



## Complementary characteristics fusion network for weakly supervised salient object detection

Liu, Y., Zhang, Y., Wang, Z., Yang, F., Qin, C., Qiu, F., Coleman, S., & Kerr, D. (2022). Complementary characteristics fusion network for weakly supervised salient object detection. *Image and Vision Computing*, 126, 1-14. [104536]. <https://doi.org/10.1016/j.imavis.2022.104536>

[Link to publication record in Ulster University Research Portal](#)

**Published in:**  
Image and Vision Computing

**Publication Status:**  
Published (in print/issue): 31/10/2022

**DOI:**  
[10.1016/j.imavis.2022.104536](https://doi.org/10.1016/j.imavis.2022.104536)

**Document Version**  
Author Accepted version

**General rights**  
Copyright for the publications made accessible via Ulster University's Research Portal is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**  
The Research Portal is Ulster University's institutional repository that provides access to Ulster's research outputs. Every effort has been made to ensure that content in the Research Portal does not infringe any person's rights, or applicable UK laws. If you discover content in the Research Portal that you believe breaches copyright or violates any law, please contact [pure-support@ulster.ac.uk](mailto:pure-support@ulster.ac.uk).

Noname manuscript No.  
(will be inserted by the editor)

# Complementary Characteristics Fusion Network for Weakly Supervised Salient Object Detection

Yan Liu · Yunzhou Zhang\* · Zhenyu Wang · Fei Yang · Cao Qin ·  
Feng Qiu · Sonya Coleman · Dermot Kerr

Received: date / Accepted: date

**Abstract** Salient object detection is a challenging and fundamental research in computer vision and image processing. Although the fully convolutional network has made a great progress in the saliency detection task, most existing methods mainly rely on dense ground truth as labels for training, which takes extensive effort and is time-consuming. This paper proposes a novel and effective scribble-based weakly supervised approach named complementary characteristics fusion network (CCFNet), which learns from easily accessible scribbles such as centerlines instead of fully pixel-wise ground truth. To be more specific, in order to deal with the fact that scribbles are always located inside the objects with lacking annotations close to the semantic boundaries, an edge fusion module is presented to equip our model with the power of aggregating edge information, which would be beneficial to generate saliency maps with more useful information. Alternatively, since scribbles are too sparse to provide enough supervision for the network, we design feature correlation modules based on low-level, high-level global and edge information, which will complement each other to obtain relatively complete salient regions using features of different ways. To further improve the results of saliency maps in foreground and background, a self-supervised saliency detection loss is designed to ensure the network

with stronger generalization ability. Extensive experiments using five benchmark datasets demonstrate that our proposed approach performs favorably against the state-of-the-art weakly supervised algorithms, and even surpasses the performance of those fully supervised.

**Keywords** Salient object detection, Weakly supervised learning, Complementary characteristics fusion network, Self-supervised saliency detection loss

## 1 Introduction

The objective of salient object detection (SOD) is to locate and segment the most dominant objects in a given image. It plays an important role in a variety of computer vision and image processing related fields, such as image manipulation [5,10], robot navigation [6], semantic segmentation [39] and object tracking [53,11].

Following the previous studies, fully deep learning methods have been developed, which broke the limits of traditional handcrafted features since their capability of extracting features at various scales. However, these deep learning based methods usually suffer from a key problem, that they strongly depend on a large volume of accurately labeled data with full pixel-wise annotations for training. It takes extensive effort and time to collect. Therefore, this paper concentrates on designing weakly supervised salient object detection methods based on the sparse labels.

In order to address a trade-off between label efficiency and model performance, some researchers attempted to develop a framework to learn saliency maps from the sparse label [31,45,49,48,21], but there still remains challenges. Image-level category labels are used in [34], which requires large scale datasets with image-level labels. A related work [21] utilized bounding box labels

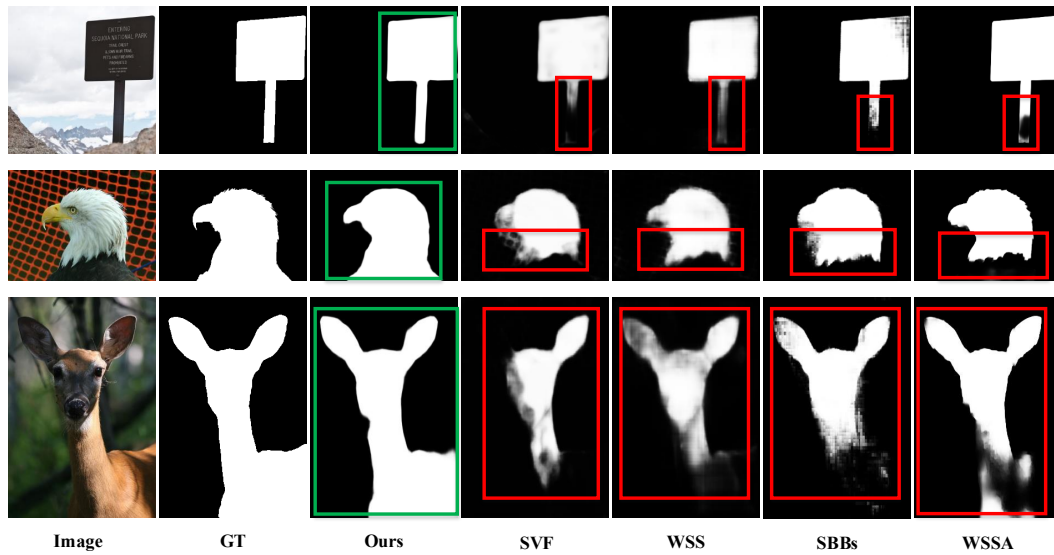
\* Corresponding author.

Y. Liu, and Z. Wang are with Faculty of Robot Science and Engineering, Northeastern University, Shenyang, China.

Y. Zhang, C. Qin and F. Qiu are with College of Information Science and Engineering, Northeastern University, Shenyang, China (E-mail: zhangyunzhou@mail.neu.edu.cn(✉)).

F. Yang is with Zhejiang Lab new campus, Zhongtai, Yuhang, Hangzhou, China.

S. Coleman and D. Kerr are with School of Computing, Engineering and Intelligent Systems, Ulster University, N. Ireland, UK.



