

Saliency Detection by Fully Learning a Continuous Conditional Random Field

Keren Fu, Irene Yu-Hua Gu, *Senior Member, IEEE*, and Jie Yang

Abstract—Salient object detection is aimed at detecting and segmenting objects that human eyes are most focused on when viewing a scene. Recently, conditional random field (CRF) is drawn renewed interest, and is exploited in this field. However, when utilizing a CRF with unary and pairwise potentials having essential parameters, most existing methods only employ manually designed parameters, or learn parameters partly for the unary potentials. Observing that the saliency estimation is a continuous labeling issue, this paper proposes a novel data-driven scheme based on a special CRF framework, the so-called continuous CRF (C-CRF), where parameters for both unary and pairwise potentials are jointly learned. The proposed C-CRF learning provides an optimal way to integrate various unary saliency features with pairwise cues to discover salient objects. To the best of our knowledge, the proposed scheme is the first to completely learn a C-CRF for saliency detection. In addition, we propose a novel formulation of pairwise potentials that enables learning weights for different spatial ranges on a superpixel graph. The proposed C-CRF learning-based saliency model is tested on 6 benchmark datasets and compared with 11 existing methods. Our results and comparisons have provided further support to the proposed method in terms of precision-recall and F-measure. Furthermore, incorporating existing saliency models with pairwise cues through the C-CRF are shown to provide marked boosting performance over individual models.

Index Terms—Continuous conditional random field (C-CRF), feature integration, learning, saliency map, salient object detection, spatial ranges.

I. INTRODUCTION

SALIENCY detection is aimed at detecting conspicuous image parts that attract human attention. It simulates and models the selective mechanism of human eyes [1], [2]. There are generally two subcategories of saliency detection: eye-fixation

Manuscript received September 5, 2016; revised January 23, 2017; accepted February 23, 2017. This work was supported in part by the National Natural Science Foundation of China under Grant 61572315 and Grant 6151101179, and in part by the 863 Plan of China under Grant 2015AA042308. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Judith Redi. (*Corresponding authors: Jie Yang and Keren Fu.*)

K. Fu is with the Institute of Image Processing and Pattern Recognition, Shanghai Jiao Tong University, Shanghai 200240, China, and also with the Department of Signals and Systems, Chalmers University of Technology, Gothenburg 41296, Sweden (e-mail: fkrshichaoren@qq.com).

I. Y.-H. Gu is with the Department of Signals and Systems, Chalmers University of Technology, Gothenburg 41296, Sweden (e-mail: irenegu@chalmers.se).

J. Yang is with the Institute of Image Processing and Pattern Recognition, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: jieyang@sjtu.edu.cn).

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/TMM.2017.2679898

prediction [3]–[6] and salient object/region detection [7]–[10]. The former task aims to detect sparse eye-fixation points where human attend in a scene, whereas the latter task is to detect and emphasize entire salient objects from an image, yielding a *saliency map* as output where the pixel-wise intensities indicate the probability of being a salient object. The recent advance in salient object detection is driven by emerging multimedia applications such as automatic object detection and segmentation [11]–[13], content-based image editing [14]–[18], image retrieval [19]–[21] and compression, image sequence and video analysis [22], [23]. In this paper, we mainly address salient object detection.

To emphasize salient objects uniformly, the conditional random field (CRF) that can provide label consistency becomes popular in this field. By utilizing CRF, high quality saliency maps that maintain well-defined object boundaries and uniformly emphasized object interior are achieved. In existing studies [24]–[27], CRFs are employed in explicit or implicit ways. However, when utilizing a CRF whose energy function consists of parameterized unary and pairwise energy potentials, most previous methods use manually designed parameters [25] or learn the parameters only for unary potentials [24], [27]. Hence for saliency detection, the full power of CRF on feature integration is hardly exploited.

Motivated by the above issues, this paper proposes to fully learn a CRF, namely to learn both unary and pairwise parameters in order to exploit the power of CRF for feature integration in saliency detection [1], [3]. More specifically, we investigate a special CRF framework—*continuous CRF (C-CRF)* [28]–[30]. This is motivated by the idea that saliency detection is conventionally treated as a *continuous labeling problem*. In this paper a novel data-driven saliency detection scheme based on C-CRF [28] is proposed, which differs from [24], [27] since ours enables learning to integrate various pairwise features. This allows the C-CRF model to capture more sophisticated interactions between image parts, leading to enhanced delineation between objects and background in the resulting saliency maps. It is worth noting that the work of Mai *et al.* [26] is closely related to ours. In [26], the unary and pairwise potentials of CRF all include parameters. However, the main difference in between is that [26] employs discrete CRF (as will be mentioned later), whereas we propose to use C-CRF which benefits from different designs of energy function, hence very different techniques for learning and inference. In addition, as shown in Section V-C the proposed method improves the performance significantly from [26]. A straightforward comparison of the

TABLE I
COMPARISON OF REPRESENTATIVE CRF-BASED METHODS IN
THE SALIENT OBJECT/REGION DETECTION COMMUNITY

Related work	CRF type	Learning unary terms	Learning pairwise terms
Liu <i>et al.</i> [24]	Discrete (D-CRF)	Yes	No
Mai <i>et al.</i> [26]	Discrete (D-CRF)	Yes	Yes
Yang <i>et al.</i> [25]	Continuous (C-CRF)	No	No
Lu <i>et al.</i> [27]	Continuous (C-CRF)	Yes	No
Ours	Continuous (C-CRF)	Yes	Yes

The abbreviation C-CRF and D-CRF stand for continuous CRF and discrete CRF, respectively.

83 proposed method to state-of-the-art CRF related works are given
84 in Table I. To the best of our knowledge, the complete C-CRF
85 learning and inference theories have not yet been applied to
86 saliency estimation.

87 C-CRF was firstly proposed for ranking documents [28],
88 and later applied to recognition [29] and depth estimation
89 [30]. It is worth noting that CRF has already been applied to
90 figure-ground segmentation [31], semantic segmentation [31]–
91 [33], and also saliency detection [24], [26], [27] (Table I).
92 However, most of them are conventional CRFs with *dis-*
93 *crete* labels. We will later call this type of CRFs as D-CRF
94 (discrete CRF). In the context of saliency detection, C-CRF
95 may suit this problem better since saliency maps are known
96 to be continuous and real-valued [3], [8], [34], revealing
97 saliency detection can be regarded as a continuous labeling
98 problem.

99 The main contributions of this paper are four-fold:

- 100 1) This study is the first to apply the complete C-CRF learn-
101 ing and inferring theories to saliency detection, leading to
102 a data-driven way for saliency feature integration.
- 103 2) As shown in Table I, our work differs from existing
104 saliency models that have explicit/implicit relation to CRF,
105 evolving from partially learning unary terms [24], [27] to
106 jointly learning both unary plus pairwise terms, and from
107 discrete field [26] to continuous field.
- 108 3) We propose a novel formulation of pairwise potentials for
109 C-CRF defined on a superpixel graph. Such a formula-
110 tion is conducted by graph topology decomposition and
111 enables learning pairwise parameters for different spatial
112 ranges of graph connections. This avoids the manual effort
113 of tuning spatial connections of a graph.
- 114 4) We show from tests and comparisons that integrat-
115 ing widely employed unary saliency features with pair-
116 wise cues in a C-CRF manner outperforms a range of
117 state-of-the-art methods. Furthermore, integrating sev-
118 eral best-performing state-of-the-art methods through
119 a C-CRF further pushes the performance to a new
120 high level.

121 The reminder of this paper is organized as follows. Section II
122 briefly reviews the fundamental theories of CRF and C-CRF.
123 Section III describes the related work. Section IV describes the
124 proposed method. Experimental results, performance evaluation
125 and comparisons are included in Section V. Finally, conclusion
126 is drawn in Section VI.

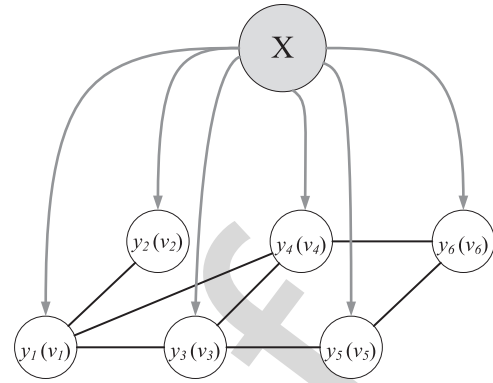


Fig. 1. General graphic model of CRF for image labeling task. A white vertex (v_i) represents a label (y_i) and the gray vertex (\mathbf{x}) represents the entire image. The gray arrows indicate the unary dependencies (conditions) while the black lines indicate the pairwise relations associating with a graph.

II. CONDITIONAL RANDOM FIELD (CRF) AND CONTINUOUS CONDITIONAL RANDOM FIELD (C-CRF): A BRIEF REVIEW

A. Probabilistic Formulation

130 Conditional random field (CRF) is originally proposed by
131 Lafferty *et al.* [35] for labeling sequence data. For the image
132 labeling task, given an image \mathbf{x} , the conditional probability
133 distribution of a label configuration \mathbf{y} (in vector form) on the
134 CRF can be defined as

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp\{-\mathbb{E}(\mathbf{y}, \mathbf{x})\} \quad (1)$$

135 where $\mathbb{E}(\mathbf{y}, \mathbf{x})$ is the energy function and $Z(\mathbf{x})$ is the partition
136 function.¹ The energy function can be expressed as unary terms
137 plus pairwise terms as

$$\mathbb{E}(\mathbf{y}, \mathbf{x}) = \sum_i \underbrace{U_\alpha(y_i, \mathbf{x})}_{\text{Unary term}} + \sum_{i,j,i \sim j} \underbrace{P_\varphi(y_i, y_j, \mathbf{x})}_{\text{Pairwise term}} \quad (2)$$

138 where y_i is the i th element of the label vector \mathbf{y} , U_α and P_φ
139 denote the unary and pairwise terms parameterized by vector
140 α and φ (vector α contains the parameters for unary poten-
141 tials, and vector φ contains the parameters for pairwise poten-
142 tials). A CRF is often coupled with the definition of an
143 undirected graph $G(V, E)$ [35], where V is the set of graph
144 nodes and E is the set of graph edges. The label assigned to
145 each graph node $v_i \in V$ is denoted as y_i . In (2), the notation
146 “ $i \sim j$ ” means that v_i and v_j are graph neighbors. Theoretically,
147 the unary term U_α represents the dependency between a
148 label and the image \mathbf{x} at a specific node, whereas the pairwise
149 term P_φ encourages neighboring graph nodes to take similar
150 labels (i.e., enforces labeling consistency). A general graphic
151 model of CRF for image labeling task is given in Fig. 1, where a
152 white vertex represents a label and the gray vertex represents the
153 entire image.

¹The partition functions for D-CRF and C-CRF are defined as $Z(\mathbf{x}) = \sum_{\mathbf{y}} \exp\{-\mathbb{E}(\mathbf{y}, \mathbf{x})\}$ and $Z(\mathbf{x}) = \int_{\mathbf{y}} \exp\{-\mathbb{E}(\mathbf{y}, \mathbf{x})\} d\mathbf{y}$, respectively.

154 B. D-CRF and C-CRF

155 In the conventional CRF, i.e., the D-CRF [31]–[33], [35], all
 156 components of \mathbf{y} range over a finite label alphabet (e.g., subject
 157 to $\mathbf{y} \in \{0, 1\}^n$ for a binary labeling problem, where n is the
 158 dimension of \mathbf{y}), whereas in the continuous CRF (C-CRF) [28],
 159 \mathbf{y} is relaxed to be continuous values ($\mathbf{y} \in \mathcal{R}^n$). Due to such re-
 160 laxing, the designs of energy functions for D-CRF and C-CRF
 161 differ. For example, D-CRF usually employs Potts model [32],
 162 [33] with indicator function for the pairwise terms, whereas in
 163 C-CRF, quadratic cost function can be used to measure the label
 164 compatibility. Besides, the techniques for learning and inference
 165 of C-CRF [28] differ significantly from D-CRF. The exact learn-
 166 ing/inference of D-CRF is usually intractable due to its discrete
 167 property, which requires approximation techniques [36] such
 168 as belief propagation, mean field, Monte Carlo approaches, to
 169 name a few. In contrast, C-CRF offers direct learning together
 170 with closed-form inference, which will be shown later in this
 171 paper.

