \leq j \leq M$ , where  $i$  and  $j$  denote the index of image and saliency method respectively. Let  $\mathcal{S}^i = \{S_j^i\}_{j=1}^M$  be a map set of the  $i^{th}$  image and  $\mathcal{S}_j = \{S_j^i\}_{i=1}^N$  correspond to the  $j^{th}$  saliency detection method. The maps in  $\mathcal{S}^i$  are computed by  $M$  different methods for the image  $I^i$  as shown in each row of Fig. 2(b), and the ones in  $\mathcal{S}_j$  are generated from  $N$  images by the  $j^{th}$  method as shown in each column of Fig. 2(b). The superscript  $i$  and subscript  $j$  are used with other terms in this paper as well.

### A. Saliency Thresholding

The saliency thresholding aims at extracting the salient region detected by a saliency map. We perform over-segmentation to image  $I^i$  by [35], and decompose  $I^i$  into a set of superpixels  $\mathcal{X}^i = \{x_k^i\}_{k=1}^{n^i}$ , where  $n^i$  is the number of all the superpixels in image  $I^i$ . A binary map  $B_j^i$  is defined by thresholding on  $S_j^i$ :

$$B_j^i(x) = \begin{cases} 1, & \text{mean}(S_j^i(x)) > T_j^i, \\ 0, & \text{otherwise} \end{cases}, \quad (1)$$

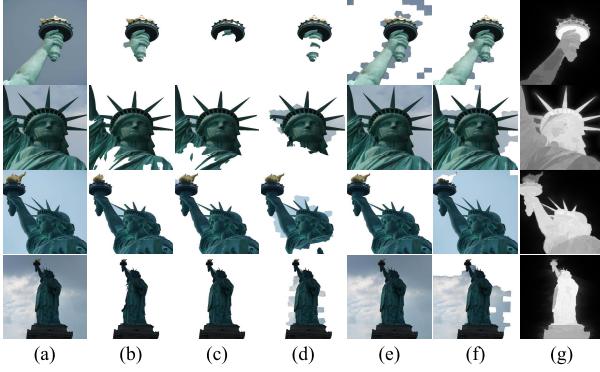


Fig. 3. Illustration of the rank constraint. (a) The input images. (b)-(f) The salient regions extracted by HS [14], MR [15], RC [3], CO [20] and SP [20], respectively. (g) Our final results.

where  $x \in \mathcal{X}^i$ ,  $\text{mean}(\cdot)$  computes the mean saliency score of all pixels in superpixel  $x$  by saliency map  $S_j^i$ , and the threshold  $T_j^i$  is computed by:

$$T_j^i = \alpha \cdot \max\{\text{mean}(S_j^i(x_k^i))\}_{k=1}^{n^i}, \quad (2)$$

where  $\alpha$  is used to control the threshold. We set  $\alpha = 0.3$  according to the evaluation in our experiment. The salient region<sup>1</sup>  $f_j^i$  (an image in Fig. 2(d)) is defined as a set of candidate salient superpixels decided by the  $j^{th}$  saliency detection method in the image  $I^i$ , and it is obtained by:

$$f_j^i = I^i \cdot B_j^i, \quad (3)$$

where  $f_j^i$  may not be spatially connected. Note that, each saliency map set  $\mathcal{S}^i$  has a group of salient regions  $\mathcal{F}^i = \{f_j^i\}_{j=1}^M$  (a row in Fig. 2(d)), and each  $\mathcal{S}_j$  is corresponding to  $\mathcal{F}_j = \{f_j^i\}_{i=1}^N$  (a column in Fig. 2(d)).

### B. Rank Constraint

As illustrated in Fig. 3(b-f), the  $i^{th}$  row is denoted as  $\mathcal{F}^i$  and the  $j^{th}$  column is corresponding to  $\mathcal{F}_j$ . First, based on the definition of co-saliency in this work, the salient regions in  $\mathcal{F}_j$  should exhibit similarity and consistency property. To measure this property, a 3-D color histogram is employed to represent each salient region,  $f_j^i$ , as a  $K$ -bin histogram  $\mathbf{h}_j^i \in \mathbb{R}^{(K \times 1)}$ . The 3-D color histogram is computed based on the three channels of the RGB color space. In detail, each color channel is uniformly quantized into 10 bins, and there is a total of  $K = 10^3$  bins. Then, we stack the  $N$  histograms of the salient regions in  $\mathcal{F}_j$  to get the matrix  $\mathbf{H}_j$ , where  $\mathbf{H}_j = [\mathbf{h}_j^1 \mathbf{h}_j^2 \dots \mathbf{h}_j^N]^T \in \mathbb{R}^{(N \times K)}$ . How much the  $j^{th}$  method contributes to the co-saliency detection could be measured by the consistency of rows in the feature matrix  $\mathbf{H}_j$ .

Since the salient regions in  $\mathcal{F}_j$  should be similar, the rank of  $\mathbf{H}_j$  is expected to be one ideally. In practice, saliency thresholding is unable to precisely separate the salient regions from original images. In addition, the color feature representation of the co-salient regions is sensitive to noises. Therefore,

<sup>1</sup>In practice, the ‘salient region’ in this paper is a set of superpixels in an image and may not be spatially connected, which slightly differs from the general concept of a connected region.

$\mathbf{H}_j$  is a low rank matrix. We formalize this property as rank constraint. An important benefit of this constraint is that it transforms the consistency measure into measuring the rank of the feature matrix. We offer the rank-one energy to calculate the consistency of  $\mathbf{H}_j$ , which decides the self-adaptive weight of each map in the map set  $\mathcal{S}_j$ .

Another crucial observation in Fig. 3 is that the similarity exists not only in  $\mathcal{F}_j$ , but also in  $\mathcal{F}^i$ . The feature matrix  $\mathbf{H}^i$  is utilized to denote  $\mathcal{F}^i$ , where  $\mathbf{H}^i = [\mathbf{h}_1^i \mathbf{h}_2^i \dots \mathbf{h}_M^i]^T \in \mathbb{R}^{(M \times K)}$ . Since the salient regions in  $\mathcal{F}^i$  are used to represent the same salient object from the image  $I^i$ , rank constraint is also applicable for  $\mathbf{H}^i$ .  $\mathcal{F} = \{\mathcal{F}^i\}_{i=1}^N$ , corresponds to all the salient regions in Fig. 3(b-f), and is described by the feature matrix:

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}^1 \\ \mathbf{H}^2 \\ \vdots \\ \mathbf{H}^N \end{bmatrix}, \quad (4)$$

where  $\mathbf{H} \in \mathbb{R}^{(NM \times K)}$ . Note that, each matrix  $\mathbf{H}^i$  is used to represent the salient object of an image  $I^i$ , and all these  $N$  observations exhibit the co-saliency property. Therefore, the matrix  $\mathbf{H}$  also satisfies rank constraint.