**Fig. 1** Sample results of our method compared with other unsupervised and weakly supervised methods.

as supervision, which first produced the initial pseudo ground truth saliency maps by unsupervised learning, then adopted post-processing to obtain the final dense predictions. As shown in Fig. 1, WSS (image-level category labels) and SBBs (bounding box labels) only can detect part of salient regions with wrong errors. For example, given an image with an eagle in Fig. 1, the forementioned approaches are only able to segment the head of the eagle (white) whereas the body is also a salient object. Moreover, scribble annotations [48, 27] are becoming more and more popular in computer vision, which belongs to a middle ground between image-level supervision and box-level supervision. The key problem for saliency detection based on scribble annotations lies in two aspects. The first one is that the scribbles are always located inside the objects with lacking annotations close to the semantic boundaries, and thus usually generate imprecise saliency maps on boundaries. The second one is the scribbles are too sparse to provide enough supervision information for the network, which can't make confident predictions.

With respect to the first issue, as illustrated in [41, 55], edge information has been widely used and has made a great progress in fully supervised saliency detection, weakly supervised SOD models rarely have such ideas. Therefore, edge fusion module (EFM) is employed to capture edge information from local and global views, instead of the simple backbone features, which can effectively improve edge performance of saliency maps. To alleviate the second issue, we propose a feature correlation module (FCM) to capture rich information. Note that different level features have different functions, such as low-level features have rich details and high-level features have rich semantics. Our FCM achieves

complementary characteristics of different input that has a large potential to exploit the relationship from different views. The work in [48] which is different from our approach, only use concatenation operation to fuse different features. Our FCM correlates low-level features, high-level global features and edge features at different stages, which is conducive to enhancing the saliency maps. Furthermore, structural information is also crucial for scribble supervised SOD except for context information. Inspired from [1], we design a self-supervised saliency detection loss to learn structural information, which ensures the network with stronger generalization ability and distinguishes the foreground and background better. As shown in Fig. 1, benefiting from the above, our proposed approach is able to detect more accurate edge information with some challenging environments compared with other methods, such as low contrast scenarios (the background behind the deer) or complex scene understanding (details of the sign).

Based on the above consideration, we propose a complementary characteristics fusion network (CCFNet) for weakly supervised salient object detection. We design edge fusion module to learn salient edge information, which can better understand edge information. In order to exploit complete salient regions with different level features, this paper proposes feature correlation modules for saliency detection. Meantime, the output of global context guiding operation is fed into feature correlation module as input, which could address the high-level features gradually diluted as the top-down pathways. To boost the performance of our proposed model, a self-supervised saliency detection loss is presented as well to distinguish foreground and background. Finally, to demonstrate the performance of our proposed

method, we conduct experiment results on five well-known datasets. Some ablation studies are reported as well to evaluate the effect of each module. From the above, our main contributions can be summarized as follows:

- We develop a novel complementary characteristics fusion network (CCFNet) based on scribble annotations for salient object detection, without resorting to laborious pixel labeling.
- We propose an edge fusion module to equip our model with the power of aggregating edge information. In addition, a feature correlation module is employed to make full use of the complementarities different features to improve saliency detection accuracy.
- We introduce a self-supervised saliency detection loss, which encourages our network to learn structural information and guides the network paying high attention to saliency objects.
- Experimental results demonstrate that the proposed approach achieves comparable performance on five common datasets compared with other state-of-the-art methods, where it even performs comparably to some of the fully supervised methods.

## 2 Related Work

**Fully supervised salient object detection** Traditional SOD approaches mainly depend on some handcrafted features [2, 4, 14, 46] to directly detect salient objects in each image while lacking in high-level semantic information, especially in the complex environments. Compared with early researches on SOD, deep learning based methods [17, 28, 56, 12, 30, 35, 36, 37, 23, 3, 54, 26, 33, 22] have become popular because of their accurate performance. On the one hand, a variety of effective fully convolutional network based (FCN-based) architectures [23, 13, 36, 37] have been proposed to enhance the generation of saliency maps in literature. For example, Hou et al. [13] utilized short connections for multi-scale feature fusion from different layers in FCN to address the scale-space problem. In [36], Wang et al. employed fixation prediction to segment salient objects in an attentive saliency network (ASNet), demonstrating that ASNet achieves more accurate results due to the computed fixation map. The F<sup>3</sup>Net was introduced in [37], to solve the problem generated by the different receptive fields of different convolutional layers. On the other hand, edge information has been attracted attention to assist the performance of saliency prediction [41, 55, 19, 38]. Zhao et al. [55] designed an edge guidance network for salient object detection with binary

cross-entropy. Liu et al. [19] adopt other edge datasets as ground truth for joint training. [55, 41] used edge ground-truth as auxiliary supervision, it proves that is helpful for saliency maps, especially object boundaries. Moreover, a number of related works [3, 54, 56] leveraged the attention mechanism to learn more distinctive features, others [12, 37] introduced multi-level features to boost the performance of saliency maps. Although these methods achieve highly-accurate results, deep models require a large number of fully annotated images when trained on datasets, which is a labor-intensive and costly process.

### Weakly supervised salient object detection

To reduce the time and the cost of labeling, weakly-supervised learning utilize weak labels for the saliency detection task, such as noisy labels [25, 47, 49], bounding boxes [31], scribble annotations [44, 48] and image-level labels [15, 34], which have received a lot of attention from researchers. Currently, Wang et al. [34] adopted foreground inference network for object saliency prediction with image-level labels, which is the first application of image-level labels to SOD. Li et al. [15] subsequently introduced a multi-task fully convolutional network (Multi-FCN) to generate saliency maps using image-level weak supervision. Piao et al [29] built a saliency network and multiple directive filters to enhance the performance of SOD, which is a multiple-pseudo-label framework. Furthermore, S-DUTS was proposed in [48] first on saliency detection, which is based on sparse labels and typically takes 1~2 seconds to label each image. They also designed a network fusing edge detection approach and a gated structure-aware loss function to maintain the accuracy of the salient prediction. Yu et al. [44] introduced a one-round end-to-end training approach using scribbles for weakly-supervised saliency detection. Unlike these methods, we cooperate the characteristics and complementarity of different features to reduce the gap between fully supervised learning and weakly-supervised learning.