172 Assuming that the parameters of a CRF are given or esti-
 173 mated by learning, theoretically the optimal labeling vector \mathbf{y}
 174 can be inferred by maximizing (1), or equivalently minimizing
 175 the negative logarithm of (1) as

$$-\log p(\mathbf{y}|\mathbf{x}) = \mathbb{E}(\mathbf{y}, \mathbf{x}) + \log Z(\mathbf{x}). \quad (3)$$

176 Since $\log Z(\mathbf{x})$ is a constant with respect to \mathbf{y} , one can directly
 177 minimize the energy function $\mathbb{E}(\mathbf{y}, \mathbf{x})$. From this viewpoint, ex-
 178 isting methods on saliency detection such as manifold ranking
 179 [25], graph regularization [37], quadratic model [27] that min-
 180 imize an energy function in the form of (2) can be viewed as
 181 special cases of inferring C-CRFs.

182 III. RELATED WORK

183 A large number of literatures on salient object detection ex-
 184 ist, see the comprehensive survey [38] and benchmarking [10].
 185 Here we review some previous works that are highly relevant to
 186 data-driven approaches or CRF-based approaches.

187 *Data driven approaches:* The concept of *learning to detect* in
 188 saliency detection originates from [24], [39]. The idea behind is
 189 to automatically discover feature integration rules from training
 190 data instead of using manually designed rules. Judd *et al.* [39]
 191 propose to learn a saliency model from eye-tracking data, where
 192 low-, middle- and high-level image features are integrated by
 193 a linear SVM. Their work is, however, focused on eye-fixation
 194 prediction. Alex *et al.* [40] learn to score windows sampled from
 195 a given image, where the Bayesian theory is adopted for cue in-
 196 tegration. The posterior constitutes the final objectness score of
 197 a window. Khuwuthyakorn *et al.* [41] learn to integrate pixel-
 198 wise saliency features via a mixture of linear SVMs. Mehrani
 199 *et al.* [42] use confidence scores from a boosting classifier to
 200 formulate a saliency map. After that, the saliency map is fed to a
 201 graph cut program for figure-ground segmentation. Jiang *et al.*
 202 [43] propose to extract abundant discriminative features from
 203 image regions. A random forest regressor is trained to map re-
 204 gional features to final saliency scores. Online saliency learning
 205 is proposed in [44], [45], where multiple kernel boosting is em-
 206 ployed to identify salient parts against non-salient parts. Some

recent data-driven methods [46], [47] consider deep learning for
 saliency detection. Due to the deep architecture of convolutional
 neural networks (CNNs), impressive performance is obtained.
 However, in CNNs there often lacks explicit modeling of neigh-
 borhood relations. Therefore, post-processing like C-CRF may
 be required.

CRF inference-based approaches: Several methods [25], [37]
 are based on inferring C-CRF without learning, where fea-
 tures and integration rules are manually specified. In [25], [37],
 though the word “CRF” or “continuous CRF” is not explicitly
 mentioned, there is a potential connection between these meth-
 ods and C-CRF. To be more specific, the employed manifold
 ranking [25] and graph regularization [37] are special cases of
 inferring C-CRFs, as aforementioned in Section II-B.

CRF learning-based approaches: Some methods on saliency
 detection are based on both learning and inferring D-CRFs
 or C-CRFs. Learning is first conducted to obtain optimal pa-
 rameters and inference is then applied on user-input images to
 achieve final saliency maps. Representative works include: Liu
et al. [24] detect and segment salient objects by aggregating
 pixel saliency cues in a D-CRF. Linear weights for those cues
 are learned under the maximized likelihood (ML) criteria by
 tree-reweighted belief propagation. Mai *et al.* [26] propose a
 saliency aggregation approach, which aggregates saliency maps
 output by existing saliency models using a D-CRF. Weights
 for aggregation are learned from images retrieved from a pre-
 defined dataset. Lu *et al.* [27] learn optimal combination of
 seeds for graph-based diffusion by maximizing figure-ground
 segregation, where the employed graph diffusion is tightly re-
 lated to C-CRF. The method boils down to learning the linear
 parameters of unary terms of the C-CRF. In summary, [24], [26]
 concern D-CRFs for saliency detection, where only unary pa-
 rameters are learned. [27] implicitly considers a C-CRF, where
 again only the unary terms are learned. In contrast, our data-
 driven scheme differs from all the above methods on learning a
 complete C-CRF.

IV. THE PROPOSED METHOD

This section describes the proposed method for saliency de-
 tection that is based on fully learning and inferring a continuous
 CRF (C-CRF). The block diagram of the proposed method is
 given in Fig. 2. An input image is first over-segmented into su-
 perpixels and a superpixel graph is established to capture intrin-
 sic image context. A C-CRF will later be defined in conjunction
 with this graph. Next, we extract various unary saliency features
 and pairwise cues, which will be used to compose the unary and
 pairwise terms in the C-CRF energy function. By utilizing the
 off-line learned C-CRF parameters for both unary and pairwise
 potentials, the inference of the C-CRF corresponds to a final
 saliency map that is continuously valued. Details of each part of
 the method are further given in the following subsections.

A. Graph Construction From an Image

We first describe the graph construction, where the C-CRF is
 defined upon. Rather than constructing CRF on the pixel level
 [24], the proposed C-CRF is constructed on superpixels, where

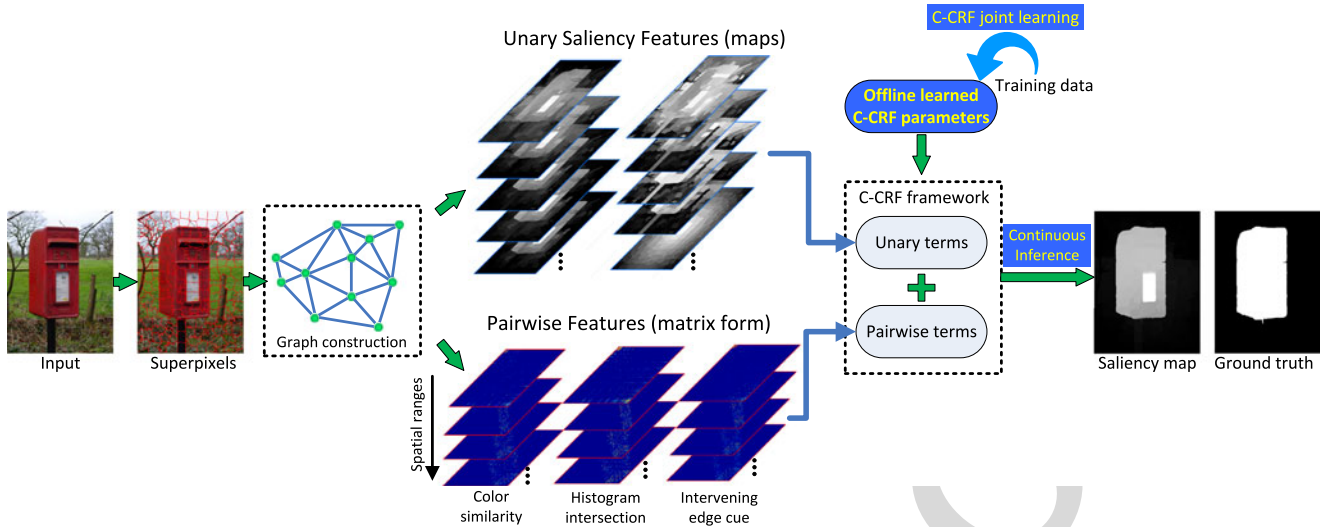


Fig. 2. Block diagram of the proposed salient object detection method.

only a small number of graph nodes are needed. An input image \mathbf{x} is first over-segmented into n superpixels by using the SLIC algorithm [48], which is very widely employed by previous work [25], [27], [37], [49]–[51] as a pre-processing step. Then superpixels are used as processing units. A graph $G = (V, E)$ is then constructed, where the node set V consists of superpixels. In this paper, the terms of “superpixels” and “graph nodes” are interchangeable, and $v_i, i \in \{1 : n\}$ indicates the i th superpixel/node. To build the connections of graph edges, we first construct an initial adjacency graph $G^0 = (V, E^0)$, where vertices V correspond to superpixels. E^0 is the edge set (weighed by value 1.0) formed between pairs of spatially adjacent superpixels. Let $D^0(v_i, v_j)$ be the length of the shortest path on G^0 between nodes v_i and v_j . Then, the edge set E is formed between pairs of superpixels that are less than T nodes away on G^0 , namely

$$e_{ij} \in E, \text{ if } D^0(v_i, v_j) \leq T \quad (4)$$

where T ($T \geq 1$) is a predefined integer that specifies the *maximum spatial range*. Further, as observed in many images that boundary superpixels are likely the same semantic background and also inspired by previous work [25], [50]–[52], we establish connection between arbitrary boundary superpixels as below

$$e_{ij} \in E, \text{ if } v_i, v_j \in \mathbb{B} \quad (5)$$

where \mathbb{B} is a set containing all boundary superpixels. Fig. 3 shows an example of the graph connections for the case of $T = 3$. By this mean, boundary superpixels are able to serve as “bridges” for labeling consistency in image background.

B. C-CRF Composition

The definitions of unary term U_α and pairwise term P_φ in our method are motivated by the work of Qin *et al.* [28]. The basic idea is that although a unary term calculates the dependency between a node label y_i and the entire image \mathbf{x} (Fig. 1), the case can be simplified by considering the dependency between y_i and a corresponding feature vector \mathbf{f}_i that derives from \mathbf{x} . In our

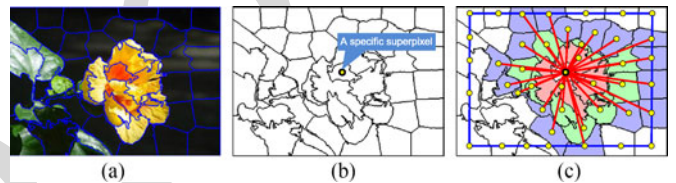


Fig. 3. Superpixels and graph construction. (a) An input image with superpixel boundaries overlapped in blue. About 50 superpixels are generated just for better illustration. (b) Superpixel boundaries in black are shown, where a superpixel is specified. (c) C-CRF graph connections, where the connections (red lines) from the specific node are shown. The maximum spatial range $T = 3$ is set as example. Superpixels filled in red/green/blue mean that they are 1/2/3 nodes away from the specific superpixel, respectively. Besides, blue lines around image boundary means that two arbitrary boundary nodes are connected, as expressed in (5).

case, \mathbf{f}_i is a feature vector that captures the saliency information in \mathbf{x} (see Section IV-C).

The unary term: Assuming a d -dimensional unary saliency feature vector \mathbf{f}_i for node v_i , the unary term is defined as a weighted sum of quadratic cost

$$U_\alpha(y_i, \mathbf{x}) = \sum_{k=1}^d \alpha_k (y_i - f_{i,k})^2 \quad (6)$$

where $\alpha_k, f_{i,k}$ are the k th components of α and \mathbf{f}_i respectively, and α_k indicates the weight of the k th component in the feature vector. The overall cost becomes larger if the label y_i deviates from the correspondent feature components with high weights. Further, in (6), \mathbf{x} is omitted for simplicity, though the unary feature vector \mathbf{f}_i is dependent on \mathbf{x} .

The pairwise term: Likewise, the pairwise term is a weighted sum of quadratic cost defined as

$$P_\varphi(y_i, y_j, \mathbf{x}) = \frac{1}{2} \sum_{k=1}^h \varphi_k S_{ij}^k (y_i - y_j)^2 \quad (7)$$

where S_{ij}^k is the k th pairwise feature defined between nodes v_i and v_j , φ_k is the k th component of φ , and h is the number of pairwise features. In the proposed method, S_{ij}^k is a positive

309 affinity (similarity) function between v_i and v_j , and it is large if
 310 v_i and v_j are similar, so that they can be assigned with similar
 311 labels by C-CRF. Similar to (6), we have omitted \mathbf{x} for simplicity
 312 in (7), although S_{ij}^k depends on \mathbf{x} as well.