1) *Rank-One Energy*: Low rank matrix approximation is employed to compute the rank-one energy, which is an application of singular value decomposition (SVD) theorem. Suppose that the rank of  $\mathbf{H}_j$  is equal to  $r$ ,  $1 \leq r \leq N$ . We first decompose the matrix  $\mathbf{H}_j$  by SVD method as:

$$\mathbf{H}_j = [\mathbf{u}_1 \mathbf{u}_2 \dots \mathbf{u}_N] \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_N \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_K^T \end{bmatrix} \\ = \sum_{t=1}^N \sigma_t \mathbf{A}_t, \quad (5)$$

where  $\mathbf{u}_t \in \mathbb{R}^{N \times 1}$ ,  $\mathbf{v}_t \in \mathbb{R}^{K \times 1}$ ,  $\mathbf{A}_t = \mathbf{u}_t \mathbf{v}_t^T \in \mathbb{R}^{(N \times K)}$ ,  $\|\mathbf{A}_t\|_2 = 1$  and  $1 \leq t \leq N$ . The singular matrix is a  $N \times K$  matrix, and we obtain the singular values in descending order, where  $\sigma_1 \geq \dots \geq \sigma_r \geq \sigma_{r+1} = \dots = \sigma_N = 0$ . The  $k$ -rank ( $1 \leq k \leq r$ ) matrix approximation of  $\mathbf{H}_j$  is calculated by:

$$\tilde{\mathbf{H}}_j^{(k)} = \sum_{t=1}^k \sigma_t \mathbf{A}_t, \quad (6)$$

where  $\tilde{\mathbf{H}}_j^{(k)}$  is the most close one to  $\mathbf{H}_j$  of all the  $k$ -rank matrixes in Frobenius norm. Under the rank constraint, the rank  $r$  should be close to 1, which indicates that the 1-rank matrix approximation  $\tilde{\mathbf{H}}_j^{(1)}$  approximates to  $\mathbf{H}_j$ . Therefore, the error

$$\|\mathbf{H}_j - \tilde{\mathbf{H}}_j^{(1)}\|_F = \sqrt{\sigma_2^2 + \dots + \sigma_r^2}$$

is close to zero, and the  $\sigma_2$  is small enough to elide the other singular values. When the set  $\mathcal{F}_j$  shows higher consistency, the rank  $r$  approaches to one more closely. Thus, we transform the consistency measure into calculating the similarity between  $\mathbf{H}_j$  and its 1-rank matrix approximation, where the similarity can

be computed by the quotient of these two matrixes. Note that,  $\|\mathbf{H}_j\|_F^2 = \sigma_1^2 + \dots + \sigma_r^2 \approx \sigma_1^2 + \sigma_2^2$ , and  $\|\tilde{\mathbf{H}}_j^{(1)}\|_F^2 = \sigma_1^2$ . Hence, the quotient of  $\|\mathbf{H}_j\|_F$  and  $\|\tilde{\mathbf{H}}_j^{(1)}\|_F$  is approximately equal to

$$\sqrt{1 + \sigma_2^2/\sigma_1^2},$$

which is mainly decided by the ratio of  $\sigma_2$  and  $\sigma_1$ . Therefore, the rank-one energy is defined as:

$$\xi_j = -\sigma_2/\sigma_1, \quad (7)$$

where  $\sigma_1$  and  $\sigma_2$  are the first two singular values of the feature matrix  $\mathbf{H}_j$ , and  $-1 \leq \xi_j \leq 0$ ,  $1 \leq j \leq M$ .

The rank-one energy is proposed to compute the self-adaptive weight of each elementary saliency map  $S_j^i$ , and we denote  $\xi_j^i$  as the energy of  $S_j^i$ . Since the energy  $\xi_j$  in Eq. 7 is based on the consistency of the set  $\mathcal{F}_j$  by the  $j^{th}$  method, each map in  $\mathcal{S}_j$  has the same rank-one energy, *i.e.*,  $\xi_j^1 = \dots = \xi_j^N = \xi_j$ .

2) *Low-Rank Energy*: However, the rank-one energy focuses on the consistency of a set  $\mathcal{F}_j$ . It is too strict to approximate the whole salient regions  $\mathcal{F}$ . Thus, we also provide the low-rank energy to describe the rank constraint as a more general version. As shown in Fig. 3, the extracted *statues* from the images by different maps have similar color feature with noisy backgrounds such as *sky*, so we consider  $\mathcal{F}$  as a combination of salient regions residing in a low dimensional space with backgrounds as sparse noises. The feature matrix  $\mathbf{H}$  is utilized to represent the feature statistics of  $\mathcal{F}$ . Hence,  $\mathbf{H}$  could be decomposed as a sum of two parts  $\mathbf{H} = \mathbf{L} + \mathbf{E}$ , where  $\mathbf{L}$  denotes the recovered low rank matrix corresponding to salient regions and  $\mathbf{E}$  represents the residual backgrounds. The analysis above is formulated as a low rank matrix recovery problem:

$$\begin{aligned} (\mathbf{L}^*, \mathbf{E}^*) &= \arg \min_{\mathbf{L}, \mathbf{E}} (\text{rank}(\mathbf{L}) + \lambda \|\mathbf{E}\|_0), \\ \text{s.t. } \mathbf{H} &= \mathbf{L} + \mathbf{E}, \end{aligned} \quad (8)$$

where  $\|\cdot\|_0$  denotes the  $\ell_0$ -norm, and  $\lambda$  is the coefficient controlling the weight of the sparsity of  $\mathbf{E}$ . Since the problem in Eq. 8 is intractable in polynomial-time, we solve its convex relaxation by Robust Principal Component Analysis (RPCA) [36], [37] as follows:

$$\begin{aligned} (\mathbf{L}^*, \mathbf{E}^*) &= \arg \min_{\mathbf{L}, \mathbf{E}} (\|\mathbf{L}\|_* + \lambda \|\mathbf{E}\|_1), \\ \text{s.t. } \mathbf{H} &= \mathbf{L} + \mathbf{E}, \end{aligned} \quad (9)$$

where  $\|\cdot\|_*$  indicates nuclear norm and  $\|\cdot\|_1$  is  $\ell_1$ -norm. The matrix  $\mathbf{E}$  is regarded as an error matrix between  $\mathbf{H}$  and the low rank matrix  $\mathbf{L}$ , where  $\mathbf{E} \in \mathbb{R}^{(NM \times K)}$ ,  $\mathbf{E}^i = [\mathbf{e}_1^i \ \mathbf{e}_2^i \ \dots \ \mathbf{e}_M^i]^T \in \mathbb{R}^{(M \times K)}$ , and  $\mathbf{e}_j^i \in \mathbb{R}^{(K \times 1)}$ ,  $1 \leq j \leq M$ ,  $1 \leq i \leq N$ .  $\mathbf{E}^i$  is the error matrix of image  $I^i$ . Based on the rank constraint, the matrix  $\mathbf{H}$  should be low rank, and the consistency of the whole salient regions  $\mathcal{F}$  could be computed by the error  $\mathbf{E}$ . In detail, the low-rank energy  $\xi_j^i$  of the saliency map  $S_j^i$  is defined as a histogram distance between the corresponding feature histogram  $\mathbf{h}_j^i$  and the recovered  $\mathbf{l}_j^i$  in  $\mathbf{L}$ . After an normalization of  $\mathbf{E}$ , we calculate this distance by:

$$\xi_j^i = -\|\mathbf{e}_j^i\|_2, \quad (10)$$

where  $(\mathbf{e}_j^i)^T$  is the  $((i-1) \times M + j)^{th}$  row of the matrix  $\mathbf{E}$ , and  $\|\cdot\|_2$  is the  $\ell_2$ -norm. The higher low-rank energy  $\xi_j^i$  indicates that the salient region detected by  $S_j^i$  is more likely to be the true co-salient region.

In summary, the rank-one energy is focused on each set  $\mathcal{F}_j$  by one detection method, and the low-rank considers all the salient regions simultaneously. These two energy both utilize the consistency among the images to explain our rank constraint. The difference is that the low-rank energy also employs the consistency property among different methods. As a result, the low-rank energy is more general than the rank-one, for it could be used in single saliency detection.