**Unsupervised salient object detection** Early methods have been proposed for predicting the saliency map, which mainly used some handcrafted features, contrast, different priors and so on [10, 14, 4]. A related work [47] proposed a deep learning framework from unsupervised methods with heuristics to produce saliency maps. Li et al. [18] developed a contour-to-saliency network based on the well-trained contour detection network. Subsequently, Nguyen et al. [25] presented a two-stage network for unsupervised saliency detection to improve prediction quality, which was updated through noisy labels generated. In conclusion, ground truth is not required for these methods. Unsupervised learning on salient object detection has been made a great and

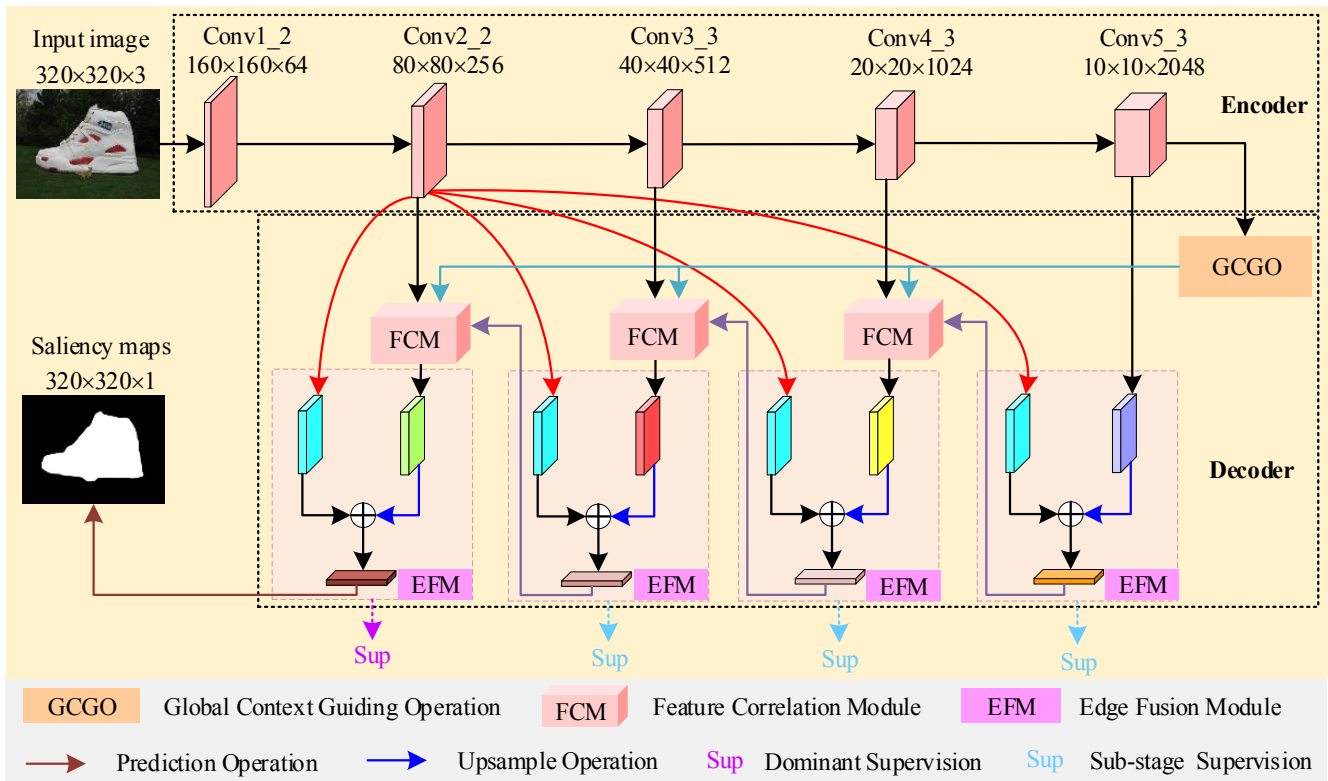


Fig. 2 Illustration of our proposed CCFNet architecture.

231 significant process, but the accuracy is limited due to  
 232 the gap between fully supervised learning and unsuper-  
 233 vised learning.

### 234 3 Methodology

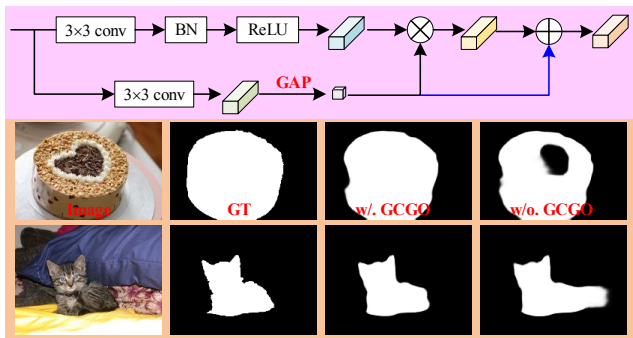
235 In this section, we first introduce the proposed comple-  
 236 mentary characteristics fusion network (CCFNet) for  
 237 weakly supervised salient object detection. Then the  
 238 details of the global context guiding operation (GC-  
 239 GO), the edge fusion module (EFM) and the feature  
 240 correlation module (FCM) are described. The network  
 241 supervision strategy is presented at the end of this sec-  
 242 tion.

#### 243 3.1 Overall pipeline

244 The overall architecture of CCFNet is illustrated in Fig.  
 245 2. Our model is designed based on FCN architecture  
 246 and chooses ResNet-50 as the backbone, which consists  
 247 of five convolutional blocks for feature extracting. Given  
 248 an input image with size  $H \times W$ , the encoder will gener-  
 249 ate different level features, denoted as  $\{f_i | i = 1, \dots, 5\}$   
 250 with resolutions  $[\frac{H}{2^i}, \frac{W}{2^i}]$ . Since the 1st level feature  $f_1$   
 251 would increase computation cost and have a lot of nois-  
 252 es, which yields limited performance improvements, we

253 choose features from  $\{f_i | i = 2, \dots, 5\}$  for later oper-  
 254 ations. Specifically, to alleviate the problem of U-  
 255 shape networks as top-down ways gradually diluted,  $f_5$   
 256 is fed into GCGO to obtain  $\{g_i | i = 1, \dots, 3\}$ , which  
 257 can guarantee global semantics delivered. Since low-  
 258 level features have more details such as boundaries,  
 259 which are useful and indispensable for generating accu-  
 260 rate saliency maps, we leverage the  $f_2$  to extract  
 261 the boundaries. In contrast, high-level features have  
 262 more semantics but lacking details. Taking account of  
 263 these considerations, we aim to explicitly notice the  
 264 salient edges where salient objects are. Hence we mod-  
 265 el EFM to strengthen edge information, denoted as  
 266  $\{e_i | i = 1, \dots, 4\}$ . Besides, in view of different type-  
 267 s of features delivering different information, to this  
 268 end, we design FCM in this paper. FCM is performed  
 269 to refining low-level features  $\{f_i | i = 2, \dots, 4\}$ , global  
 270 high-level features  $\{g_i | i = 1, \dots, 3\}$  and edge features  
 271  $\{e_i | i = 1, \dots, 4\}$ . This way enables the network to un-  
 272 derstand scenarios from different views, which will gen-  
 273 erate the discriminative features. It may limit the cap-  
 274 ability of the network due to only choosing scribbles  
 275 to train our network. To address this limitation, we al-  
 276 so propose a self-supervised saliency detection loss for  
 277 joint training to enrich structural information. More de-  
 278 tails of CCFNet are described as follows.