313 *The energy function:* According to (2), the energy function
 314 $\mathbb{E}(\mathbf{y}, \mathbf{x})$ has the following form:

$$\mathbb{E}(\mathbf{y}, \mathbf{x}) = \sum_{i=1}^n \sum_{k=1}^d \alpha_k (y_i - f_{i,k})^2 + \sum_{i,j,i \sim j} \frac{1}{2} \sum_{k=1}^h \varphi_k S_{ij}^k (y_i - y_j)^2 \quad (8)$$

315 where $\alpha_k > 0$ and $\varphi_k \geq 0$ are needed to ensure the partition
 316 function $Z(\mathbf{x})$ analytically computable (will be clear later).

317 Let \mathbf{F} denote the stacked *feature matrix* whose row is \mathbf{f}_i^T , and
 318 let \mathbf{S}^k denote the matrix whose entry is S_{ij}^k . With some math-
 319 ematic derivation, the matrix form of (8) can be equivalently
 320 expressed as

$$\begin{aligned} \mathbb{E}(\mathbf{y}, \mathbf{x}) &= \mathbf{e}^T \boldsymbol{\alpha} \mathbf{y}^T \mathbf{I} \mathbf{y} - 2\mathbf{y}^T \mathbf{F} \boldsymbol{\alpha} + \text{Tr}\{\mathbf{F} \text{diag}(\boldsymbol{\alpha}) \mathbf{F}^T\} \\ &\quad + \sum_{k=1}^h \varphi_k \mathbf{y}^T \mathbf{L}^k \mathbf{y} \end{aligned} \quad (9)$$

321 where \mathbf{e} is an all-one vector, \mathbf{I} is an identity matrix, $\text{Tr}\{\cdot\}$ is the
 322 trace, $\text{diag}(\boldsymbol{\alpha})$ is the diagonal matrix with $\boldsymbol{\alpha}$ in the diagonal, and
 323 \mathbf{L}^k is the Laplacian matrix of \mathbf{S}^k . The definition of Laplacian
 324 matrix is $\mathbf{L}^k := \mathbf{D}^k - \mathbf{S}^k$, where \mathbf{D}^k is the degree matrix whose
 325 i th diagonal entry is $D_{ii}^k = \sum_j S_{ij}^k$.

326 *The partition function:* The partition function $Z(\mathbf{x})$ in the
 327 proposed scheme is integrable due to the continuous property
 328 of C-CRF. Firstly we introduce the below notation \mathbf{A} , \mathbf{b} , c :

$$\mathbf{A} = \mathbf{e}^T \boldsymbol{\alpha} \mathbf{I} + \sum_{k=1}^h \varphi_k \mathbf{L}^k, \quad \mathbf{b} = \mathbf{F} \boldsymbol{\alpha}, \quad c = \text{Tr}\{\mathbf{F} \text{diag}(\boldsymbol{\alpha}) \mathbf{F}^T\}. \quad (10)$$

329 Then according to the Gaussian integration [53], we have

$$\begin{aligned} Z(\mathbf{x}) &= \int \exp(-\mathbb{E}(\mathbf{y}, \mathbf{x})) d\mathbf{y} \\ &= \exp(-c) \int \exp(-\mathbf{y}^T \mathbf{A} \mathbf{y} + 2\mathbf{y}^T \mathbf{b}) d\mathbf{y} \\ &= \exp(-c) \int \exp\left(-\frac{1}{2} \mathbf{y}^T (2\mathbf{A}) \mathbf{y} + (2\mathbf{b})^T \mathbf{y}\right) d\mathbf{y} \\ &= \frac{\pi^{\frac{n}{2}}}{|\mathbf{A}|^{\frac{1}{2}}} \exp(\mathbf{b}^T \mathbf{A}^{-1} \mathbf{b} - c) \end{aligned} \quad (11)$$

330 where n equals to the dimension of \mathbf{A} , and $|\mathbf{A}|$ is the determi-
 331 nant. The invertibility of \mathbf{A} is guaranteed, as $\alpha_k > 0$, $\varphi_k \geq 0$,
 332 and \mathbf{L}^k is positive semi-definite.

333 *The negative log-likelihood:* Substitute (11) and (9) into (3)
 334 meanwhile notice the notations in (10), the negative log likeli-
 335 hood in (3) can be re-written as

$$\begin{aligned} &-\log p(\mathbf{y}|\mathbf{x}) \\ &= \mathbf{y}^T \mathbf{A} \mathbf{y} - 2\mathbf{y}^T \mathbf{b} + \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b} + \frac{n}{2} \log \pi - \frac{1}{2} \log |\mathbf{A}|. \end{aligned} \quad (12)$$

TABLE II
COLUMNS OF THE UNARY SALIENCY FEATURE MATRIX \mathbf{F}

Column	Categorization	Description
$\mathbf{F}_{:,1 \sim 4}$	Connectivity-based	Geodesic distance to each side of image borders
$\mathbf{F}_{:,5}$	Connectivity-based	Minimum geodesic distance to four image borders
$\mathbf{F}_{:,6}$	Connectivity-based	Normalized soft region area subtracted by 1
$\mathbf{F}_{:,7}$	Contrast-based	Spatially weighted color contrast to other superpixels
$\mathbf{F}_{:,8}$	Contrast-based	Color contrast to all boundary superpixels (backgroundness)
$\mathbf{F}_{:,9}$	Distribution heuristic	Normalized color spatial variances subtracted by 1
$\mathbf{F}_{:,10}$	Distribution heuristic	Image center bias map
$\mathbf{F}_{:,11}$	Clarity-based	Normalized singular value feature subtracted by 1

TABLE III
PAIRWISE FEATURES (IN MATRIX FORM) BETWEEN SUPERPIXELS FROM EDGE SETS E_B AND $E_x |_{x \in \{1:T\}}$

Notation	Categorization	Description
$\mathbf{S}^1, \mathbf{S}^{2 \sim T+1}$	Color-based	Color similarity ($S_{ij}^{(c)}$) from E_B and $E_x _{x \in \{1:T\}}$
$\mathbf{S}^{T+2}, \mathbf{S}^{T+3 \sim 2T+2}$	Color-based	Histogram intersection ($S_{ij}^{(h)}$) from E_B and $E_x _{x \in \{1:T\}}$
$\mathbf{S}^{2T+3}, \mathbf{S}^{2T+4 \sim 3T+3}$	Edge-based	Intervening edge cue ($S_{ij}^{(e)}$) from E_B and $E_x _{x \in \{1:T\}}$

C. Definition of Unary and Pairwise Features

This subsection describes the unary saliency features (\mathbf{f}_i) and the pairwise features (S_{ij}^k) in the proposed C-CRF model. The proposed formulation of pairwise potentials enables learning importance for different spatial ranges of graph connections. All features used are summarized in Tables II and III. Details are given below:

1) *Unary Saliency Features:* Unary saliency feature vector $\mathbf{f}_i \in \mathbb{R}^d$ is initial description for the saliency level of v_i . Each component of \mathbf{f}_i is a type of pre-computed saliency feature, where regions correspond to larger components are more salient. Recall that $\mathbf{F} \in \mathbb{R}^{n \times d}$ is the feature matrix whose row is \mathbf{f}_i^T . Thereby a certain column of \mathbf{F} can be regarded as a type of *feature map*, denoted as $\mathbf{F}_{:,k}$, $k \in \{1:d\}$ (Fig. 4). The unary saliency features considered in this paper fall into four types: connectivity-based, contrast-based, distribution heuristics, and clarity-based features, as given in Table II.

Connectivity-based features: Connected regions tend to be perceived as one entity by human eyes, and regions that easily connect to the image boundary are likely to be the background [49]. The boundary connectivity is hence defined based on the geodesic distance [49]. Computing geodesic distance between superpixels and four image borders separately leads to four feature maps, denoted as $\mathbf{F}_{:,1 \sim 4}$. The minimum geodesic distance between superpixels and image boundary leads to a single feature map $\mathbf{F}_{:,5}$. Since salient objects usually occupy small regions as comparing to large areas of background, we compute

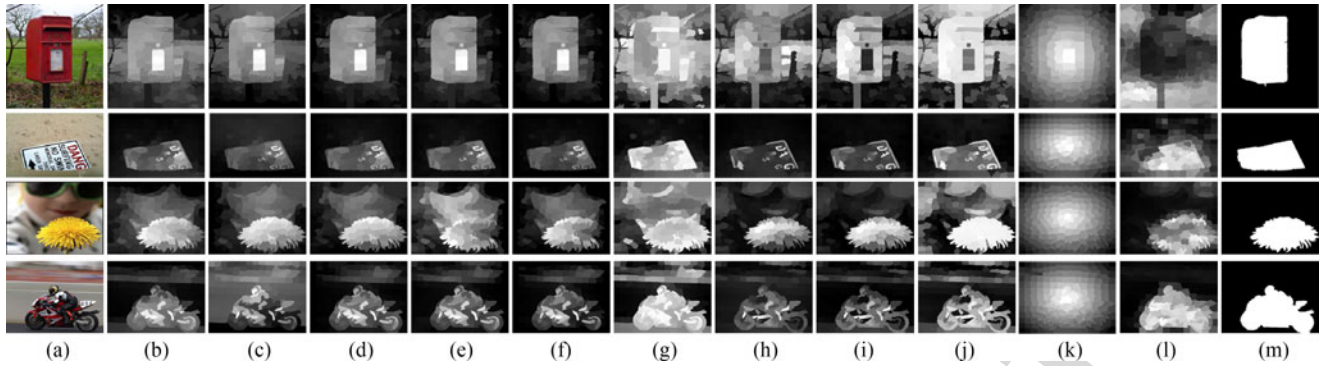


Fig. 4. Feature maps of unary saliency features. From left to right: (a) Input images. (b)–(g) Features $\mathbf{F}_{:,1\sim6}$ (connectivity-based features). (h)–(i) Features $\mathbf{F}_{:,7\sim8}$ (contrast-based features). (j)–(k) Features $\mathbf{F}_{:,9\sim10}$ (distribution-based features). (l) Feature $\mathbf{F}_{:,11}$ (clarity-based feature). (m) Ground truth masks.

363 a type of feature which takes the region size into account. Let
 364 $d_{geo}(v_i, v_j)$ be the geodesic distance between superpixel v_i and
 365 v_j , the geodesic affinity [50] between v_i and v_j can be defined
 366 as $\mathcal{A}_{ij} = \exp(-\frac{d_{geo}^2(v_i, v_j)}{2\sigma_g^2})$ (σ_g is set according to [50]). Then
 367 we compute the spanning area associated with v_i as $\sum_{j=1}^n \mathcal{A}_{ij}$.
 368 This definition of the region area avoids an explicit hard seg-
 369 mentation of image and is “soft”. Meanwhile it takes advantage
 370 of superpixels. Finally, $\mathbf{F}_{:,6}$ is formed by normalizing the span-
 371 ning area value within the range $[0, 1.0]$ and then subtracting the
 372 result from 1.0. This fits the intuition that small object regions
 373 tend to be salient (Fig. 4(g)).

374 *Contrast-based features:* Global color contrast is an indicator
 375 for saliency [8], [54]. $\mathbf{F}_{:,7}$ is computed similarly to [54] by
 376 comparing the color contrast of a superpixel to other superpixels,
 377 where spatially nearer superpixels are rendered larger weights.
 378 Furthermore, we formulate $\mathbf{F}_{:,8}$ by computing the contrast of a
 379 superpixel v_i to all the boundary superpixels as $\sum_{v_j \in \mathbb{B}} \|c_i -$
 380 $c_j\|_2$, where c_i and c_j are the average colors of superpixel v_i and
 381 v_j , since boundary superpixels are likely to be the background
 382 [43], [49].