### C. Co-Saliency Assignment

Our goal is to generate the co-saliency maps for all the input images  $\{I^i\}_{i=1}^N$ . The co-saliency map of  $I^i$  is obtained by combining each map  $S_j^i$  in the map set  $\mathcal{S}^i$  with its own self-adaptive weight. The consistency energy  $\xi_j^i$  (*i.e.*, rank-one energy or low-rank energy) of  $S_j^i$  is utilized to define the self-adaptive weight as:

$$w_j^i = \frac{\exp(\xi_j^i)}{\sum_{j=1}^M \exp(\xi_j^i)}, \quad (11)$$

where  $w_j^i$  is the weight of  $S_j^i$ , and  $\sum_{j=1}^M w_j^i = 1$ ,  $0 \leq w_j^i \leq 1$ . An exponential function  $\exp(\cdot)$  is used to emphasize the  $\xi_j^i$  for it is of high significance and discriminative power. The saliency map with higher energy value is weighted more by Eq. 11. We denote  $CS^i$  as the co-saliency map of image  $I^i$ , and it is obtained by:

$$CS^i = \sum_{j=1}^M w_j^i \cdot S_j^i. \quad (12)$$

The co-saliency map  $CS^i$  is used to represent the saliency score of each pixel in image  $I^i$ .

Moreover, the proposed method is also applicable for single image saliency detection. Suppose  $N = 1$ , and we only have the map set  $\mathcal{S}^1 = \{S_j^1\}_{j=1}^M$ . Since the rank-one energy is computed by the consistency of a group of images, it is inapplicable for a single image. However, we can utilize the set  $\mathcal{F}^1$  to calculate the low-rank energy of each map in  $\mathcal{S}^1$ . In our fusion framework, single saliency map is obtained by recovering the common salient region shared by all the  $M$  elementary saliency maps, which could be also computed by Eq. 12 with the low-rank energy and  $i = 1$ .

### D. Spatial Refinement

As suggested by [7], the areas that are close to the foci of attention should be explored more significantly than far away regions, which indicates that the salient pixels gather together within an image. Thus, we refine the raw co-saliency map  $CS^i$  in two steps. First, we obtain a set of salient pixels for image  $I^i$  by a threshold, which is set to be  $0.3 * \max(CS^i)$ . Then, inspired by [7] and [38], the raw saliency score of each pixel  $p$  is refined as:

$$CS^{i*}(p) = CS^i(p) \cdot \exp(-\beta \cdot D(p)), \quad (13)$$

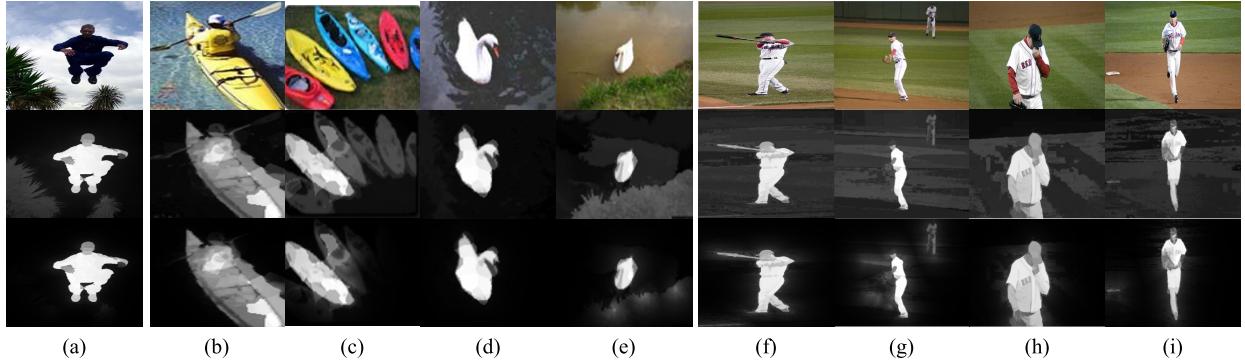


Fig. 4. Visual comparison between our raw saliency detection result and the refined one. The input images, our raw saliency/co-saliency maps and the maps after spatial refinement are shown in the first, the second and the third row, respectively. (a) Saliency detection on a single image. (b-e) Co-saliency detection on two image pairs. (f-i) Co-saliency detection on a set of four images.

where  $\beta$  is a scale parameter with  $\beta = 6$  according to [38].  $D(\cdot)$  is defined as the closest spatial positional distance between pixel  $p$  and the salient pixel set, which is normalized to the range  $[0,1]$ . After the spatial refinement,  $CS^{i*}$  is the final co-saliency map for the image  $I^i$ .

Fig. 4 demonstrates the power of the spatial refinement. Our raw saliency/co-saliency maps and the refined ones are shown in the second and third rows, respectively. Compared to the raw maps, the refined maps alleviate the noises of the backgrounds effectively. For example, the background *tree* of the input image in Fig. 4(a) is detected by our raw saliency map with low saliency score. However, the corresponding refined map eliminates this noise.

### III. EXPERIMENTS

In this section, We test our method on single image saliency detection, image pair co-saliency detection, and multiple image co-saliency detection, respectively. Qualitative and quantitative analyses of our results are presented. We also compare our framework with the state-of-the-art methods on a variety of benchmark datasets. Our following results (*i.e.*, visual results) are based on the low-rank energy by default, since it is more general than the rank-one energy.

#### A. Experimental Settings

1) *Datasets:* We test our algorithm on three public datasets: MSRA-B [10], [39], Image Pair [18] and CMU-Cornell iCoseg dataset [40]. The first one is mainly used for single image saliency detection (called single saliency for short) and the other two focus on co-saliency detection. MSRA-B dataset (5000 images) is one of the most largest saliency image databases, where accurate pixel-wise annotations have been provided by [39] to replace the labeled rectangles of it. The Image Pair dateset [18] collects 105 image pairs (*i.e.*, 210 images), and the ground truth is labeled manually. Each image pair contains one or more similar objects with different backgrounds. The CMU-Cornell iCoseg dataset [40], including 38 groups of totally 643 images, is provided along with pixel-wise ground truth annotations, where each group contains one or multiple salient objects.

2) *Baselines:* In our experiment, we adopt twelve state-of-the-art saliency algorithms for comparison, which consists of ten single image saliency methods: the Itti model (IT) [1], graph-based visual saliency (GB) [29], saliency using spatiotemporal cues (LC) [12], frequency-tuned saliency (FT) [11], spectral residual saliency (SR) [28], histogram based saliency (HC) [3], global-contrast saliency (RC) [3], saliency using low rank matrix recovery (LR) [13], hierarchical saliency model (HS) [14], graph-based manifold ranking saliency (MR) [15], and two co-saliency algorithms: co-saliency model of image pair (Li) [18] and cluster-based co-saliency model (Fu) [20]. Moreover, our method also compares with four multiple image saliency methods, which are denoted as CC, CP, SP, and CO. CC and CP provided by [18] are computed with the similarity of an image pair by color (CC) and texture (CP), respectively. SP and CO are generated by two saliency cues (*i.e.*, Spatial and Corresponding) proposed in [20]. LC, FT, SR, HC and RC are implemented by the software provided in [3]. For the others, we run the authors' codes, and use the default parameter settings.

3) *Elementary Maps:* The proposed method is based on integrating multiple elementary saliency maps self-adaptively. By employing different maps, various saliency cues are added into the fusion framework, which makes our method more robust. These elementary saliency maps are the outputs of various saliency methods. Our framework uses these outputs directly, which is helpful to integrate more methods into the framework. Three different groups of elementary maps are used in our algorithm:

- Single image saliency detection (G1): MR [15], HS [14], LR [13], RC [3], HC [3] and GB [29].
- Image pair co-saliency detection (G2): MR [15], HS [14], RC [3], HC [3], CC [18], CP [18] and SP [20].
- Multiple image co-saliency detection (G3): MR [15], HS [14], RC [3], SP [20] and CO [20].