**Fig. 3** Illustration of our proposed global context guiding operation (GCGO).

### 279 3.2 Global context guiding operation

Regarding the U-shape architecture exists an issue that the high-level features will be gradually diluted as the top-down pathways. Therefore, we propose a global context guiding operation to strengthen high-level information and obtain global information, as shown in Fig. 3. Specifically, we apply a combination of  $3 \times 3$  convolutional  $\rightarrow$  batch normalization  $\rightarrow$  ReLU operation for input feature  $f_5$ . After that, a global average pooling (GAP) layer is embedded on these features, denoted as  $f_g$ , which can capture a more robust spatial translations of the input and the strongest global context. The refined feature  $f_g$  is denoted as follows:

$$f_g = GAP(\sigma(\phi(Conv(f_5, \theta)))), \quad (1)$$

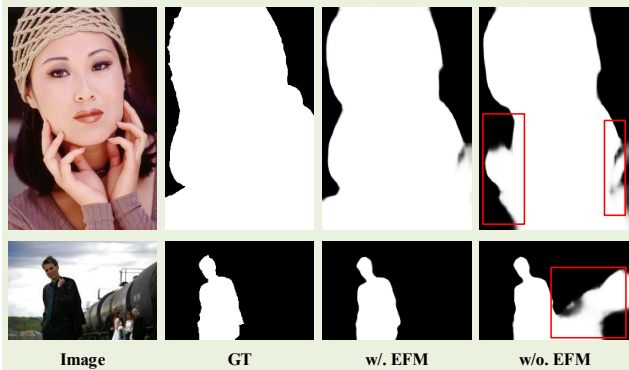
280 where each of  $Conv(\cdot, \theta)$ , denotes the convolution with  
 281 parameter  $\theta$ ,  $\phi(\cdot)$ ,  $\sigma(\cdot)$  and  $GAP(\cdot)$  denotes the batch  
 282 normalization, Relu and global average pooling, respec-  
 283 tively. Meantime, we apply  $3 \times 3$  convolutional opera-  
 284 tion to input features to squeeze the input feature  $f_5$   
 285 and adopt a upsample operation, which retain useful  
 286 information. Finally, we generate the mask  $\mathbf{W}$  and bias  
 287  $\mathbf{b}$  for multiplication and addition operation. The whole  
 288 process is formally formulated as follows.

$$g_1 = \sigma(\mathbf{W} * f_g + \mathbf{b}), \quad (2)$$

289 where  $*$  is element-wise multiplication and  $g_1$  is the final  
 290 output. From Fig. 3, we can clearly see that with GCGO  
 291 strategy achieves better performance than without it.  
 292 Detailed quantitative studies of GCGO can be found in  
 293 Section 4.

### 294 3.3 Edge fusion module

295 Ideally, a good weakly supervised salient object detec-  
 296 tion algorithm should have the ability to capture accu-  
 297 rate edge information. In other words, salient edge re-  
 298 sult is able to help salient object detection tasks in both



**Fig. 4** Visual results by applying EFM and without EFM.

segmentation and localization. To this end, we propose  
 a series of edge fusion modules (EFM) to model the  
 salient edge information. As stated before, the  $f_2$  re-  
 tains edge information, even so, it is still local infor-  
 mation and not enough. We take account of high-level  
 semantics, which are essential and necessary for obtain-  
 ing salient edge information as well.

To be more specific, taking the first EFM as an ex-  
 ample, we take a  $3 \times 3$  convolutional layer after extract-  
 ing the edge feature from  $f_2$ . In order to increase the  
 reliability of salient edge information, we fuse high-level  
 semantic information from  $f_5$ . We add a convolutional  
 operation with kernel size  $3 \times 3$  after fusing both  
 features, it is able to effectively reduce the aliasing ef-  
 fect of upsampling. For the other EFMs, our goal is the  
 high level cue mined is applied over the corresponding  
 feature, which is further propagated to the next EFM  
 for edge generation. That is to say, the feature maps  
 from the corresponding feature correlation module are  
 replaced of the backbone feature from  $f_5$ . The whole  
 process is formally formulated as follows.

$$h_i = \begin{cases} Conv(f_2, \theta) + Up(f_5), & \text{if } i = 1 \\ Conv(f_2, \theta) + Up(\Psi_i), & \text{if } i = 2, 3, 4 \end{cases} \quad (3)$$

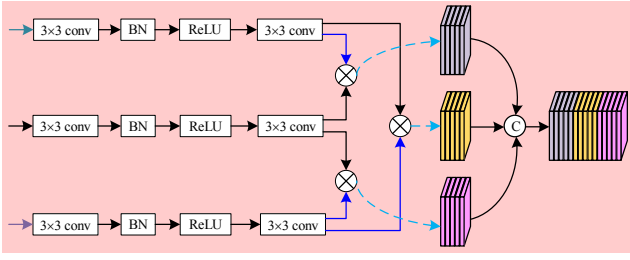
where  $\Psi_i$  is the output of feature correlation module.  
 Hence the final output of EFM can be described as  
 follows.

$$e_i = \theta(\phi(Conv(h_i, \theta))) \quad i = 1, 2, 3, 4 \quad (4)$$

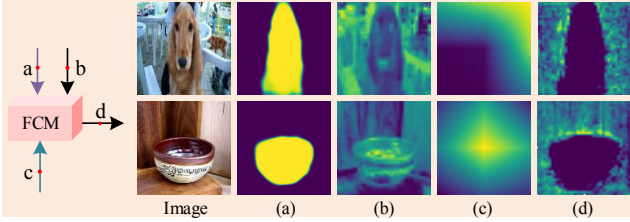
where  $e_i$  is the final output. To verify the effectiveness  
 of our designed EFM, we visualize the saliency maps  
 by applying EFM or not in Fig 4, it is clearly observed  
 that our model with EFM has high quality in edge.

### 310 3.4 Feature correlation module

Given three pathway features: low-level features, global  
 high-level features and edge features in CCFNet, which



**Fig. 5** Illustration of our proposed feature correlation module (FCM).



**Fig. 6** Visualization of the feature maps around FCM. (a) Results of applying EFM. (b) Results of applying backbone. (c) Results of applying GCGO. (d) The output results of FCM.

313 can better preserve details, global semantics and edge  
 314 information, respectively. However, there exists a issue  
 315 that a single feature provides locally limited informa-  
 316 tion. To solve this deficiency, considering that these fea-  
 317 tures are complementary to each other, it is essential to  
 318 form an effective decoder to strengthen the quality of  
 319 saliency maps.