383 *Distribution heuristics:* Salient objects tend to present com-
 384 pact color distribution [54], [55]. Taking into consideration of
 385 this, we compute a color distribution map [54], where spatial
 386 variances of colors are normalized and subtracted by 1.0 to form
 387 $\mathbf{F}_{:,9}$. Furthermore, to describe the center-bias in human atten-
 388 tion [39], [56], $\mathbf{F}_{:,10}$ in our case is a parameter-free center-bias
 389 map computed by

$$f_{i,10} = 1 - \frac{\|\mathbf{p}_i - \mathbf{p}_c\|_2}{\sqrt{(l_h/2)^2 + (l_w/2)^2}} \quad (13)$$

390 where l_h, l_w are the height and width of the image, $\mathbf{p}_i, \mathbf{p}_c$ are
 391 the spatial coordinates of v_i and image center, respectively.

392 *Clarity-based feature:* Photographers tend to put objects of
 393 interest in focus meanwhile defocus irrelevant background when
 394 making high quality photos. To characterize this, we consider
 395 the Singular Value Feature (SVF) [57], [58] that models the
 396 degree of blur. An input image is first split into $l \times l$ number
 397 of grids, and the SVF [57] is then computed from each grid
 398 and further assigned to the pixels in the grid. A superpixel-
 399 based map is obtained by averaging SVF of pixels in every

superpixel. After normalizing all SVF values into $[0, 1.0]$, they
 are subtracted by 1 to achieve the final feature map $\mathbf{F}_{:,11}$. This
 describes focused objects as more salient. It is worth noting
 that $l \in \{10, 20, 30\}$ are used to consider different scales, and
 feature maps are averaged to form $\mathbf{F}_{:,11}$ [Fig. 4(l)].

405 *Remarks:* In total, unary saliency feature vector \mathbf{f}_i has 11
 406 components (i.e., $d = 11$), with each dimension normalized into
 407 $[0, 1.0]$. Fig. 4 shows examples of all feature maps visually.
 408 Noting that some of the features mentioned above are employed
 409 in existing work, however with different application context. We
 410 reformulate and modify the above features to constitute a unary
 411 feature matrix \mathbf{F} for our C-CRF model. It is worthy noting our
 412 model is generic and not limited to the above features. If needed,
 413 more features can be easily integrated in such a way.

414 2) *Pairwise Features:* As summarized in Table III, we con-
 415 sider color-based and edge-based pairwise features to capture
 416 the interaction between superpixels.

417 *Color-based features:* For color-based pairwise features,
 418 we consider the average-color similarity $S_{ij}^{(c)} = e^{-\lambda_c \|c_i - c_j\|_2}$
 419 and the histogram intersection $S_{ij}^{(h)} = \sum_{k=1} \min\{\tilde{h}_{i,k}, \tilde{h}_{j,k}\}$,
 420 where c_i and c_j are the average colors of v_i and v_j , and \tilde{h}_i, \tilde{h}_j
 421 are the normalized color histograms from v_i and v_j . We obtain
 422 quantized color histograms similarly to Cheng’s work [59] by
 423 first dividing the color space into $8^3 = 512$ bins. Color bins
 424 that are occupied by 99% of image pixels are kept, whereas
 425 pixels with discarded colors are then replaced by their nearest
 426 colors. This reduces the dimension of histograms and makes the
 427 computation more efficient.

428 *Edge-based features:* The edge-based feature is defined as
 429 $S_{ij}^{(e)} = e^{-\lambda_e \max_{i' \in \tilde{i}\tilde{j}} \|f_{i'}\|}$, where $\tilde{i}\tilde{j}$ is a straight line connecting
 430 centers of v_i and v_j on the image plane, i' is a pixel on $\tilde{i}\tilde{j}$,
 431 and $\|f_{i'}\|$ is the edge magnitude at i' that can be derived from
 432 some edge detector. The rationale behind this feature is that
 433 $S_{ij}^{(e)}$ becomes small when there exists strong intervening edges
 434 between two superpixels, meaning v_i and v_j are less likely to
 435 have similar labels. We adopt the structured random forest-based
 436 edge detector proposed in [60] as it produces multi-scale edges
 437 with fast speed. It is worthy noting that in practice, although most
 438 superpixels in an image have their centroids inside due to the
 439 spatial compactness of SLIC superpixels [48], there may be few

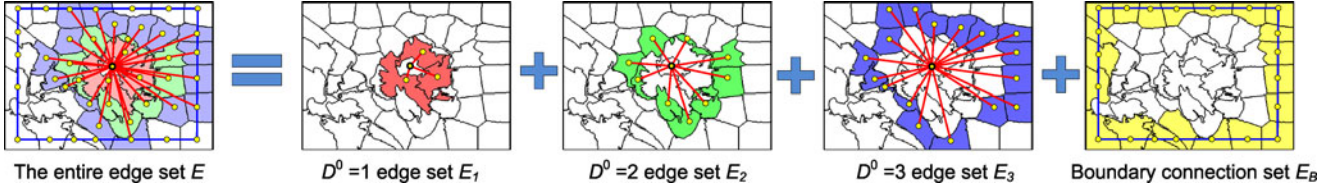


Fig. 5. Decomposition of graph edge connections in Fig. 3. The graph topology is decomposed into a boundary set E_B and $D^0 = 1, 2, 3$ sets E_1, E_2, E_3 , which indicates that superpixels which are exactly one, two, and three nodes away are connected.

440 superpixels whose centroids are located out of them, resulting
 441 in less accurate intervening edge cue. Hence for such a case,
 442 a pixel location is sampled randomly inside each superpixel,
 443 and is used instead of the centroid to compute the edge-based
 444 features.

445 *Graph connection decomposition:* To enable learning dif-
 446 ferent importance of different spatial ranges, the initial
 447 graph topology of G is partitioned into $(T + 1)$ edge sets
 448 $E_B, E_1, E_2, \dots, E_T$, where E_B contains only boundary con-
 449 nections, whereas $E_{x|x \in \{1:T\}}$ contains connections between super-
 450 pixels that are exactly x nodes away. Such a graph decom-
 451 position is designed for properly representing different spatial
 452 ranges meanwhile avoiding an individual edge being counted
 453 multiple times. An example of such topology decomposition is
 454 shown in Fig. 5, where $T = 3$. For a specific type of connec-
 455 tions $E_{x|x \in \{B, 1, 2, \dots, T\}}$, three aforementioned pairwise features
 456 are calculated, leading to $3 \times (T + 1)$ pairwise potentials (see
 457 Table III). For example, when specifying the maximum range
 458 $T = 3$ (Figs. 3 and 5), it typically results in 12 pairwise fea-
 459 tures corresponding to 12 matrices ($\mathbf{S}^{1 \sim 12}$). The proposed graph
 460 decomposition enables C-CRF to automatically learn different
 461 weights for different ranges of connections. $\varphi_k \rightarrow 0$ is equiva-
 462 lent to discarding a type of connections if their contribution is
 463 very little during learning.

464 *Remarks:* Although some of the pairwise information above
 465 is employed by existing saliency work to build graph weights,
 466 they are usually used in an unsupervised fashion. In contrast, we
 467 combine the above features in a supervised way through learning
 468 a complete C-CRF. Besides, the advantage of our formulation of
 469 pairwise potentials is that it avoids the manual effort of tuning
 470 spatial connections. It has been observed in recent work [25],
 471 [27], [61], [62] that the ranges of spatial connections impact
 472 the final detection performance. Most of those models typically
 473 adopt non-local graph connections which are manually deter-
 474 mined. Choosing appropriate graph connections, however, is
 475 a non-trivial task and the optimal connection ranges can de-
 476 pend on the coarseness of superpixels in the image. By contrast,
 477 our technique enables one to specify a relatively large maximum
 478 range T and then automatically learn the corresponding weights
 479 of connections within T . By checking the weights, one can fur-
 480 ther decide whether extension or pruning of spatial ranges is
 481 needed.

482 D. C-CRF Learning and Inference

483 We formulate the C-CRF learning as follows: given N
 484 training images $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N$ with their ground truth labels

$\mathbf{y}^1, \mathbf{y}^2, \dots, \mathbf{y}^N$, learn C-CRF parameters α and φ . The reg- 485
 ularized maximum conditional likelihood (RMCL) training is 486
 adopted for C-CRF learning, which is equivalent to minimizing 487
 (3) summed over all training images 488

$$\min_{\alpha, \varphi} \sum_{i=1}^N \left\{ -\log p(\mathbf{y}^i | \mathbf{x}^i) + \frac{\lambda_1}{2} \|\alpha\|_2^2 + \frac{\lambda_2}{2} \|\varphi\|_2^2 \right\} \quad (14)$$

s.t. $\alpha_k > 0, \varphi_k \geq 0$

where λ_1 and λ_2 are regularization parameters (pre-tuned). The 489
 optimal solution can be found by using gradient descent [28], 490
 [29]. Due to the constraints $\alpha_k > 0$ and $\varphi_k \geq 0$, we apply gradi- 491
 ent descent iteratively on $\log \alpha_k$ and $\log \varphi_k$ during the optimiza- 492
 tion. Let the gradient of the energy loss in (14) w.r.t. $\log \alpha_k$ and 493
 $\log \varphi_k$ be $\nabla_{\log \alpha_k}$ and $\nabla_{\log \varphi_k}$, respectively. Here by dropping 494
 the summation operation for notation simplicity, the derivation 495
 of $\nabla_{\log \alpha_k}$ and $\nabla_{\log \varphi_k}$ is written as 496

$$\nabla_{\log \alpha_k} = \alpha_k \left\{ \sum_i (y_i - f_{i,k})^2 + \frac{\partial \log Z(\mathbf{x})}{\partial \alpha_k} + \lambda_1 \alpha_k \right\} \quad (15)$$

$$\nabla_{\log \varphi_k} = \varphi_k \left\{ \mathbf{y}^T \mathbf{L}^k \mathbf{y} + \frac{\partial \log Z(\mathbf{x})}{\partial \varphi_k} + \lambda_2 \varphi_k \right\} \quad (16)$$

where further according to (11) and use the notations in (10), 497
 $\frac{\partial \log Z(\mathbf{x})}{\partial \alpha_k}$ can be computed 498

$$\frac{\partial \log Z(\mathbf{x})}{\partial \alpha_k} = -\frac{1}{2|\mathbf{A}|} \frac{\partial |\mathbf{A}|}{\partial \alpha_k} + \frac{\partial \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}}{\partial \alpha_k} - \frac{\partial c}{\partial \alpha_k} \quad (17)$$

$$\frac{\partial |\mathbf{A}|}{\partial \alpha_k} = |\mathbf{A}| \text{Tr}(\mathbf{A}^{-1}) \quad (18)$$

$$\frac{\partial \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}}{\partial \alpha_k} = \mathbf{F}_{:,k}^T \mathbf{A}^{-1} \mathbf{b} - \mathbf{b}^T \mathbf{A}^{-1} \mathbf{A}^{-1} \mathbf{b} + \mathbf{b}^T \mathbf{A}^{-1} \mathbf{F}_{:,k} \quad (19)$$

$$\frac{\partial c}{\partial \alpha_k} = \|\mathbf{F}_{:,k}\|_2^2 \quad (20)$$

and $\frac{\partial \log Z(\mathbf{x})}{\partial \varphi_k}$ can be computed 499

$$\frac{\partial \log Z(\mathbf{x})}{\partial \varphi_k} = -\frac{1}{2|\mathbf{A}|} \frac{\partial |\mathbf{A}|}{\partial \varphi_k} + \frac{\partial \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}}{\partial \varphi_k} \quad (21)$$

$$\frac{\partial |\mathbf{A}|}{\partial \varphi_k} = |\mathbf{A}| \text{Tr}(\mathbf{A}^{-1} \mathbf{L}^k) \quad (22)$$

$$\frac{\partial \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}}{\partial \varphi_k} = -\mathbf{b}^T \mathbf{A}^{-1} \mathbf{L}^k \mathbf{A}^{-1} \mathbf{b}. \quad (23)$$