MR, HS, LR, RC, HC, and GB are used to generate the single-image saliency maps, while CC, CP, SP and CO are used to generate the multi-image saliency maps. Since CC and CP are only applicable for a pair of images, they cannot be employed for multiple image co-saliency detection. For the convenience of statement, these three map groups are denoted as G1, G2, and G3, respectively.

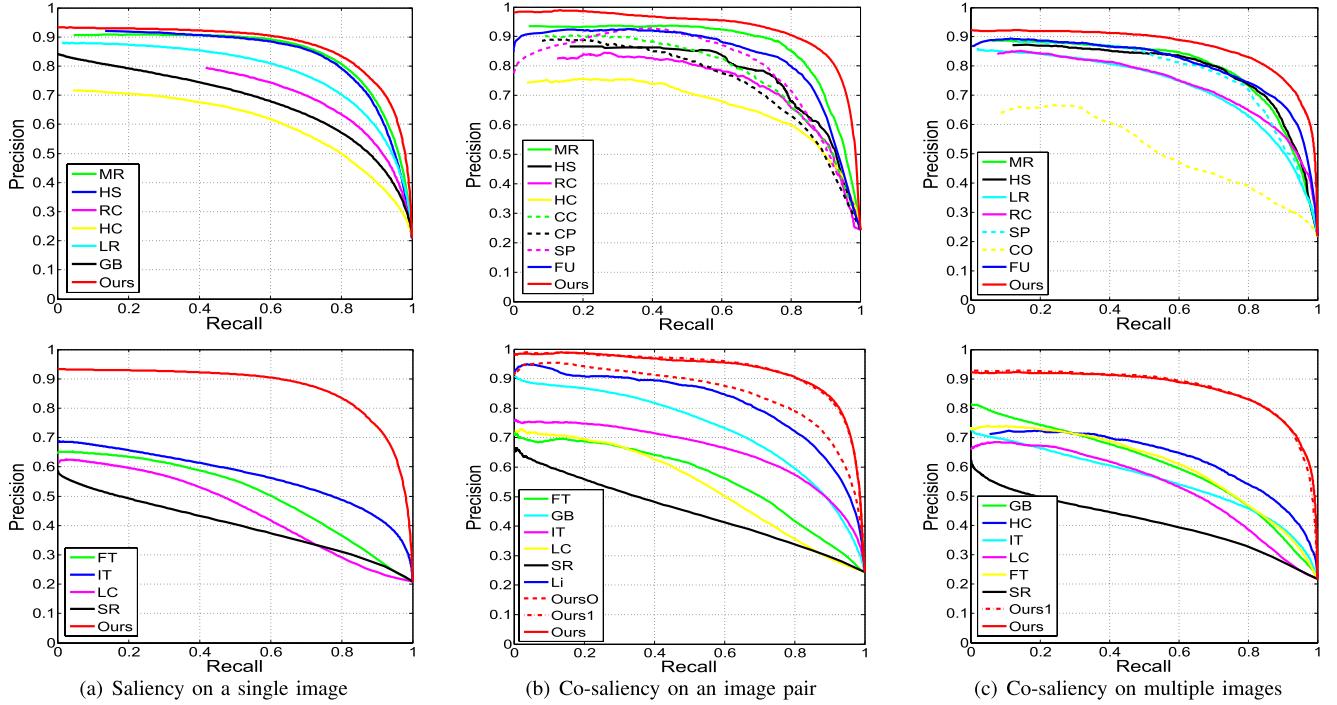


Fig. 5. Comparison of the PR curves between our results and other state-of-the-art works. To better illustrate the results, we separate the curves into two rows since the compared state-of-the-art methods are quite a few. (a) The PR curves of single saliency detection on MSRA-B dataset [10], [39]. We compare our method with ten single image saliency detection algorithms including MR [15], HS [14], LR [13], RC [3], HC [3], GB [29], FT [11], IT [1], LC [12] and SR [28]. (b) The PR curves of co-saliency detection on Image Pair dataset [18]. The other two co-saliency methods Li [18] and Fu [20] are added in the comparison. Our result generated by the same maps in [18] is denoted as OursO. (c) The PR curves of co-saliency detection on iCoseg dataset [40]. For better view, the multi-image saliency maps, CC [18], CP [18], SP [20] and CO [20], are shown as dotted lines in (b) and (c). Our results with the rank-one energy (Ours1) are demonstrated as dash-dotted lines in (b) and (c). (a) Saliency on a single image. (b) Co-saliency on an image pair. (c) Co-saliency on multiple images.

**4) Evaluation Metrics:** To evaluate the performance of our method, we test our results based on four widely used criteria, including PR (precision-recall) curve, ROC (receiver operating characteristic) curve, F-score [11], and AUC (Area Under ROC Curve) score. The precision is defined as the ratio of correctly detected salient pixels to all the detected pixels, and the recall is the fraction of detected salient pixels belonging to the salient pixels in ground truth. In our experiment, precision and recall are calculated by a fixed threshold. The PR curve is drawn by 256 precision-recall pairs, which are obtained by varying the threshold from 0 to 255. High recall always appears at the expense of reducing the precision, and vice-versa. Therefore, it is crucial to consider them together. F-score evaluates precision and recall simultaneously, and it is calculated by:

$$\text{F-score} = \frac{(1 + \gamma^2)\text{Precision} \times \text{Recall}}{\gamma^2 \times \text{Precision} + \text{Recall}}, \quad (14)$$

where we use  $\gamma^2 = 0.3$  as in [11] to weight precision more than recall. The ROC curve indicates the similarity between the predicted saliency map and the ground truth, whose performance is quantized by the AUC scores. We compute AUC scores based on the [41, Algorithm 3].

### B. Quantitative Analysis

We evaluate our algorithm on three aspects: saliency detection for single image, co-saliency detection for image pair,

and co-saliency detection for multiple images. As mentioned before, four criteria are employed in our quantitative analysis.

**1) Single Image Saliency Detection:** We first verify the performance of our method for single saliency detection on the MSRA-B dataset [10], [39]. Fig. 5(a) shows the PR curves of our algorithm and ten single image saliency detection methods. It can be seen that our method achieves the highest precision for any given recall value, which indicates that we have the best performance. The similar ROC curve is shown in Fig. 6(a), where our framework also performs better than the others. We employ the maximum F-score and AUC score to quantize the PR curve and ROC curve, respectively. From Table I, we observe that our method has the highest F-score and AUC score on the MRSA-B dataset. Among the ten single saliency methods, MR [15] and HS [14] are the top performers. MR [15] has the highest AUC score (0.9386), and the second F-score (0.8157). HS [14] has the best F-score (0.8207) and the second AUC score (0.9297). However, our results outperform these two methods with F-score = 0.8335 and AUC = 0.9625, which are 1.28% and 2.39% better than HS [14] and MR [15], respectively. Note that, our curve dominates the six saliency methods of the map group G1 in both PR space and ROC space. It proves that our rank constraint is effective to inherit the various merits of different saliency maps.