To this end, we define the feature correlation mod-  
 320 ule (FCM) to get rich features from different pathway  
 321 in this section, which is able to produce saliency maps  
 322 with accurate segmentation. The details of FCM is il-  
 323 lustrated in Fig. 5. Formally, we first conduct a series  
 324 of operations:  $Conv(3 \times 3, 256) \rightarrow BN \rightarrow ReLU \rightarrow$   
 325  $Conv(3 \times 3, 256)$ , denoted as  $F_i$ ,  $G_i$  and  $E_i$  respec-  
 326 tively. Inspired by the attention mechanism [3], we evalu-  
 327 ate the interaction between the any embedding features  
 328 that is used to generate a global feature by aggregating  
 329 every local feature. Therefore, we adopt upsample to  
 330 feature  $G_i$  (cyan line as input) firstly so that it has the  
 331 same size as feature  $F_i$  (black line as input), then the  
 332 mutual influence of feature  $G_i$  and feature  $F_i$  is achieved  
 333 by element-wise multiplication, that is  $Up(G_i) * F_i$ . This  
 334 way is used to capture more discriminative characteras-  
 335 tics representation from global context and details. On  
 336 the one hand, it is able to obtain features from details  
 337 and boundaries, on the other hand, it can gain fea-  
 338 tures in global dimension and in edge. Note that these  
 339 features complement each other to form feature correla-  
 340 tion module with more discerning capabilities, which  
 341 are critical for salient detection. The whole process is

defined as follows.

$$\begin{cases} \mathcal{Y}_i = Cat(Up(G_i) * F_i, \\ Up(E_i) * F_i, Up(E_i * G_i)), & \text{if } i = 1, 2, 3 \\ \Psi_i = Conv(\mathcal{Y}_i, \theta), \end{cases} \quad (5)$$

320 where  $Cat(\cdot, \cdot)$  denotes concatenate operation and  $\Psi_i$   
 321 denoted the output of  $i$ th FCM, respectively. Further-  
 322 more, to verify the rationality of our proposed FCM,  
 323 we visualize the feature maps near the FCM in Fig.  
 324 6, which can see that FCM are helpful and combin-  
 325 ing them together are able to remedy for the deficiency  
 326 of each branch feature. Detailed quantitative studies of  
 327 FCM can be found in Section 4.

### 3.5 Self-supervised saliency detection loss

328 In view of the fact that we only choose scribble an-  
 329 notations to train our network, which contains a large  
 330 number of unlabeled pixels and thus may limit the capa-  
 331 bility of the network. The proposed modules focus  
 332 on obtaining context information, whereas structural  
 333 information also plays an important role in scribble su-  
 334 pervised saliency object detection. Partial cross-entropy  
 335 (PCE) loss [32] is widely used to weakly supervised  
 336 learning. However, it only calculates binary cross-entropy  
 337 loss between the scribbles and the predicted map while  
 338 not comprehensive for saliency detection. Based on this  
 339 consideration, to encourage better saliency maps with  
 340 more structural information from the network, we pro-  
 341 pose a novel self-supervised saliency detection (SSD)  
 342 loss to help the network better distinguish foreground  
 343 and background.

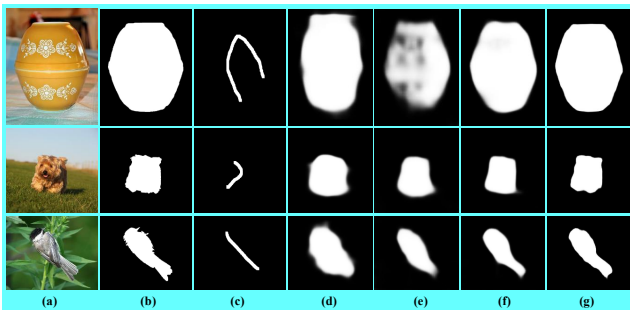
As shown in Fig. 7, it can be clearly seen that the  
 344 results of the first two phases (d & e) are not as good  
 345 as that of the third (f). Considering that the later EFM  
 346 maintains more information than the former, and the  
 347 results of the last EFM is the prediction result of the w-  
 348 hole network. To learn more structural information and  
 349 guide the network paying high attention to saliency ob-  
 350 jects, we produce pseudo ground truth masks from the  
 351 penult EFM by considering confidence  $> \%60$  (that is  
 352 MAE scores  $< 0.6$ ), which generate saliency maps that  
 353 are closer near to ground truth, rather than scribbles.  
 354 Note that the pixels with low confidence are ignored by  
 355 the loss function. To this end, we design a gate func-  
 356 tion to judge if it satisfy the needs of the above. The details  
 357 are described as follows.

$$g(x, y) = \begin{cases} 1, & \text{if } PSE_{MAE}(x, y) < 0.6 \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

where  $PSE_{MAE}(x, y)$  is the MAE scores of the results  
 358 from the penult EFM. Motivated by semantic segmen-  
 359 tation method [1], we use pixels cross-entropy loss, but

**Table 1** Main characteristics of the datasets used in the experiments.

Name	Stage	Size	Description
S-DUTS [48]	Train	10553	Relabel salient object detection dataset DUTS-TR with scribbles.
ECSSD [42]	Test	1000	A dataset includes semantically meaningful and complex structures.
DUT-OMRON [43]	Test	5168	A dataset with high quality and challenging images has one or more salient objects with complex background scenes.
PASCAL-S [7]	Test	850	A dataset selected from the PASCAL VOC 2010 segmentation, which contains 20 object categories and complex scenes.
HKU-IS [16]	Test	4447	Contain multiple salient objects with overlapping objects touching the image boundary or with low color contrast.
DUTS-TE [34]	Test	5019	It selected from the largest salient object detection benchmark dataset DUTS, which contains complex scens in different scales.

**Fig. 7** Intermediate results at training time. (a) Input image. (b) Per-pixel wise ground truth. (c) Scribble annotations. (d) Results of applying first EFM. (e) Results of applying second EFM. (f) Results of applying third EFM (pseudo ground truth). (g) Ours.

the loss for saliency objects are normalized according to the number of corresponding pixels contained in the pseudo ground truth. Hence, SSD loss can be described as follows.

$$L_{ssd} = g(x, y) L_{bce}, \quad (7)$$

where

$$L_{bce} = - \sum_{(x,y)} [p(x, y) \log(q(x, y)) + (1 - p(x, y)) \log(1 - q(x, y))], \quad (8)$$

345 where  $p(x, y)$  and  $q(x, y)$  denote the pseudo ground  
346 truth masks and the predicted saliency maps, respec-  
347 tively.