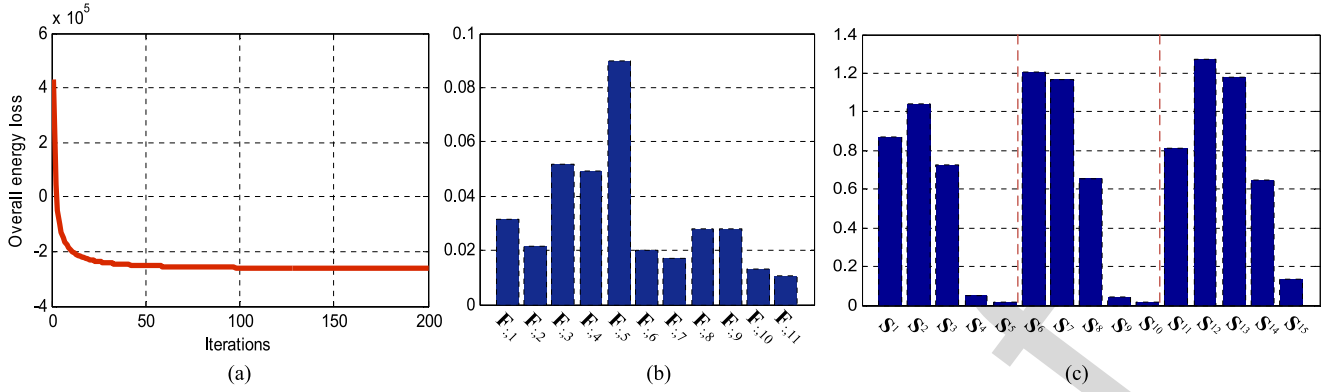


Fig. 6. C-CRF learning outcomes. (a) The overall energy changes of (14) regarding to iterations. (b) The learned α . (c) The learned φ ($T = 4$ case). The notations of features in (b) and (c) are consistent with those in Tables II and III.

500 When α and φ are learned, the saliency inference of C-CRF
 501 on a new test image is achieved by minimizing (9). Setting
 502 $\partial \mathbb{E}(\mathbf{y}, \mathbf{x}) / \partial \mathbf{y} = 0$ leads to the closed-form solution

$$\mathbf{y} = \left(\mathbf{e}^T \alpha \mathbf{I} + \sum_{k=1}^h \varphi_k \mathbf{L}^k \right)^{-1} \mathbf{F} \alpha \quad (24)$$

$$= \mathbf{A}^{-1} \mathbf{b}.$$

503 In (24), the invertibility is guaranteed, since $\alpha_k > 0$, $\varphi_k \geq 0$,
 504 and \mathbf{L}^k is positive semi-definite. Finally, we normalize \mathbf{y} into
 505 $[0, 1.0]$ to render a final saliency map. After the normalization,
 506 we ensure at least one superpixel has value 1 and at least 10%
 507 superpixels have values 0.

508 *A diffusion perspective of C-CRF:* Interestingly, the above
 509 closed-form solution of the C-CRF inference coincides with the
 510 unified formulation of diffusion-based saliency methods raised
 511 in [27], where $(\mathbf{e}^T \alpha \mathbf{I} + \sum_{k=1}^h \varphi_k \mathbf{L}^k)^{-1} = \mathbf{A}^{-1}$ is the diffu-
 512 sion matrix and $\mathbf{F} \alpha$ is the integrated saliency “seed vector” that
 513 leads to a *raw saliency map*. The optimal solution (inference) is
 514 the product of the diffusion matrix and a saliency seed vector,
 515 leading to the equilibrium state vector as the diffused saliency
 516 detection results. However, comparing to [27], we have investi-
 517 gated a different framework—C-CRF, with a totally different
 518 learning strategy.

519 V. EXPERIMENTS AND RESULTS

520 A. Setup

521 Six benchmark datasets were used for our tests, including:
 522 ASD [7] (1000 images), MSRA-B [24] (5000 images), ECSSD
 523 [56] (1000 images), SOD [63] (300 images), SED1 (one-object
 524 image set having 100 images) [64], and SED2 (two-objects
 525 image set having 100 images). Training images were chosen from
 526 MSRA-B. Since the ASD dataset is a subset of MSRA-B, in order
 527 to evaluate the performance on ASD, we first exclude images
 528 that belong to ASD from MSRA-B, resulting in 4000 images
 529 remained. Then we randomly select 3000 images for training and
 530 leave the other 1000 images as the MSRA-B test set. The
 531 trained C-CRF on this dataset is then applied to other datasets.
 532 C-CRF parameters α and φ to be learned were initialized as all-

one vectors. The regularization parameters $\lambda_1 = 1$ and $\lambda_2 = 5$
 533 were set. Since performing the gradient descent on 3000 train-
 534 ing samples are tractable, we used the gradient descent instead
 535 of stochastic gradient descent for learning, in order to achieve
 536 more stable convergence. The learning rate was set as 1×10^{-5} ,
 537 and the convergence was achieved after 200 iterations.
 538

539 During feature extraction, each image was segmented into
 540 $n \approx 200$ superpixels. The maximum graph range T was ini-
 541 tially set to 4 but then pruned to 3 according to the learning
 542 outcomes (see Section V-B for details). Besides, all parameters
 543 for individual unary and pairwise features were empirically set
 544 ($\lambda_e = 10$ and $\lambda_c = 10$ were set for pairwise features).

545 B. Learning Outcomes

546 Fig. 6 shows the learning results of the C-CRF. From Fig. 6(a),
 547 one can see that the overall energy decreases monotonously as
 548 gradient descent proceeds and has reached a stable minimum.
 549 Due to the continuous property of C-CRF, (1) computed on
 550 some images might be larger than 1.0 and it would result in
 551 a negative log-likelihood. This is why in a) the overall energy
 552 turns negative as iteration proceeds. This phenomenon on C-
 553 CRF is different from D-CRF since the solution space of the
 554 latter is finite and countable. Hence for D-CRF, (1) will result
 555 in a probability value instead of a probability density value.

556 Fig. 6(b) and 6(c) show the learned α and φ , respectively. The
 557 learning results in Fig. 6(b) indicate that the geodesic features
 558 $F_{:,1 \sim 5}$ are the most informative ones, which have gained large
 559 weights. Among them $F_{:,5}$ is the most important one. Follow-
 560 ing that, the contrast to image boundary ($F_{:,8}$) and color spatial
 561 distribution ($F_{:,9}$) gain larger weight than the global contrast
 562 ($F_{:,7}$) and center-bias ($F_{:,10}$). This observation is somewhat
 563 consistent with [27] where the center bias feature does not ap-
 564 pear in the top among the listed features. The last feature $F_{:,11}$
 565 (clarity-based) has obtained the lowest weight. The cause of
 566 this is that like most blur detectors, the SVF is based on local
 567 gradient and has limitation in distinguishing between smooth
 568 object surfaces and real blurred image regions [e.g., the 1st row
 569 of Fig. 4(l)].

570 Fig. 6(c) shows the learned φ when setting maximum spa-
 571 tial range $T = 4$. One can see the learned weights of pairwise

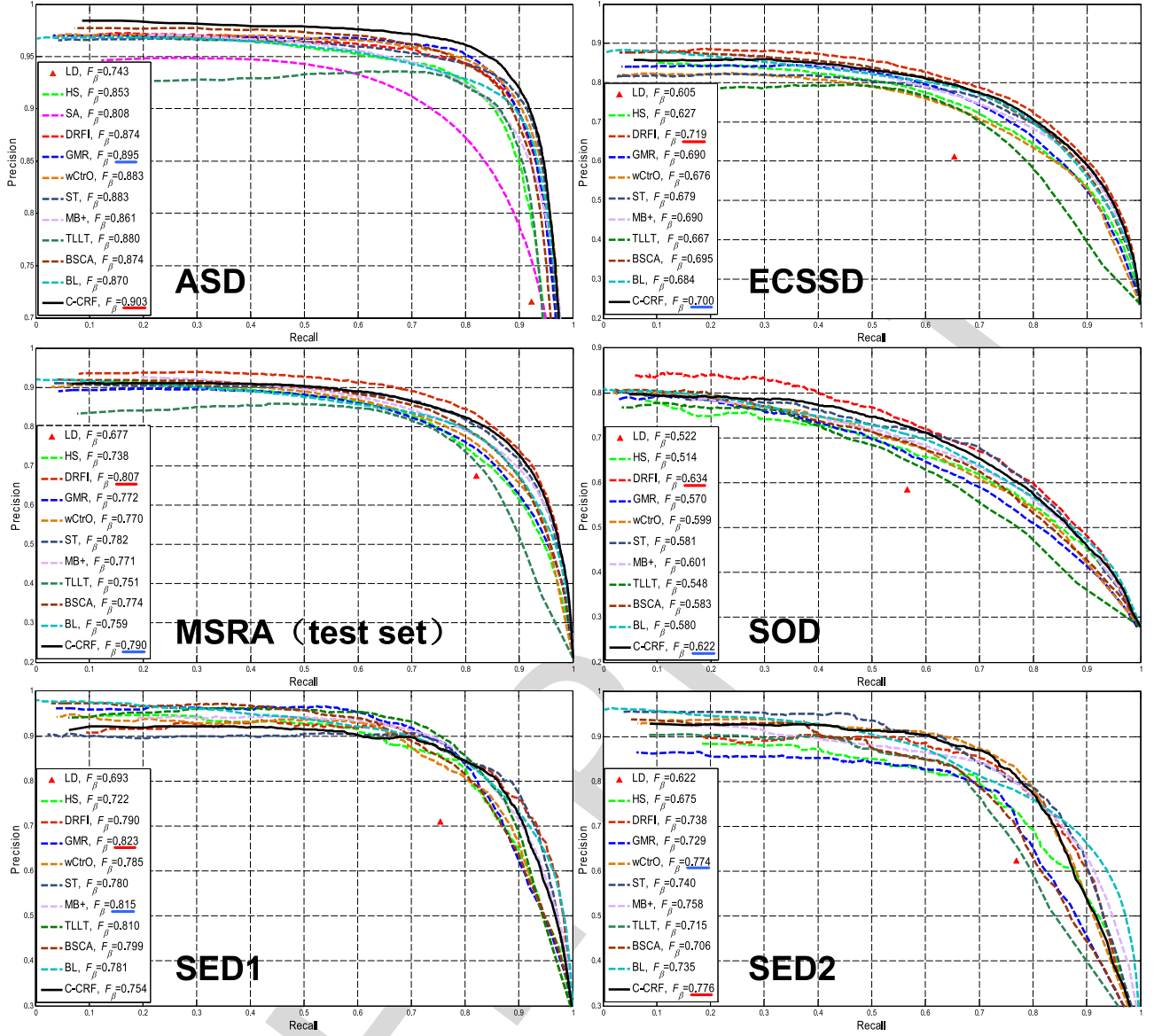


Fig. 7. Quantitative comparisons (precision-recall curves and F_β scores) of the proposed method (C-CRF) to the state-of-the-art methods on six benchmark datasets. The best and the second best F_β are underlined by red and blue.