**2) Image Pair Co-Saliency Detection:** Our method achieves co-saliency detection via integrating various single-image

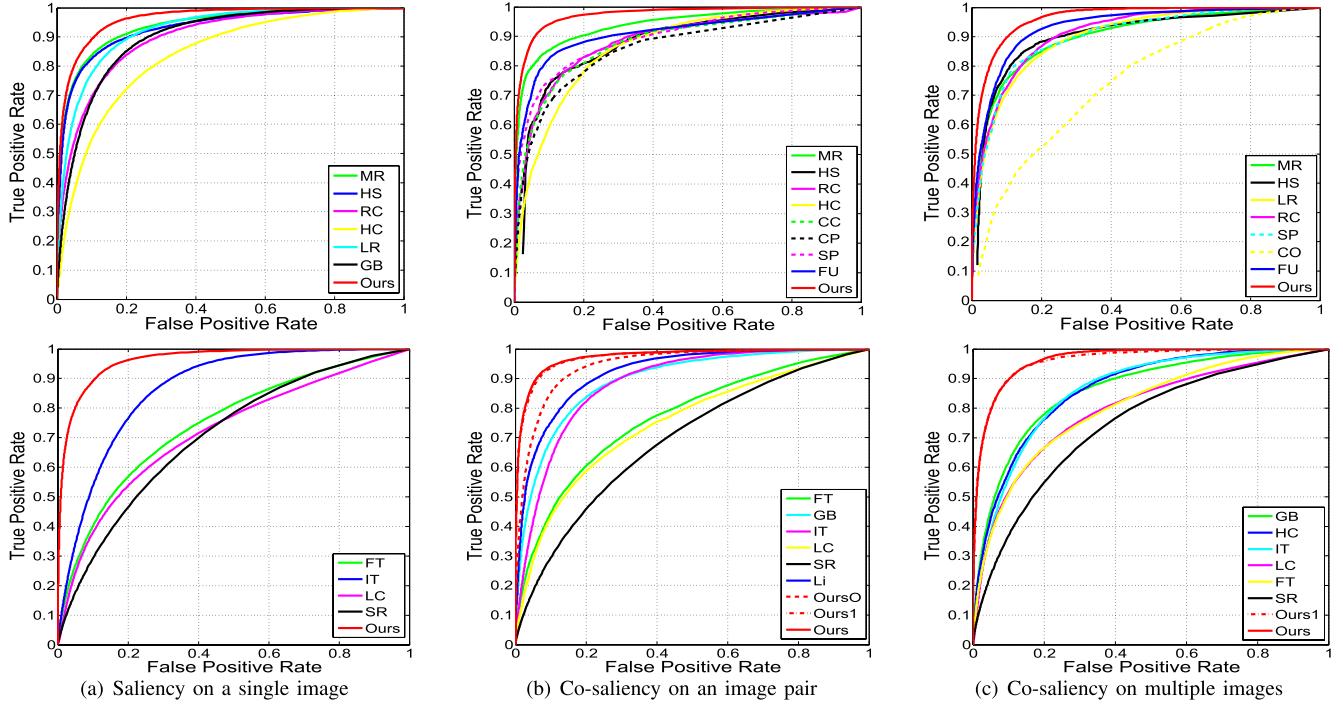


Fig. 6. Comparison of the ROC curves between our results and other state-of-the-art works. To better illustrate the results, we separate the curves into two rows since the compared state-of-the-art methods are quite a few. (a) The ROC curves of single saliency detection on MSRA-B dataset [10], [39]. We compare our method with ten single image saliency detection algorithms including MR [15], HS [14], LR [13], RC [3], HC [3], GB [29], FT [11], IT [1], LC [12] and SR [28]. (b) The ROC curves of co-saliency detection on Image Pair dataset [18]. The other two co-saliency methods Li [18] and Fu [20] are added in the comparison. Our result generated by the same maps in [18] is denoted as OursO. (c) The ROC curves of co-saliency detection on iCoseg dataset [40]. For better view, the multi-image saliency maps, CC [18], CP [18], SP [20] and CO [20], are shown as dotted lines in (b) and (c). Our results with the rank-one energy (Ours1) are demonstrated as dash-dotted lines in (b) and (c).

TABLE I  
COMPARISON OF F-SOCRE AND AUC BETWEEN OUR METHOD AND OTHER STATE-OF-THE-ART WORKS

Criteria	Datasets	MR	HS	LR	RC	HC	GB	Li	Fu	Ours
F-score	MSRA-B	0.8157	0.8207	0.7538	0.7529	0.6153	0.6592	-	-	<b>0.8335</b>
	Image Pair	0.8527	0.7684	-	0.7393	0.6607	0.6971	0.7765	0.8132	<b>0.8795</b>
	iCoseg	0.7807	0.7759	0.7089	0.7111	0.6362	0.6016	-	0.7688	<b>0.8249</b>
AUC	MSRA-B	0.9386	0.9297	0.9238	0.8964	0.8374	0.8984	-	-	<b>0.9625</b>
	Image Pair	0.9389	0.8782	-	0.8825	0.8654	0.8932	0.9207	0.9098	<b>0.9730</b>
	iCoseg	0.9017	0.9038	0.8998	0.9125	0.8643	0.8620	-	0.9348	<b>0.9639</b>
Criteria	Datasets	FT	IT	LC	SR	CC	CP	SP	CO	
F-score	MSRA-B	0.5403	0.5707	0.4939	0.4257	-	-	-	-	
	Image Pair	0.5825	0.6491	0.5590	0.4611	0.7603	0.7298	0.7939	-	
	iCoseg	0.6114	0.5562	0.5594	0.4374	-	-	0.7626	0.5464	
AUC	MSRA-B	0.7459	0.8653	0.7174	0.7050	-	-	-	-	
	Image Pair	0.7655	0.8833	0.7470	0.6897	0.8790	0.8595	0.8955	-	
	iCoseg	0.8030	0.8607	0.7961	0.7461	-	-	0.9008	0.7450	

saliency maps and multi-image saliency maps under the rank constraint. For a better view, we plot all the multiple image saliency methods using dotted lines in Fig. 5 and Fig. 6. We compare our proposed method with nine single image saliency methods, three multiple image saliency methods, and two co-saliency methods on the Image Pair dataset [18]. As demonstrated in Fig. 5(b), our PR curve outperforms the other methods and boosts the performance of co-saliency significantly. Fig. 6(b) shows the comparison result of the ROC curves, where our proposed algorithm also excels. The co-saliency model of Li [18] combines two multi-image saliency maps (*i.e.*, CC and CP) and three single-image saliency maps (*i.e.*, IT, FT, and SR) with the fixed weights

(0.4, 0.4, 0.067, 0.067, 0.067). Different from Li [18], our method weights each elementary map self-adaptively by the rank constraint. Thanks to the self-adaptive weight, the maps extracting the common salient regions are more accurately highlighted. We generate an additional result, named as OursO, by using the same five maps (*i.e.*, IT, FT, SR, CC and CP) as in [18]. The bottom rows in Fig. 5(b) and Fig. 6(b) show that OursO also performs better than Li. The comparison results demonstrate that our self-adaptive weight is superior to the fixed one.

3) *Multiple Image Co-Saliency Detection*: Most existing co-saliency detection methods focus on a pair of images, and only few co-saliency models are applicable for the multiple