### 348 3.6 Objective Function

Given an input image, we utilize the loss of each sub-stage and the dominant loss to train our model. First, the loss of sub-stage is defined as follows.

$$L_{sub} = L_{pce} + L_{lsc}, \quad (9)$$

where

$$L_{pce} = \sum_{i \in S} -s_i \log \hat{s}_i + (1 - s_i) \log(1 - \hat{s}_i), \quad (10)$$

$$L_{lsc} = \sum_i \sum_{j \in G_i} F(i, j) D(i, j), \quad (11)$$

Here, Eq. (10) is partial cross-entropy (PCE) loss, which is widely used to weakly supervised learning, where  $s$  is the scribble annotations,  $\hat{s}$  is the predicted values and  $S$  is the labeled pixel set. However, due to PCE loss only calculating binary cross-entropy loss between the scribbles and the predicted map, it is not comprehensive for saliency detection. In order to further use scribble annotations, local saliency coherence (LSC) loss Eq. (11) is adopted in previous work [44], where  $F(i, j)$  is Gaussian kernels and  $D(i, j)$  is  $L1$  distance. Second, self-supervised saliency detection will joint with PCE loss and LSC loss to supervise in this paper. The dominant loss can be described as follows.

$$L_{dom} = L_{pce} + L_{lsc} + L_{ssd}. \quad (12)$$

Hence, the total loss in the whole network can be expressed as follows.

$$L = L_{dom} + \gamma_i \sum_{i=1}^3 L_{sub}^i, \quad (13)$$

349 where  $\gamma_i$  is the a coefficient to balance the dominant  
350 loss and the different sub-stage loss. Because different  
351 sub-stage provides various extend of information, we set  
352  $\gamma_1 = 0.8, \gamma_2 = 0.6, \gamma_3 = 0.4$  in this paper.

## 353 4 Experiments

### 354 4.1 Implementation details

355 The proposed approach was implemented on the Py-  
356 torch platform using a RTX3090 GPU. The batch size



is set as 16. The whole network is optimized by stochastic gradient descent (SGD), where the weight decay is set to  $5e-4$ , the momentum is set to 0.9 and the initial learning rate  $1e-5$ . We resize each image to  $320 \times 320$  and then feed into the network to obtain saliency maps. Additionally, the characteristics of each dataset are summarized in Table 1.

## 4.2 Evaluation metrics

For the salient object detection task, six popular evaluation metrics are used to evaluate the effectiveness of our CCFNet including precision-recall curve (PR curve), F-measure curve, F-measure score ( $F_\beta$ ), mean absolute error (MAE), E-measure score ( $E_\phi$ ) and S-measure score ( $S_\alpha$ ).

**PR curve** can be determined by generated pairs of precision and recall values. Precision and recall are computed as:

$$Precision = \frac{TP}{TP + FP}, \quad Recall = \frac{TP}{TP + FN}, \quad (14)$$

where  $TP$ ,  $FP$  and  $FN$  denote true-positive, false-positive and false-negative, respectively.

**F-measure score** ( $F_\beta$ ) is an overall performance measurement, which is calculated by the weighted harmonic mean of precision and recall:

$$F_\beta = \frac{(1 + \beta^2) \times Precision \times Recall}{\beta^2 \times Precision + Recall}, \quad (15)$$

where  $Precision$  and  $Recall$  are given by thresholding the predicted saliency map, and  $\beta^2$  is set to 0.3 in accordance with [3]. Then the obtained pairs (threshold,  $F_\beta$ ) is employed to plot the F-measure curve.

**MAE** reflects the average pixel-wise absolute difference between the saliency map  $S(x, y)$  and ground-truth maps  $G(x, y)$ :

$$MAE = \frac{1}{W \times H} \sum_{x=1}^W \sum_{y=1}^H |S(x, y) - G(x, y)|, \quad (16)$$

where  $W$  and  $H$  represent width and height of the saliency maps respectively. This is an appropriate metric for evaluating the applicability of a saliency module in a task such as image segmentation.

**Enhanced-alignment measure**  $E_\phi$  [9] is applied to evaluate both local and global similarity between the predicted map and the ground-truth:

$$E_\phi = \frac{1}{W \times H} \sum_{x=1}^W \sum_{y=1}^H F_\phi(x, y), \quad (17)$$

$F_\phi$  denotes the enhanced alignment matrix.

**Structure measure**  $S_\alpha$  [8] is utilized as the structure similarity of the predicted non-binary saliency map and the ground-truth, which is defined as follows:

$$S_\alpha = (1 - \alpha)S_r + \alpha S_o, \quad (18)$$

where  $S_r$  and  $S_o$  denote region-aware and object-aware structural similarity respectively, and  $\alpha$  is typically set to 0.5.

## 4.3 Comparison with the State-of-the-Art Methods

We compare our model with state-of-the-art nineteen methods, including eleven fully supervised methods (Amulet [51], UCF [52], NLDF [24], RAS [3], PAGR [54], BMPM [50], DSS [12], EGNet [55], CPD [40], MINet [28] and VST [20]), two unsupervised methods (SVF [47] and C2S [18]), six weakly supervised methods (WSS [34], ASMO [15], MWS [45], WSSA [48], MFNet [29] and SBBs [21]). For fair comparison, all the saliency maps are provided by the authors. In addition, our results are diametrically produce by CCFNet without relying on any post-processing.

**Quantitative comparison** The detailed F-measure, MAE, E-measure and S-measure values are provided in Table 2 and Table 3 on five common datasets, in which our approach performs favorably against other state-of-the-art unsupervised and weakly supervised approaches by a large margin, and even superior to some fully supervised methods, like Amulet, UCF and NLDF. It is worth noting that we also achieve the best results for saliency detection using challenging datasets, such as DUT-OMRON and DUTS-TE. For fairness, in terms of the average of each metric, we can conclude that our proposed approach shows a preferred average  $F_\beta$  (0.809 vs. 0.779), MAE (0.061 vs. 0.069),  $E_\phi$  (0.880 vs. 0.875) and  $S_\alpha$  (0.847 vs. 0.818) across five datasets than SBBs, which is the latest competitive algorithm. We have to admit that there are some gaps compared with the fully supervised algorithm in terms of performance even though it is reasonable and logical. Generally speaking, our approach is superior to other counterparts across all datasets using these evaluation metrics. Besides, Figure 8 and Figure 9 show the PR curves and F-measure curves with other weakly supervised state-of-the-art methods on the five benchmark datasets, respectively. It can be observed that our method achieves a better performance than the other ones in most cases. To further analyze the overall difference between our algorithm and other methods, we compare the quantitative results including average F-measure, average E-measure, average S-measure and average MAE on five common datasets, which are calculated by the average

**Table 2** Comparison with other state-of-the-art approaches on ECSSD and DUT-OMRON datasets. 'F' means fully supervised, 'W' means weakly supervised and 'Un' is for unsupervised.  $\uparrow$  &  $\downarrow$  denote larger and smaller is better, respectively.