572 features decrease as the spatial range D^0 increases. This meets
 573 the common sense since spatially close superpixels should
 574 have strong interaction, but noting that such relationship in
 575 our method is automatically learned rather than handcrafted.
 576 Besides, highly degraded weights (close to zeros) for features
 577 S^4 , S^9 , S^{14} that correspond to $D^0 = 3$ and for features S^5 ,
 578 S^{10} , S^{15} that correspond to $D^0 = 4$ reveal further extending
 579 the maximum spatial range to $T \geq 4$ under the current exper-
 580 imental setup is not essential. Considering when $D^0 = 4$, the
 581 corresponding weights are very low. In practice we prune the
 582 spatial range and use $T = 3$, as shown in Figs. 3 and 5. As
 583 shown in Fig. 6(c), the intervening edge cues ($S^{11} \sim S^{15}$) are
 584 the most informative ones among all pairwise features. They
 585 generally gain larger weights than color similarity ($S^1 \sim S^5$)
 586 and histogram intersection ($S^6 \sim S^{10}$). This validates that in-
 587 corporating the edge cues makes contribution. Finally, large

weights of boundary connections (S^1 , S^6 , S^{11}) reveal connect- 588
 ing boundary superpixels is useful. 589

C. Comparison to Existing Methods 590

We compare the proposed saliency detection to 11 exist- 591
 ing methods including: LD (Learning to Detect) [24], 592
 HS (Hierarchical Saliency) [56], SA (Saliency Aggregation) 593
 [26], DRFI (Discriminative Regional Feature Integration) 594
 [43], GMR (Graph-based Manifold Ranking) [25], wCtrO 595
 (background weighted Contrast with Optimization) [50], ST 596
 (Saliency Tree) [65], MB+ (Minimum Barrier Saliency) [66], 597
 TLLT (Teaching-to-Learn and Learning-to-Teach saliency) 598
 [51], BSCA (Background-based Single-layer Cellular Au- 599
 tomata) [52], BL (Bootstrap Learning) [44]. Among them, LD 600
 [24], SA [26], GMR [25] are CRF-related methods listed in 601

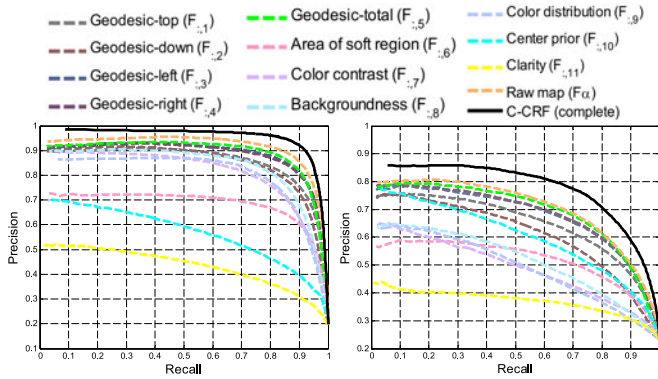


Fig. 8. Precision-recall curves of individual unary features on ASD (left) and ECSSD (right).

602 Table I. Unfortunately, the authors of SA only provide their
 603 results on ASD dataset. Therefore, we can only evaluate SA
 604 on ASD. Besides, the code of [27] is not publicly released, so
 605 the results cannot be compared. For all the compared methods,
 606 we use the public available implementations/results provided by
 607 the authors. Precision-recall curve and F_{β} -measure are used for
 608 evaluating the overall performance [7], [25].

609 Fig. 7 shows the results of precision-recall curves and F_{β}
 610 scores. The proposed method (C-CRF) is comparable to state-
 611 of-the-art methods on both criteria, which has validated the
 612 effectiveness of learning a C-CRF for saliency detection. No-
 613 tably, our method outperforms C-CRF related methods LD, SA,
 614 GMR together with other state-of-the-art methods with notice-
 615 able margins. Regarding to the F_{β} , our method consistently
 616 achieves 1st on ASD, SED2, and the 2nd on ECSSD, MSRA-B
 617 (test set) and SOD. Another data-driven method DRFI some-
 618 times performs better than our method, which may be due
 619 to different feature extraction and learning strategies. Visual
 620 comparisons are shown in Fig. 11.

621 D. Integration of Features/State-of-the-Art Models

622 Since C-CRF is employed in this study as a principled feature
 623 integrating framework, its performance on integrating various
 624 unary and pairwise features should be evaluated. Fig. 8 shows
 625 the precision-recall curves of unary features on ASD and EC-
 626 SSD. One can see that the individual features vary widely on
 627 performance, and among them $F_{.5}$ (which computes the mini-
 628 mum geodesic distance to image borders) achieves the best
 629 results. This coincides with the learning outcomes from MSRA-
 630 B, where $F_{.5}$ gains the highest weight. Observing Fig. 8, the
 631 weighted sum of features (the raw map computed by F_{α})
 632 outperforms all individual features, but the improvement is
 633 relatively marginal. In contrast, the performance is boosted
 634 drastically by a complete C-CRF.

635 To validate the effectiveness of learning for pairwise features,
 636 we treat the C-CRF inference stage (24) as a diffusion process
 637 and replace its diffusion matrix A^{-1} with the propagation matrix
 638 used in GMR [25]. Note GMR is related to C-CRF but with-
 639 out learning (Table I). Its propagation matrix merely consid-
 640 ers the similarity of average colors between superpixels, which

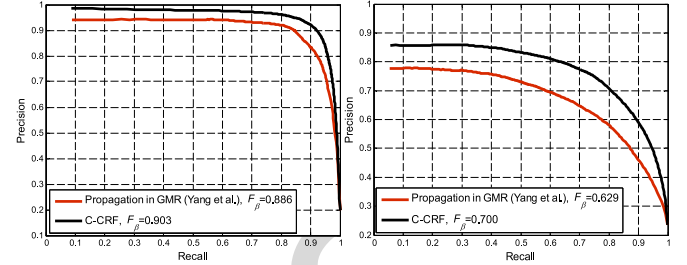


Fig. 9. Effectiveness of integrating pairwise features, validated on ASD (left) and ECSSD (right). In this test, the same “seed” vector F_{α} is used.

intuitively is less effective on representing more sophisticated
 641 interaction between neighboring superpixels, such as fine-
 642 grained color information and texture differences. Fig. 9 shows
 643 the results of this experiment, which validate such an intuition. It
 644 can be seen that by using the same “seed vector” (F_{α}), the diffu-
 645 sion technique employed by [25] is inferior to C-CRF exploited
 646 in this paper.
 647

648 Besides, we validate the power of integrating state-of-the-art
 649 methods by C-CRF, where 5 models are considered: HS, DRFI,
 650 GMR, wCtrO, and MB+. The resulting saliency maps from
 651 these five models are used as the unary feature maps, which are
 652 converted into superpixel-wise maps by averaging pixel-wise
 653 saliency. The C-CRF is then re-trained. Fig. 10 shows the C-
 654 CRF integration performance on ASD, MSRA (test set) and
 655 ECSSD, where the performance boost over individual methods
 656 can be observed on all three datasets. Some visual results from
 657 this experiment are in Fig. 11.

658 E. Effectiveness of Graph Topology Decomposition

659 To show the advantages of learning weights adaptively for
 660 different spatial ranges, we compare to the C-CRF variants
 661 without graph topology decomposition but with manually spec-
 662 ified graph ranges. Here 1-ring graph, 2-ring graph, and 3-ring
 663 graph are considered. Noting an x -ring graph means a superpixel
 664 (graph node) is connected to superpixels within its x -ring neigh-
 665 borhood [25], [27], [61]. Besides, in each graph, the boundary
 666 superpixels are connected with each other as in this paper. For
 667 each one of the three graphs, the three types of pairwise features
 668 namely two color-based ($S_{ij}^{(c)}$, $S_{ij}^{(h)}$) and one image edge-based
 669 ($S_{ij}^{(e)}$) as described in Section IV-C are calculated, resulting
 670 in 3 matrices ($S^{(c)}$, $S^{(h)}$, $S^{(e)}$) for each graph. Except for the
 671 graph ranges, all other C-CRF configurations including unary
 672 features and parameters are kept consistent with Section V-A.
 673 Then for each graph, C-CRF is re-trained and used for saliency
 674 prediction. Fig. 12 shows the quantitative comparison between
 675 the above three graphs and our graph topology decomposition.
 676 It can be observed that the proposed strategy performs more
 677 robustly than an x -ring graph which is manually specified.

678 F. Robustness to The Number of Superpixels

679 Experiments were done by varying superpixel number from
 680 100 to 300, and meanwhile keeping other setup the same as

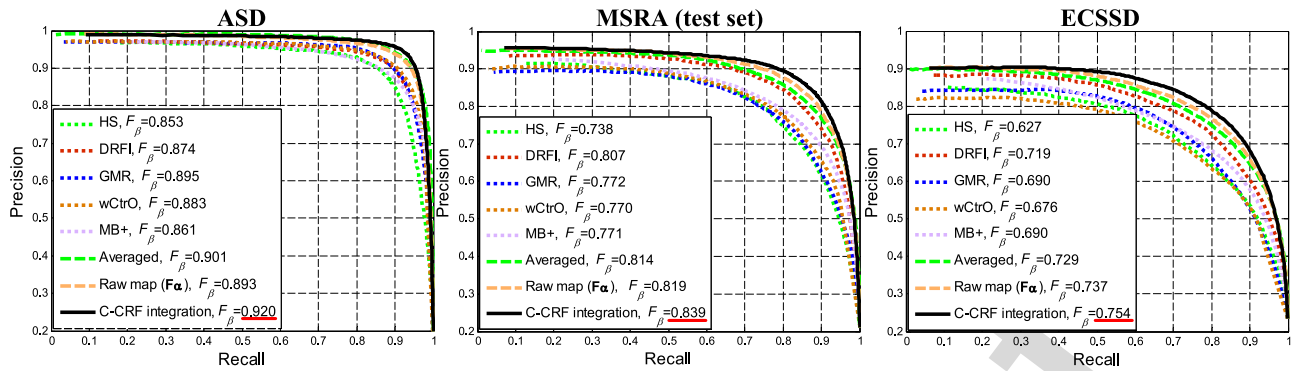


Fig. 10. Integrating five state-of-the-art methods including HS, DRFI, GMR, wCtrO, and MB+ by the proposed C-CRF based framework. The best F_β are underlined by red.

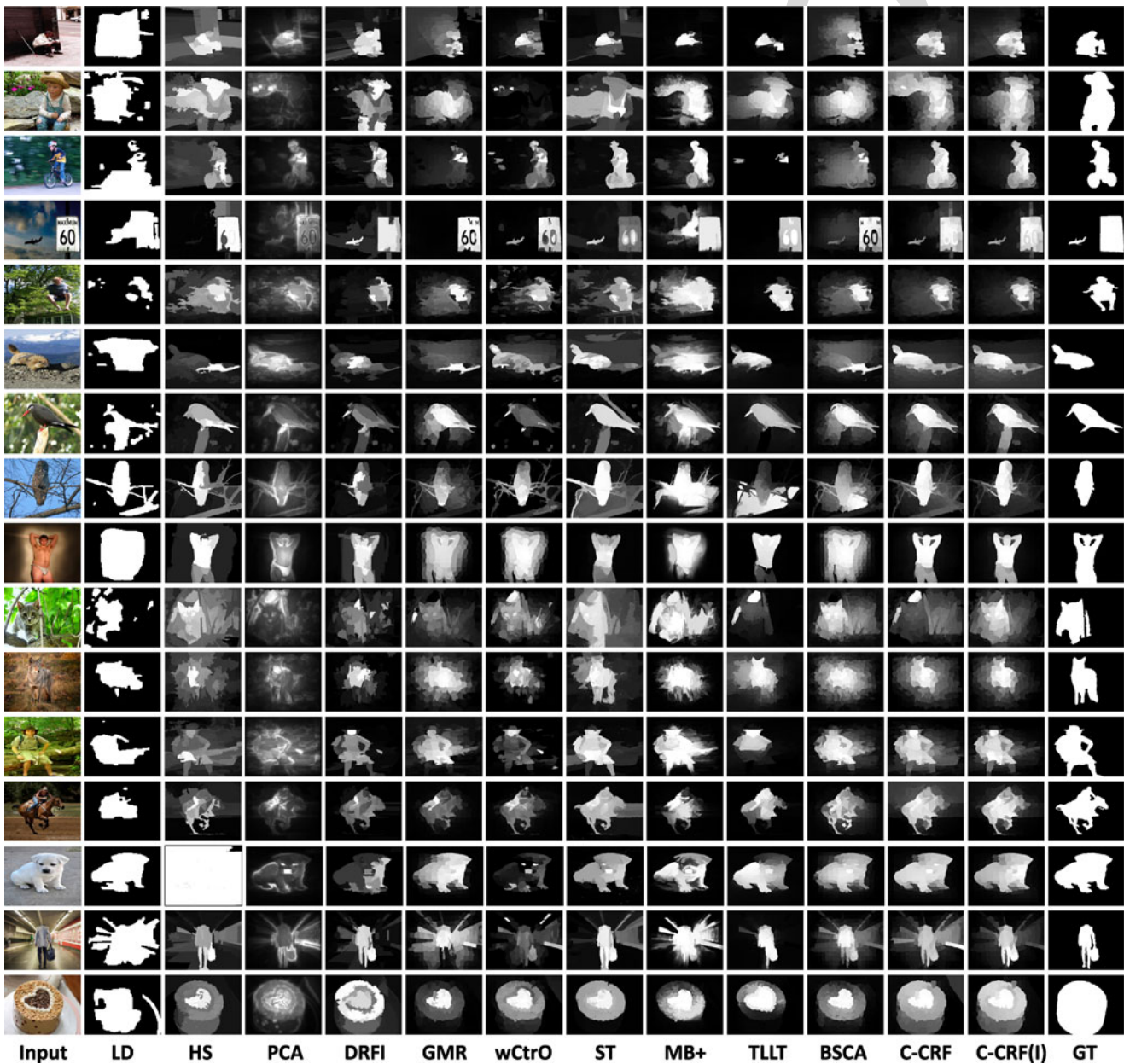


Fig. 11. Qualitative comparison of the proposed method with existing methods on some challenging images with textured background. C-CRF refers to results by integrating different unary saliency features. C-CRF(I) refers to results obtained by integrating five prior models (HS, DRFI, GMR, wCtrO, and MB+), as demonstrated in Section V-D.