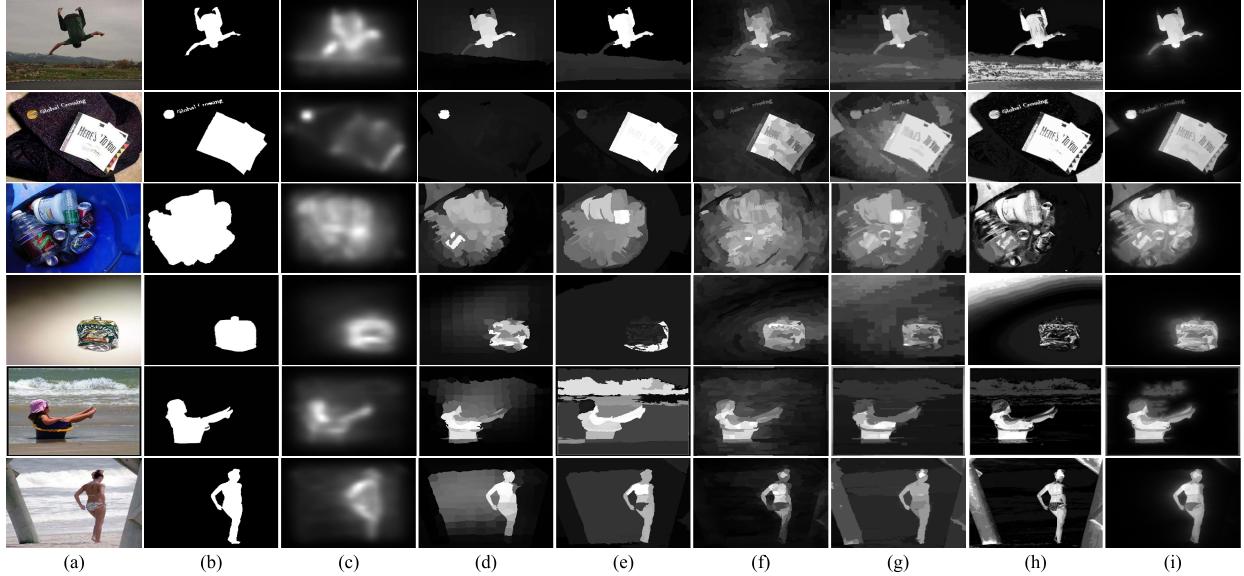


Fig. 7. Visual comparison of single image saliency detection on MSRA-B dataset [10], [39]. (a) Images. (b) Ground truth. (c) GB [13]. (d) MR [17]. (e) HS [16]. (f) LR [15]. (g) RC [3] (h) HC [3] (i) Ours.

images. Similar to ours, the recent cluster-based saliency model of Fu [20] is also not limited to the image number. We compare our method with the co-saliency method Fu [20], two multiple image saliency methods and ten single image saliency methods on the iCoseg dataset [40]. The results are shown in Fig. 5(c) and Fig. 6(c). Without surprise, our method has the highest PR curve and ROC curve. Moreover, the power of our co-saliency model ( $F\text{-score} = 0.8249$ ,  $AUC = 0.9639$ ) is shown in Table I, which is better than Fu ( $F\text{-score} = 0.7688$ ,  $AUC = 0.9348$ ) and excels the best two performers in G3 (*i.e.*, MR [15] with  $F\text{-score} = 0.7807$ , and HS [14] with  $AUC = 0.9038$ ). One may observe that SP [20] has some satisfactory curves. This agrees with the center prior that the objects near the image center are more attractive to people [13], [30], [42]. The other multiple image saliency method CO [20] is unsatisfactory in both Fig. 5(c) and Fig. 6(c). In fact, CO is used to describe how frequent the salient object recurs among the multiple images, which is an important global property of the common saliency [20]. However, the iCoseg dataset is provided for the co-segmentation, where the objects are always along with the alike backgrounds. Therefore, CO is confused by the recurred backgrounds.

In summary, benefiting from our rank constraint, our co-saliency detection model inherits the advantages of all the fusion maps effectively. It utilizes multiple saliency cues to generate our final results. Besides, the PR and ROC curves of our results with the rank-one energy are shown in Fig. 5 and Fig. 6 as dash-dotted lines, which indicate that our two energy are both effective for co-saliency detection.

### C. Qualitative Analysis

In this subsection, We give the visual comparison results between our framework and other competitors on three aspects.

1) *Single Image Saliency Detection:* We first show the results on the MSRA-B dataset [10], [39]. As shown in Fig. 7,

our results provide visually pleasing saliency, and approach the ground truth more closely than the other six saliency methods in G1. By reviewing the Section III-B, we find that all these methods have satisfactory performances, where MR [15] and HS [14] are the best two. However, various methods have quite different results for all kinds of scenarios. For example, GB [29] and LR [13] perform stably on the second and fourth images, while the detected salient objects are not outstanding. For another example, similar to RC [3], HC [3] has an encouraging result on images whose foregrounds have high-contrast to surrounding (*e.g.*, the second row), however, they are affected by the low-contrast situation (*e.g.*, the fourth row). Although MR [15] and HS [14] are overall better than the other four methods, they sometimes also lose their power. In detail, MR [15] only detects a small circle dot at the top left corner of the second image, and wrongly predicts many non-salient regions in the center of the fifth and sixth rows. On the other side, HS [14] obtains the major salient object in the second row, but fails in the fourth image. Without surprise, every saliency cue has its own available precondition. In contrast, our visual result has shown its robustness by combining multiple saliency cues self-adaptively. Overall, our method performs better than the others. The last row in Fig. 7 shows our superiority to other methods obviously.

2) *Image Pair Co-Saliency Detection:* Fig. 8 provides the visual co-saliency detection results of MR [15], Li [18], Fu [20] and ours. We choose seven pairs of images from the Image Pair dataset [18] to compare our method with the other three models. The single saliency method MR [15] detects salient objects through manifold ranking on a graph, which is constructed by each individual image. Since MR [15] predicts the saliency for one of the pairs independently, it ignores the corresponding information of the image pair. Therefore, its performance degenerates when the image has textured backgrounds or the low contrast foregrounds. For example, MR [15] misses the *dog* in the fifth row of Fig. 8. In contrast, co-saliency methods solve this problem by the multi-image

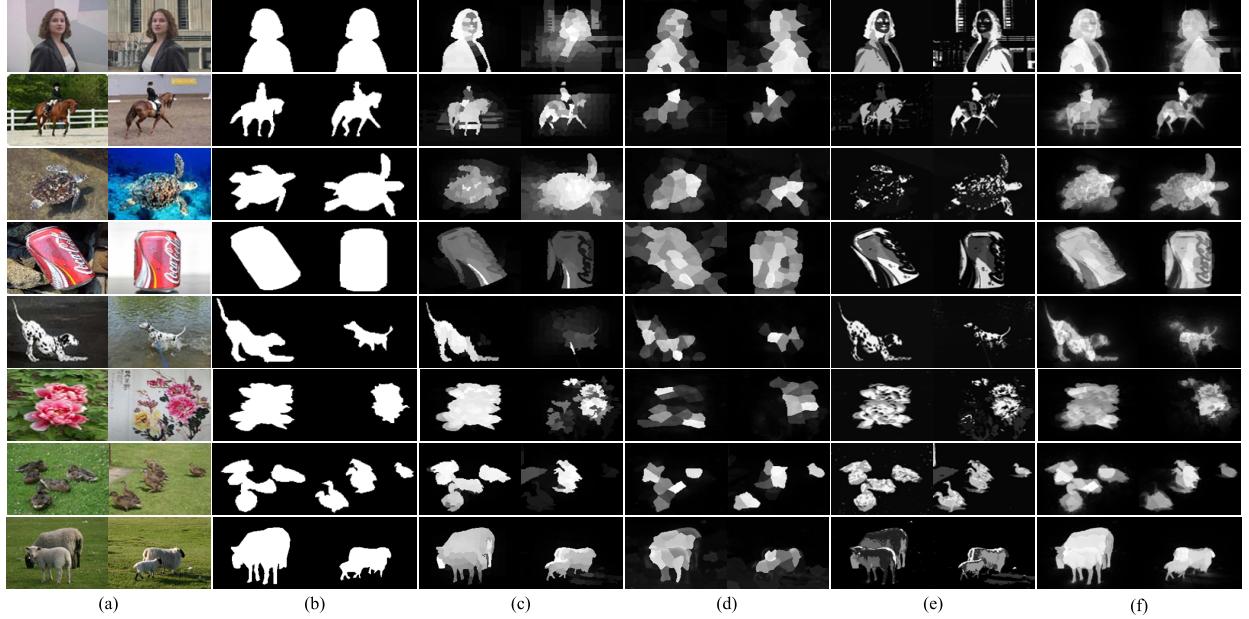


Fig. 8. Visual results of co-saliency detection on Image Pair dataset [18]. (a) Image pair. (b) Ground truth. (c) MR. (d) Li. (e) Fu. (f) Ours.