Methods	Year	Sup.	ECSSD 1000 images				DUT-OMRON 5168 images			
			$F_\beta \uparrow$	MAE $\downarrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$	$F_\beta \uparrow$	MAE $\downarrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$
Amulet [51]	ICCV 2017	F	0.868	0.059	0.901	0.894	0.647	0.098	0.779	0.781
UCF [52]	ICCV 2017	F	0.844	0.069	0.892	0.883	0.621	0.120	0.765	0.760
NLDF [24]	CVPR 2017	F	0.878	0.063	0.910	0.875	0.684	0.080	0.816	0.770
RAS [3]	ECCV 2018	F	0.889	0.056	0.914	0.893	0.713	0.062	0.846	0.814
PAGR[54]	CVPR 2018	F	0.894	0.061	0.914	0.889	0.711	0.071	0.842	0.775
BMPM [50]	CVPR 2018	F	0.868	0.045	0.914	0.911	0.692	0.064	0.837	0.809
DSS [12]	TPAMI 2019	F	0.904	0.052	0.912	0.882	0.740	0.063	0.842	0.790
EGNet [55]	ICCV 2019	F	0.920	0.037	0.927	0.925	0.755	0.053	0.868	0.841
CPD [40]	CVPR 2019	F	0.917	0.037	0.925	0.918	0.747	0.056	0.866	0.825
MINet [28]	CVPR 2020	F	0.924	0.033	0.927	0.925	0.755	0.056	0.865	0.833
VST [20]	ICCV 2021	F	0.920	0.033	0.918	0.932	0.756	0.058	0.861	0.850
SVF [47]	ICCV 2017	Un	0.809	0.088	0.875	0.832	0.608	0.108	0.768	0.747
C2S [18]	ECCV 2018	Un	0.853	0.059	0.906	0.882	0.664	0.079	0.817	0.780
WSS [34]	CVPR 2017	W	0.823	0.104	0.869	0.811	0.603	0.109	0.768	0.725
ASMO [15]	AAAI 2018	W	0.798	0.110	0.853	0.802	0.622	0.101	0.776	0.752
MWS [45]	CVPR 2019	W	0.840	0.096	0.884	0.827	0.609	0.109	0.763	0.756
WSSA [48]	CVPR 2020	W	0.870	0.059	0.901	0.865	0.703	0.068	0.840	0.785
MFNet [29]	ICCV 2021	W	0.844	0.084	0.877	0.837	0.621	0.098	0.783	0.726
SBBs [21]	TIP 2021	W	0.855	0.072	0.894	0.851	0.695	0.074	0.835	0.776
Ours	–	W	<b>0.890</b>	<b>0.050</b>	<b>0.912</b>	<b>0.882</b>	<b>0.720</b>	<b>0.069</b>	<b>0.848</b>	<b>0.796</b>

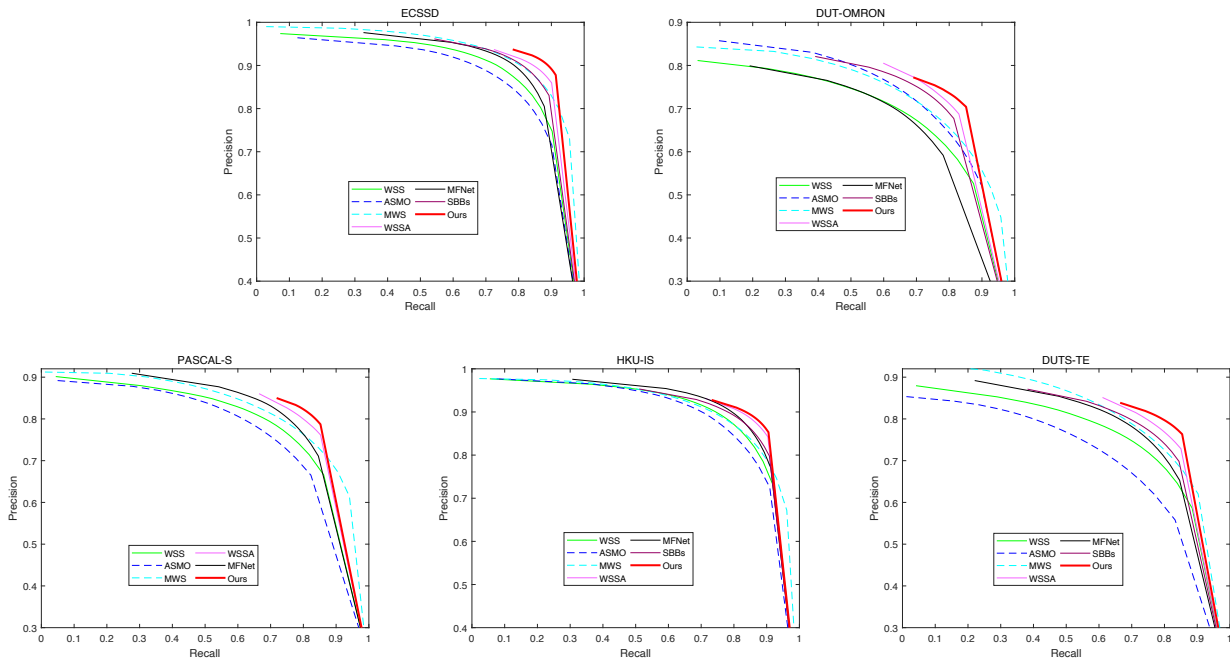
**Table 3** Comparison with other state-of-the-art approaches on PASCAL-S, HKU-IS and DUTS-TE datasets. 'F' means fully supervised, 'W' means weakly supervised and 'Un' is for unsupervised.  $\uparrow$  &  $\downarrow$  denote larger and smaller is better, respectively. "–" means the authors did not release the code, and they just provided the saliency maps, thus reporting the total number of parameters of this method is not possible.