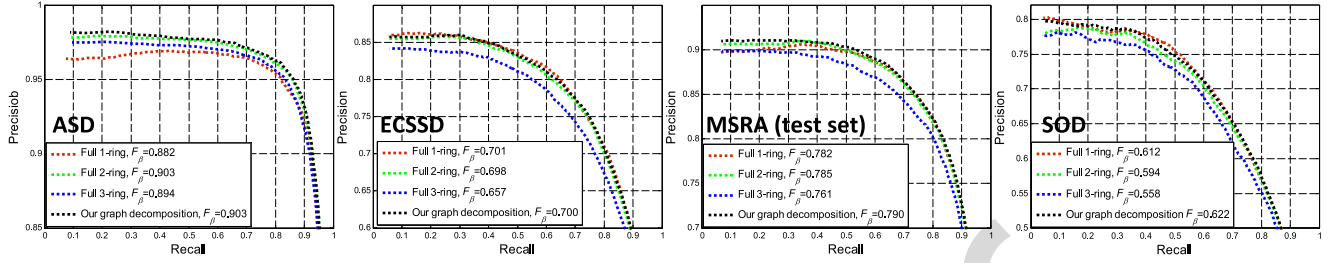


Fig. 12. Quantitative comparisons of one-ring, two-ring, and three-ring graphs and our graph topology decomposition on four benchmark datasets.

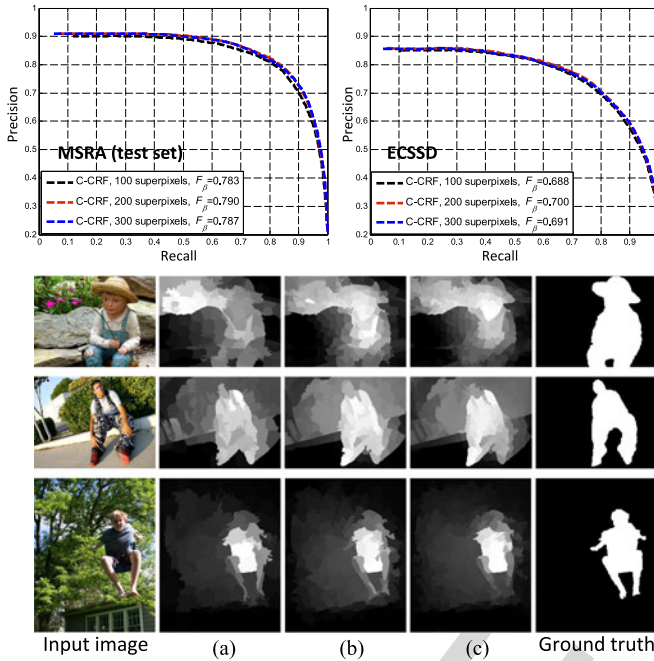


Fig. 13. Quantitative evaluation on MSRA test set (top left) and ECSSD (top right) by using different superpixel numbers. In the bottom some visual comparisons are shown: (a) 100 superpixel case, (b) 200 superpixel case, and (c) 300 superpixel case.

those mentioned in Section V-A. Next, we re-trained the C-CRF model on the training set and then tested it on MSRA-test and ECSSD. Note there are two sides of effect when varying superpixel number: 1) It somewhat affects the computed unary and pairwise features. Fewer superpixels lead to coarser image representation. Fortunately, we observed some robustness of computing unary features $F_{i,j,1 \sim 11}$ to such a change. 2) It also affects the “scale” of the C-CRF objective function, because the dimension of all vectors and matrices involved will change accordingly.

Observing the learning outcomes, we find the overall distribution (or tendency) of learned α and φ is still similar to that in Fig. 6. Fig. 13 shows the evaluation on MSRA-test and ECSSD by using different numbers of superpixels, where robustness to such change can be observed. Using 100 superpixels leads to slightly worse performance as the superpixels become coarser and hence the pre-segmentation is less accurate. Using 200 superpixels and 300 superpixels almost leads to identical performance. Some visual comparisons are shown in Fig. 13. In all, the C-CRF learning and inference is somewhat robust to superpixel

number, therefore graph node numbers. No matter what setup is adopted, C-CRF will learn the optimal feature combination under the current setup.

G. Efficiency

Though the training based on gradient descent from the offline extracted features on 3000 images from MSRA-B took about 4 h, the C-CRF prediction was very fast due to the closed-form solution. It only took 2s in average to process an image from ASD dataset. The superpixel segmentation and attribute extraction (e.g., superpixel colors and histograms) took 0.4 s. The unary feature extraction took 0.45 s, and the pairwise feature extraction took 1.1s including edge detection. The running time was reported on an i7-4720HQ 2.6 GHz laptop with 8 GB memory by Matlab code without optimization.

H. Discussion About the Limitation

Though our C-CRF learning-based method enables effective feature integration and meanwhile boosts the performance from individual saliency features (Fig. 8), the major limitation is its final detection somewhat relies on the quality of input features. If none of the unary saliency features provide reasonable initial saliency estimation, the C-CRF inference will still be bad. Conversely, good features will improve the final detection. This phenomenon can be observed by comparing the quantitative results in Figs. 7 and 10, where employing the state-of-the-art results as unary features leads to better C-CRF inference. A visual example can be found in the 10th row of Fig. 11. One potential solution to this is to enrich features in the feature pool and let the C-CRF discover useful, effective ones through learning.

VI. CONCLUSION

This paper applies the complete learning and inference theories of continuous conditional random field (C-CRF) to salient object detection. The regularized maximum conditional likelihood training by gradient descent optimization is used for parameter learning, and the inference is achieved by an efficient closed-form solution. The power of the proposed method on integrating various unary and pairwise features is tested and evaluated comprehensively. In addition, we propose a novel formulation of pairwise features by graph topology decomposition. The effectiveness on enabling learning weights of different spatial ranges is validated with reasonable learning outcomes. Experimental results and comparison with 11 existing methods show that the proposed method achieves state-of-the-art

743 performance on precision-recall curves with comparable F_{β} -
 744 measure scores. Since the proposed method enables principled
 745 feature integration, in the future some high-level features such as
 746 the category-dependent or semantic features may be incorpor-
 747 ated into the proposed method as top-down influences.