Fig. 9. Some visual results of our co-saliency detection for multiple images on iCoseg dataset [40]. For each block, the original images, the saliency maps by Fu [20] and our co-saliency maps are shown in the first, the second and the third row, respectively.

saliency maps efficiently. That is verified by the observation in Fig. 8 that all the three co-saliency models strengthen the common salient parts, such as the *dog* in the fifth row. However, our method outperforms Li [18] and Fu [20] for

obtaining more accurate salient objects with more clean backgrounds. The co-saliency model of Li [18] fixes each map with a constant weight. As a result, it is affected by some “outlier” saliency maps who lose their power. Compared to

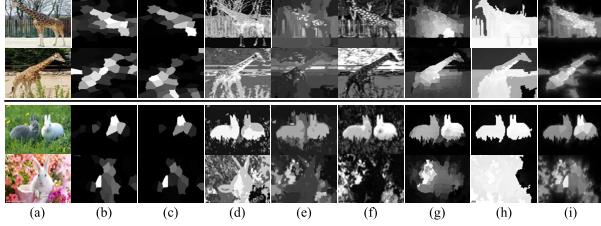


Fig. 10. Some challenging examples for our co-saliency detection. (a) Image pairs: giraffe (upper) and rabbit (bottom). (b)-(h) The elementary maps. (i) Our results for the image pairs.

Li [18], our model highlights the common salient regions by adjusting the contribution of the multiple saliency cues. As demonstrated in the sixth row in Fig. 8, the saliency value of the red *flower* in our map is obviously higher than Li [18]. The similar comparison result is found between ours and Fu [20]. The last two rows of Fig. 8 show that our method is also valid for detecting multiple co-salient objects in a set of images.

3) *Multiple Image Co-Saliency Detection*: Fig. 9 shows sample results obtained by our algorithm and Fu [20] on the iCoseg dataset [40]. Three image groups including *bear*, *cheetah* and *statue* are selected to show the differences between ours and Fu [20]. Each group consists of 10 images. The cluster-based co-saliency model of Fu [20] utilizes three saliency cues (*i.e.*, contrast, spatial, and corresponding) to generate multi-image saliency, and combines the contrast cue and the spatial cue to obtain single-image saliency. To get a more clean result, Fu *et al.* [20] generate the final co-saliency map by multiplying multi-image saliency map with their single-image saliency map. Under the way of multiplication, Fu [20] achieves high precision at the expense of reducing the recall, which results in the incomplete salient objects and the low saliency values. Different from Fu [20], although disturbed by some small non-salient regions, our method segments almost all the salient objects from the backgrounds successfully. Moreover, as argued in Section I, co-saliency has been the preprocessing of many computer vision applications, which prefers the whole salient regions. From this point, our results are more appropriate than Fu [20]. In summary, our results are more intact, conspicuous and effective.

#### D. Discussion

1) *The Degenerated Cases*: Our method is based on fusing multiple saliency maps. Therefore, the major scenario causes our method to fail is when most of the elementary maps lose their power. Fig. 10 provides two failure examples. In the first image pair, background regions are extracted by most elementary maps (*e.g.*, methods (d-h)), which leads our method to extract these non-salient background pixels. Another example is the *gray rabbit* in the second image pair. The *gray rabbit* is not a common salient object. However, our method inaccurately predicts it since five of seven elementary maps wrongly select it. Moreover, as shown in the last row of Fig. 10, our method fails to detect the *white rabbit* when almost all the maps fail.

#### IV. CONCLUSION

In this paper, we provided a general fusion framework for both single saliency and co-saliency detection. The framework inherited the merits of the multiple saliency cues, and was more robust to few wrong detections. This was achieved by utilizing the consistency properties among various saliency detection results. We formalized the consistency property among the salient regions as the rank constraint. The rank-one energy and the low-rank energy were proposed to explain it, where the low-rank is more general than the rank-one. But, they had the similar pleasing performances on the applicable co-saliency detection situations. Different from the most existing methods, we employed our rank constraint to integrate the cues self-adaptively. The experimental results demonstrated the effectiveness and superiority of our proposed algorithm.

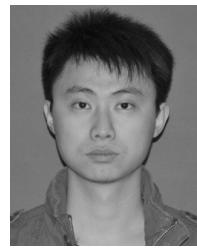
#### REFERENCES

- [1] L. Itti, C. Koch, and E. Niebur, "A model of saliency-based visual attention for rapid scene analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 20, no. 11, pp. 1254–1259, Nov. 1998.
- [2] W. Lee, T. Huang, S. Yeh, and H. Chen, "Learning-based prediction of visual attention for video signals," *IEEE Trans. Image Process.*, vol. 20, no. 11, pp. 3028–3038, Nov. 2011.
- [3] M. Cheng, G. Zhang, N. J. Mitra, X. Huang, and S. Hu, "Global contrast based salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2011, pp. 409–416.
- [4] A. Toet, "Computational versus psychophysical bottom-up image saliency: A comparative evaluation study," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 11, pp. 2131–2146, Nov. 2011.
- [5] R. Valenti, N. Sebe, and T. Gevers, "What are you looking at? Improving visual gaze estimation by saliency," *Int. J. Comput. Vis.*, vol. 98, no. 3, pp. 324–334, 2012.
- [6] U. Rutishauser, D. Walther, C. Koch, and P. Perona, "Is bottom-up attention useful for object recognition?" in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2004, pp. 37–44.
- [7] S. Goferman, L. Zelnik-manor, and A. Tal, "Context-aware saliency detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2010, pp. 2376–2383.
- [8] L. Itti, "Automatic foveation for video compression using a neurobiological model of visual attention," *IEEE Trans. Image Process.*, vol. 13, no. 10, pp. 1304–1318, Oct. 2004.
- [9] H. Liu, X. Xie, X. Tang, Z. Li, and W. Ma, "Effective browsing of web image search results," in *Proc. 6th ACM SIGMM Int. Workshop Multimedia Inform. Retr. (MIR)*, 2004, pp. 84–90.
- [10] T. Liu *et al.*, "Learning to detect a salient object," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 2, pp. 353–367, Feb. 2011.
- [11] R. Achanta, S. Hemami, F. Estrada, and S. Süsstrunk, "Frequency-tuned salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2009, pp. 1597–1604.
- [12] Y. Zhai and M. Shah, "Visual attention detection in video sequences using spatiotemporal cues," in *Proc. 14th Annu. ACM Int. Conf. Multimedia (MULTIMEDIA)*, 2006, pp. 815–824.
- [13] X. Shen and Y. Wu, "A unified approach to salient object detection via low rank matrix recovery," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2012, pp. 853–860.
- [14] Q. Yan, L. Xu, J. Shi, and J. Jia, "Hierarchical saliency detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2013, pp. 1155–1162.
- [15] C. Yang, L. Zhang, H. Lu, X. Ruan, and M.-H. Yang, "Saliency detection via graph-based manifold ranking," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2013, pp. 3166–3173.
- [16] D. Jacobs, D. Goldman, and E. Shechtman, "Cosaliency: Where people look when comparing images," in *Proc. ACM Symp. User Inter. Softw. Technol. (USIT)*, 2010, pp. 219–228.
- [17] H. Chen, "Preattentive co-saliency detection," in *Proc. 17th IEEE Int. Conf. Image Process. (ICIP)*, Sep. 2010, pp. 1117–1120.
- [18] H. Li and K. Ngan, "A co-saliency model of image pairs," *IEEE Trans. Image Process.*, vol. 20, no. 12, pp. 3365–3375, Dec. 2011.

- [19] K. Chang, T. Liu, and S. Lai, "From co-saliency to co-segmentation: An efficient and fully unsupervised energy minimization model," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2011, pp. 2129–2136.
- [20] H. Fu, X. Cao, and Z. Tu, "Cluster-based co-saliency detection," *IEEE Trans. Image Process.*, vol. 22, no. 10, pp. 3766–3778, Oct. 2013.
- [21] D. Hochbaum and V. Singh, "An efficient algorithm for co-segmentation," in *Proc. IEEE 12th Int. Conf. Comput. Vis. (ICCV)*, Oct. 2009, pp. 269–276.
- [22] S. Vicente, C. Rother, and V. Kolmogorov, "Object cosegmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2011, pp. 2217–2224.
- [23] H. Tan and C. Ngo, "Common pattern discovery using earth mover's distance and local flow maximization," in *Proc. IEEE 10th Int. Conf. Comput. Vis. (ICCV)*, Oct. 2005, pp. 1222–1229.
- [24] J. Yuan and Y. Wu, "Spatial random partition for common visual pattern discovery," in *Proc. IEEE 11th Int. Conf. Comput. Vis. (ICCV)*, Oct. 2007, pp. 1–8.
- [25] M. Cho, Y. M. Shin, and K. M. Lee, "Co-recognition of image pairs by data-driven Monte Carlo image exploration," in *Proc. 10th Eur. Conf. Comput. Vis. (ECCV)*, 2008, pp. 144–157.
- [26] A. Toshev, J. Shi, and K. Daniilidis, "Image matching via saliency region correspondences," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2007, pp. 1–8.
- [27] L. Yang, B. Geng, Y. Cai, A. Hanjalic, and X.-S. Hua, "Object retrieval using visual query context," *IEEE Trans. Multimedia*, vol. 13, no. 6, pp. 1295–1307, Dec. 2011.
- [28] X. Hou and L. Zhang, "Saliency detection: A spectral residual approach," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2007, pp. 1–8.
- [29] J. Harel, C. Koch, and P. Perona, "Graph-based visual saliency," in *Advances in Neural Information Processing Systems*. Cambridge, MA, USA: MIT Press, 2007, pp. 545–552.
- [30] T. Judd, K. Ehinger, F. Durand, and A. Torralba, "Learning to predict where humans look," in *Proc. IEEE 12th Int. Conf. Comput. Vis. (ICCV)*, Oct. 2009, pp. 2106–2113.
- [31] A. Borji, D. N. Sihite, and L. Itti, "Salient object detection: A benchmark," in *Proc. 12th Eur. Conf. Comput. Vis. (ECCV)*, 2012, pp. 414–429.
- [32] A. Joulin, F. R. Bach, and J. Ponce, "Discriminative clustering for image co-segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2010, pp. 1943–1950.
- [33] A. Joulin, F. Bach, and J. Ponce, "Multi-class cosegmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2012, pp. 542–549.
- [34] J. Yan, M. Zhu, H. Liu, and Y. Liu, "Visual saliency detection via sparsity pursuit," *IEEE Signal Process. Lett.*, vol. 17, no. 8, pp. 739–742, Aug. 2010.
- [35] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Sü, "SLIC superpixels compared to state-of-the-art superpixel methods," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 11, pp. 2274–2282, Nov. 2012.
- [36] E. J. Candès, X. Li, Y. Ma, and J. Wright, "Robust principal component analysis?" *J. ACM*, vol. 58, no. 3, p. 11, 2011.
- [37] J. Wright, A. Ganesh, S. Rao, Y. Peng, and Y. Ma, "Robust principal component analysis: Exact recovery of corrupted low-rank matrices via convex optimization," in *Advances in Neural Information Processing Systems*. Red Hook, NY, USA: Curran Associates, Inc., 2009, pp. 2080–2088.
- [38] F. Perazzi, P. Krahenbuhl, Y. Pritch, and A. Hornung, "Saliency filters: Contrast based filtering for salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2012, pp. 733–740.
- [39] H. Jiang, J. Wang, Z. Yuan, Y. Wu, N. Zheng, and S. Li, "Salient object detection: A discriminative regional feature integration approach," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2013, pp. 2083–2090.
- [40] D. Batra, A. Kowdle, D. Parikh, J. Luo, and T. Chen, "Interactively co-segmenting topically related images with intelligent scribble guidance," *Int. J. Comput. Vis.*, vol. 93, no. 3, pp. 273–292, 2011.
- [41] T. Fawcett, "Roc graphs: Notes and practical considerations for data mining researchers," HP Laboratories, Palo Alto, CA, USA, Tech. Rep. HPL-2003-4, 2004.
- [42] B. W. Tatler, "The central fixation bias in scene viewing: Selecting an optimal viewing position independently of motor biases and image feature distributions," *J. Vis.*, vol. 7, no. 14, pp. 1–17, 2007.



**Xiaochun Cao** (SM'14) has been a Professor with the Institute of Information Engineering, Chinese Academy of Sciences, Beijing, since 2012. He received the B.E. and M.E. degrees from Beihang University, Beijing, and the Ph.D. degree from the University of Central Florida, Orlando, FL, USA, all in computer science, with his dissertation nominated for the university-level Outstanding Dissertation Award. After graduation, he spent about three years at ObjectVideo, Inc., Reston, VA, USA, as a Research Scientist. From 2008 to 2012, he was a Professor with Tianjin University, Tianjin, China. He has authored and co-authored over 80 journal and conference papers. He was a recipient of the Piero Zamperoni Best Student Paper Award at the International Conference on Pattern Recognition in 2004 and 2010.



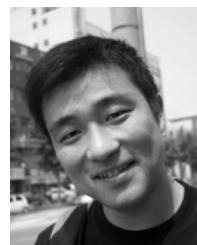
**Zhiqiang Tao** received the B.E. degree in software engineering from the School of Computer Software, Tianjin University, Tianjin, China, in 2012, where he is currently pursuing the M.E. degree with the School of Computer Science and Technology. His current research interests include computer vision, multimedia analysis, and image processing.



**Bao Zhang** received the B.E. degree in software engineering from the School of Computer Software, Tianjin University, Tianjin, China, in 2010, and the M.E. degree from the School of Computer Science and Technology, Tianjin University, in 2013. His current research interests include computer vision, scene classification, video processing, image saliency detection, and segmentation.



**Huazhu Fu** is a Research Fellow with the School of Computer Engineering, Nanyang Technological University, Singapore. He received the B.S. degree from Nankai University, Tianjin, China, in 2006, the M.E. degree from the Tianjin University of Technology, Tianjin, in 2010, and the Ph.D. degree from Tianjin University, Tianjin, in 2013. His current research interests include computer vision, medical image processing, image saliency detection, and segmentation.



**Wei Feng** (M'10) received the B.S. and M.Phil. degrees in computer science from Northwestern Polytechnical University, Xi'an, China, in 2000 and 2003, respectively, and the Ph.D. degree in computer science from the City University of Hong Kong, Hong Kong, in 2008.

He was a Research Fellow with the Chinese University of Hong Kong and then with the City University of Hong Kong from 2008 to 2010. He is currently an Associate Professor with the School of Computer Science and Technology, Tianjin University, Tianjin, China. His major research interest is media computing, in particular, general Markov random fields modeling, discrete/continuous energy minimization, image segmentation, superpixel-level image processing, weakly supervised learning, structural authentication, and generic pattern recognition.

Dr. Feng was a recipient of the Support of the Program Award for New Century Excellent Talents in University, China, in 2011.