REFERENCES

- 749 [1] A. M. Triesman and G. Gelade, "A feature-integration theory of attention,"
 750 *Cognitive Psychol.*, vol. 12, no. 1, pp. 97–136, 1980.
- 751 [2] C. Koch *et al.*, "Shifts in selective visual attention: Towards the under-
 752 lying neural circuitry," in *Matters of Intelligence*. New York, NY, USA:
 753 Springer-Verlag, 1987, pp. 115–141.
- 754 [3] L. Itti, C. Koch, and E. Niebur, "A model of saliency-based visual attention
 755 for rapid scene analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 20,
 756 no. 11, pp. 1254–1259, Nov. 1998.
- 757 [4] X. Hou and L. Zhang, "Saliency detection: A spectral residual approach,"
 758 in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2007, pp. 1–8.
- 759 [5] N. Bruce and J. Tsotsos, "Saliency based on information maximization,"
 760 in *Proc. Adv. Neural Inf. Process. Syst.*, 2005, pp. 155–162.
- 761 [6] A. Borji and L. Itti, "State-of-the-art in visual attention modeling," *IEEE*
 762 *Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 1, pp. 185–207, Jan. 2013.
- 763 [7] R. Achanta, S. Hemami, F. Estrada, and S. Süsstrunk, "Frequency-tuned
 764 salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern*
 765 *Recog.*, Jun. 2009, pp. 1597–1604.
- 766 [8] M. Cheng, G. Zhang, N. Mitra, X. Huang, and S. Hu, "Global contrast
 767 based salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern*
 768 *Recog.*, Jun. 2011, pp. 409–416.
- 769 [9] Z. Liu, O. Le Meur, S. Luo, and L. Shen, "Saliency detection using regional
 770 histograms," *Opt. Lett.*, vol. 38, no. 5, pp. 700–702, 2013.
- 771 [10] A. Borji, M. Cheng, H. Jiang, and J. Li, "Salient object detection: A
 772 benchmark," *IEEE Trans. Image Process.*, vol. 24, no. 12, pp. 5706–5722,
 773 Dec. 2015.
- 774 [11] E. Rahtu, J. Kannala, M. Salo, and J. Heikkilä, "Segmenting salient ob-
 775 jects from images and videos," in *Proc. Eur. Conf. Comput. Vis.*, 2010,
 776 pp. 366–379.
- 777 [12] L. Wang, J. Xue, N. Zheng, and G. Hua, "Automatic salient object ex-
 778 traction with contextual cue," in *Proc. IEEE Int. Conf. Comput. Vis.*, Nov.
 779 2011, pp. 105–112.
- 780 [13] Z. Liu *et al.*, "Unsupervised salient object segmentation based on kernel
 781 density estimation and two-phase graph cut," *IEEE Trans. Multimedia*,
 782 vol. 14, no. 4, pp. 1275–1289, Aug. 2012.
- 783 [14] F. Stentiford, "Attention based auto image cropping," in *Proc. ICVS Work-*
 784 *shop Comput. Attention Appl.*, 2007, pp. 1–9.
- 785 [15] L. Marchesotti *et al.*, "A framework for visual saliency detection with
 786 applications to image thumbnailing," in *Proc. IEEE Int. Conf. Comput.*
 787 *Vis.*, Sep./Oct. 2009, pp. 2232–2239.
- 788 [16] Y. Ding, X. Jing, and J. Yu, "Importance filtering for image retargeting,"
 789 in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2011, pp. 89–96.
- 790 [17] S. Goferman, L. Zelnik-Manor, and A. Tal, "Context-aware saliency
 791 detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 10,
 792 pp. 1915–1926, Oct. 2012.
- 793 [18] S. Lin, I. Yeh, C. Lin, and T. Lee, "Patch-based image warping for content-
 794 aware retargeting," *IEEE Trans. Multimedia*, vol. 15, no. 2, pp. 359–368,
 795 Feb. 2013.
- 796 [19] T. Chen, M. Cheng, P. Tan, A. Shamir, and S.-M. Hu, "Sketch2photo:
 797 Internet image montage," *ACM Trans. Graph.*, vol. 28, no. 5, 2009,
 798 Art. no. 124.
- 799 [20] Y. Gao *et al.*, "Database saliency for fast image retrieval," *IEEE Trans.*
 800 *Multimedia*, vol. 17, no. 3, pp. 359–369, Mar. 2015.
- 801 [21] Y. Zhang, X. Qian, X. Tan, J. Han, and Y. Tang, "Sketch-based image re-
 802 trieval by salient contour reinforcement," *IEEE Trans. Multimedia*, vol. 18,
 803 no. 8, pp. 1604–1615, Aug. 2016.
- 804 [22] Z. Liu, X. Zhang, S. Luo, and O. Le Meur, "Superpixel-based spatiotem-
 805 poral saliency detection," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 24,
 806 no. 9, pp. 1522–1540, Sep. 2014.
- 807 [23] Z. Liu, W. Zou, L. Li, L. Shen, and O. Le Meur, "Co-saliency detection
 808 based on hierarchical segmentation," *IEEE Signal Process. Lett.*, vol. 21,
 809 no. 1, pp. 88–92, Jan. 2014.
- 810 [24] T. Liu, Z. Yuan, J. Sun, J. Wang, and N. Zheng, "Learning to detect a
 811 salient object," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 2,
 812 pp. 353–367, Feb. 2011.
- [25] C. Yang, L. Zhang, H. Lu, X. Ruan, and M.-S. Yang, "Saliency detection
 813 via graph-based manifold ranking," in *Proc. IEEE Conf. Comput. Vis.*
 814 *Pattern Recog.*, Jun. 2013, pp. 3166–3173.
- [26] L. Mai, Y. Niu, and F. Liu, "Saliency aggregation: A data-driven approach,"
 815 in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2013, pp. 1131–
 816 1138.
- [27] S. Lu, V. Mahadevan, and N. Vasconcelos, "Learning optimal seeds for
 817 diffusion-based salient object detection," in *Proc. IEEE Conf. Comput.*
 818 *Vis. Pattern Recog.*, Jun. 2014, pp. 2790–2797.
- [28] T. Qin, T. Liu, X. Zhang, D. Wang, and H. Li, "Global ranking using
 819 continuous conditional random fields," in *Proc. Adv. Neural Inf. Process.*
 820 *Syst.*, 2008, pp. 1281–1288.
- [29] T. Baltrušaitis, N. Banda, and P. Robinson, "Dimensional affect recogni-
 821 tion using continuous conditional random fields," in *Proc. IEEE Int. Conf.*
 822 *Workshops Automat. Face Gesture Recog.*, Apr. 2013, pp. 1–8.
- [30] F. Liu, C. Shen, and G. Lin, "Deep convolutional neural fields for depth
 823 estimation from a single image," in *Proc. IEEE Conf. Comput. Vis. Pattern*
 824 *Recog.*, Jun. 2015, pp. 5162–5170.
- [31] A. Kolesnikov, M. Guillaumin, V. Ferrari, and C. Lampert, "Closed-form
 825 approximate CRF training for scalable image segmentation," in *Proc. Eur.*
 826 *Conf. Comput. Vis.*, 2014, pp. 550–565.
- [32] J. Shotton, J. Winn, C. Rother, and A. Criminisi, "Textronboost for image
 827 understanding: Multi-class object recognition and segmentation by jointly
 828 modeling texture, layout, and context," *Int. J. Comput. Vis.*, vol. 81, no. 1,
 829 pp. 2–23, 2009.
- [33] P. Krähenbühl and V. Koltun, "Efficient inference in fully connected CRFs
 830 with Gaussian edge potentials," in *Proc. Adv. Neural Inf. Process. Syst.*,
 831 2011, pp. 109–117.
- [34] R. Margolin, L. Zelnik-Manor, and A. Tal, "How to evaluate foreground
 832 maps?" in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2014,
 833 pp. 248–255.
- [35] J. Lafferty, A. McCallum, and F. Pereira, "Conditional random fields:
 834 Probabilistic models for segmenting and labeling sequence data," in *Proc.*
 835 *Int. Conf. Mach. Learn.*, 2001, pp. 282–289.
- [36] S. Nowozin and C. Lampert, "Structured learning and prediction in
 836 computer vision," *Found. Trends Comput. Graph. Vis.*, vol. 6, no. 3/4,
 837 pp. 185–365, 2011.
- [37] C. Yang, L. Zhang, and H. Lu, "Graph-regularized saliency detection
 838 with convex-hull-based center prior," *Signal Process. Lett.*, vol. 20, no. 7,
 839 pp. 637–640, 2013.
- [38] A. Borji, M.-M. Cheng, H. Jiang, and J. Li, "Salient object detection: A
 840 survey," *CoRR*, 2014. [Online]. Available: <http://arxiv.org/abs/1411.5878>
- [39] T. Judd, K. Ehinger, F. Durand, and A. Torralba, "Learning to predict
 841 where humans look," in *Proc. Int. Conf. Comput. Vis.*, Sep./Oct. 2009,
 842 pp. 2106–2113.
- [40] B. Alexe *et al.*, "What is an object?" in *Proc. IEEE Conf. Comput. Vis.*
 843 *Pattern Recog.*, Jun. 2010, pp. 73–80.
- [41] P. Khuvuthyakorn, A. Robles-Kelly, and J. Zhou, "Object of interest
 844 detection by saliency learning," in *Proc. Eur. Conf. Comput. Vis.*, 2010,
 845 pp. 636–649.
- [42] P. Mehrani and O. Veksler, "Saliency segmentation based on learning and
 846 graph cut refinement," in *Proc. Brit. Mach. Vis. Conf.*, 2010, pp. 1–12.
- [43] H. Jiang *et al.*, "Salient object detection: A discriminative regional feature
 847 integration approach," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*,
 848 Jun. 2013, pp. 2083–2090.
- [44] N. Tong, H. Lu, X. Ruan, and M.-H. Yang, "Salient object detection via
 849 bootstrap learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun.
 850 2015, pp. 1884–1892.
- [45] X. Zhou, Z. Liu, G. Sun, L. Ye, and X. Wang, "Improving saliency de-
 851 tection via multiple kernel boosting and adaptive fusion," *IEEE Signal*
 852 *Process. Lett.*, vol. 23, no. 4, pp. 517–521, Apr. 2016.
- [46] X. Wang, L. Zhang, L. Lin, Z. Liang, and W. Zuo, "Deep joint task
 853 learning for generic object extraction," in *Proc. Adv. Neural Inf. Process.*
 854 *Syst.*, 2014, pp. 523–531.
- [47] R. Zhao, W. Ouyang, H. Li, and X. Wang, "Saliency detection by multi-
 855 context deep learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*,
 856 Jun. 2015, pp. 1265–1274.
- [48] R. Achanta *et al.*, "SLIC superpixels compared to state-of-the-art super-
 857 pixel methods," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 11,
 858 pp. 2274–2282, Nov. 2012.
- [49] Y. Wei, F. Wen, W. Zhu, and J. Sun, "Geodesic saliency using background
 859 priors," in *Proc. Eur. Conf. Comput. Vis.*, 2012, pp. 29–42.
- [50] W. Zhu, S. Liang, and Y. Wei, "Saliency optimization from robust back-
 860 ground detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun.
 861 2014, pp. 2814–2821.

- 888 [51] C. Gong *et al.*, "Saliency propagation from simple to difficult," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2015, pp. 2531–2539.
- 890 [52] Y. Qin *et al.*, "Saliency detection via cellular automata," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2015, pp. 110–119.
- 892 [53] K. Petersen and M. Pedersen, *The Matrix Cookbook*. Lund, Sweden: Tech. Univ. Denmark, 2008.
- 894 [54] F. Perazzi, P. Krahenbul, Y. Pritch, and A. Hornung, "Saliency filters: Contrast based filtering for salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2012, pp. 733–740.
- 896 [55] V. Gopalakrishnan, Y. Hu, and D. Rajan, "Salient region detection by modeling distributions of color and orientation," *IEEE Trans. Multimedia*, vol. 11, no. 5, pp. 892–905, Aug. 2009.
- 900 [56] Q. Yan, L. Xu, J. Shi, and J. Jia, "Hierarchical saliency detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2013, pp. 1155–1162.
- 902 [57] B. Su, S. Lu, and C. L. Tan, "Blurred image region detection and classification," in *Proc. ACM Conf. Multimedia*, 2011, pp. 1397–1400.
- 904 [58] J. Kim, D. Han, Y. Tai, and J. Kim, "Salient region detection via high-dimensional color transform," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2014, pp. 883–890.
- 907 [59] M. Cheng *et al.*, "Efficient salient region detection with soft image abstraction," in *Proc. IEEE Int. Conf. Comput. Vis.*, Dec. 2013, pp. 1529–1536.
- 909 [60] P. Dollár and C. L. Zitnick, "Structured forests for fast edge detection," in *Proc. IEEE Int. Conf. Comput. Vis.*, Dec. 2013, pp. 1841–1848.
- 911 [61] J. Sun, H. Lu, and X. Liu, "Saliency region detection based on Markov absorption probabilities," *IEEE Trans. Image Process.*, vol. 24, no. 5, pp. 1639–1649, May 2015.
- 914 [62] K. Fu, C. Gong, I. Gu, and J. Yang, "Normalized cut-based saliency detection by adaptive multi-level region merging," *IEEE Trans. Image Process.*, vol. 24, no. 12, pp. 5671–5683, Dec. 2015.
- 917 [63] V. Movahedi and J. Elder, "Design and perceptual validation of performance measures for salient object segmentation," in *Proc. IEEE Comput. Vis. Pattern Recog. Workshops*, Jun. 2010, pp. 49–56.
- 920 [64] S. Alpert, M. Galun, A. Brandt, and R. Basri, "Image segmentation by probabilistic bottom-up aggregation and cue integration," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 2, pp. 315–327, Feb. 2012.
- 922 [65] Z. Liu, W. Zou, and O. Le Meur, "Saliency tree: A novel saliency detection framework," *IEEE Trans. Image Process.*, vol. 23, no. 5, pp. 1937–1952, May 2014.
- 924 [66] J. Zhang *et al.*, "Minimum barrier salient object detection at 80 fps," in *Proc. IEEE Int. Conf. Comput. Vis.*, Dec. 2015, pp. 1404–1412.



Irene Yu-Hua Gu (M'94–SM'03) received the Ph.D. degree in electrical engineering from the Eindhoven University of Technology, Eindhoven, The Netherlands, in 1992.

From 1992 to 1996, she was a Research Fellow with the Philips Research Institute IPO, Eindhoven, The Netherlands; a Postdoc with Staffordshire University, Staffordshire, U.K.; and a Lecturer with the University of Birmingham, Birmingham, U.K. Since 1996, she has been with the Department of Signals and Systems, Chalmers University of Technology, Gothenburg, Sweden, where she has been a full Professor since 2004. Her research interests include statistical image and video processing, object tracking and video surveillance, pattern classification, and signal processing with applications to electric power systems.

Prof. Gu was an Associate Editor of the IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS, PART A: SYSTEMS AND HUMANS and the IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS, PART B: CYBERNETICS from 2000 to 2005, and an Associate Editor of the *EURASIP Journal on Advances in Signal Processing* from 2005 to 2016. She was the Chair of the IEEE Swedish Signal Processing Chapter from 2001 to 2004. She has been on the Editorial Board of the *Journal of Ambient Intelligence and Smart Environments* since 2011.



Jie Yang received the Ph.D. degree from the Department of Computer Science, Hamburg University, Hamburg, Germany, in 1994.

He is currently a Professor with the Institute of Image Processing and Pattern Recognition, Shanghai Jiao Tong University, Shanghai, China. He has led many research projects (e.g., National Science Foundation, 863 National High Technical Plan), published one book in Germany, and authored more than 200 journal papers. His major research interests include object detection and recognition, data fusion and data mining, and medical image processing.

928
929
930
931
932
933
934
935
936
937
938
939



Keren Fu received the B.S. degree from the Huazhong University of Science and Technology, Wuhan, China, in 2011, and dual Ph.D. degrees from Shanghai Jiao Tong University, Shanghai, China, and Chalmers University of Technology, Gothenburg, Sweden, in 2016, under the joint supervision of Prof. J. Yang and Prof. I. Y.-H. Gu.

His current research interests include visual computing, saliency analysis, and machine learning.

Dr. Fu was the recipient of the National Scholarship by the Ministry of Education in 2013 and 2015.