Methods	Year	Sup.	PASCAL-S 850 images				HKU-IS 4447 images				DUTS-TE 5019 images			
			$F_\beta \uparrow$	MAE $\downarrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$	$F_\beta \uparrow$	MAE $\downarrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$	$F_\beta \uparrow$	MAE $\downarrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$
Amulet [51]	ICCV 2017	F	0.757	0.100	0.802	0.818	0.841	0.051	0.912	0.886	0.678	0.085	0.794	0.804
UCF [52]	ICCV 2017	F	0.726	0.115	0.804	0.805	0.823	0.062	0.902	0.875	0.631	0.112	0.763	0.782
NLDF [24]	CVPR 2017	F	0.769	0.098	0.839	0.805	0.874	0.048	0.929	0.887	-	-	-	-
RAS [3]	ECCV 2018	F	0.777	0.101	0.836	0.799	0.871	0.045	0.929	0.887	0.751	0.059	0.861	0.839
PAGR[54]	CVPR 2018	F	0.798	0.089	0.853	0.822	0.886	0.048	0.939	0.887	0.784	0.056	0.880	0.838
BMPM [50]	CVPR 2018	F	0.758	0.074	0.842	0.845	0.871	0.039	0.937	0.907	0.745	0.049	0.860	0.862
DSS [12]	TPAMI 2019	F	0.801	0.093	0.847	0.798	0.902	0.040	0.934	0.878	-	-	-	-
EGNet [55]	ICCV 2019	F	0.817	0.074	0.854	0.852	0.902	0.031	0.949	0.918	0.815	0.039	0.891	0.887
CPD [40]	CVPR 2019	F	0.820	0.071	0.855	0.848	0.891	0.034	0.944	0.905	0.805	0.043	0.886	0.869
MINet [28]	CVPR 2020	F	0.829	0.064	0.857	0.856	0.909	0.029	0.953	0.919	0.828	0.037	0.898	0.884
VST [20]	ICCV 2021	F	0.829	0.061	0.844	0.872	0.900	0.029	0.953	0.928	0.818	0.037	0.892	0.896
SVF [47]	ICCV 2017	Un	0.695	0.131	0.789	0.758	-	-	-	-	-	-	-	-
C2S [18]	ECCV 2018	Un	0.754	0.087	0.838	0.826	0.839	0.051	0.919	0.873	0.710	0.066	0.841	0.817
WSS [34]	CVPR 2017	W	0.715	0.139	0.791	0.744	0.821	0.079	0.896	0.822	0.654	0.100	0.795	0.748
ASMO [15]	AAAI 2018	W	0.693	0.149	0.772	0.717	0.806	0.086	0.878	0.804	0.614	0.116	0.772	0.697
MWS [45]	CVPR 2019	W	0.713	0.133	0.790	0.768	0.814	0.084	0.895	0.818	0.684	0.091	0.814	0.759
WSSA [48]	CVPR 2020	W	0.774	0.092	0.837	0.797	0.860	0.047	0.927	0.865	0.742	0.062	0.857	0.804
MFNet [29]	ICCV 2021	W	0.746	0.112	0.818	0.782	0.839	0.058	0.917	0.852	0.692	0.079	0.830	0.778
SBBs [21]	TIP 2021	W	-	-	-	-	0.843	0.056	0.920	0.854	0.722	0.073	0.851	0.789
Ours	–	W	<b>0.794</b>	<b>0.084</b>	<b>0.840</b>	<b>0.808</b>	<b>0.870</b>	<b>0.044</b>	<b>0.934</b>	<b>0.871</b>	<b>0.770</b>	<b>0.057</b>	<b>0.873</b>	<b>0.816</b>

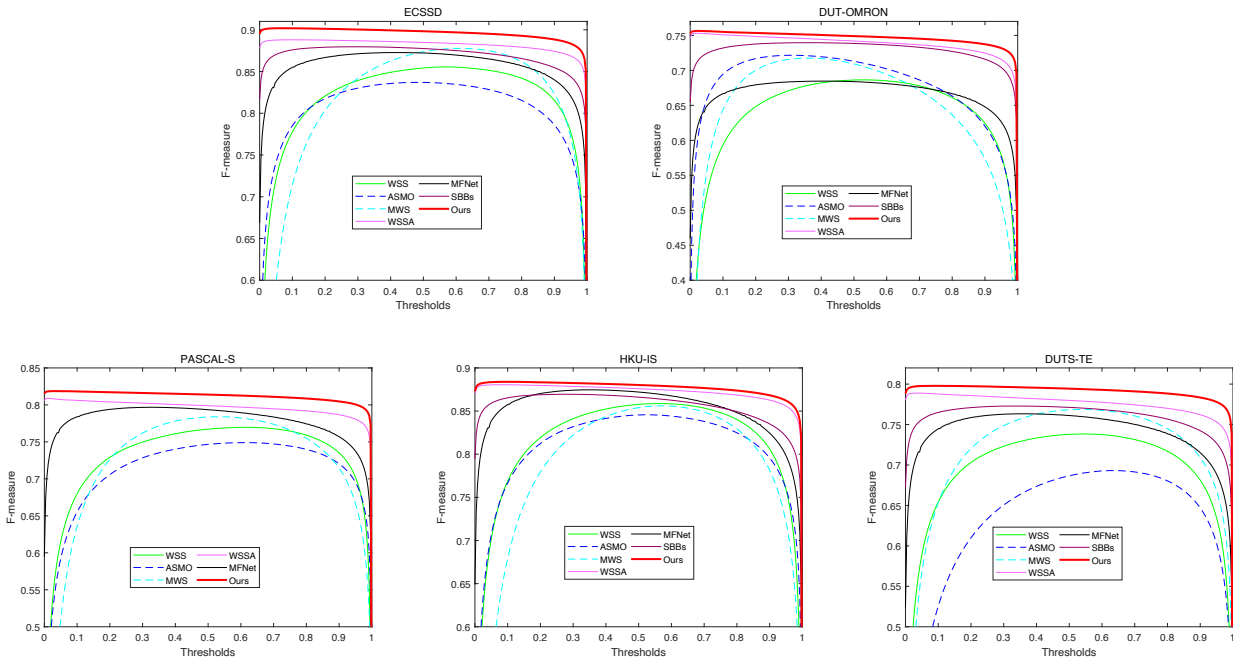
427 scores for each metric. From Fig. 10, we observe that  
 428 our method achieves the best performance in all four  
 429 average metrics.

430 **Qualitative comparison** We also show some ex-  
 431 amples of saliency maps from our proposed model and  
 432 other state-of-the-art methods using some challenging

433 cases in Figure 11. For example, we use instances of  
 434 large objects (1st and 2nd rows), multiple targets (3rd  
 435 and 6th rows), complex scenes (4th and 7th rows), s-  
 436 small objects (5th row), cluttered backgrounds (8th row)  
 437 and low contrast (9th row). Specifically, the 1st shows  
 438 a woman in an image and almost all methods are un-



**Fig. 8** PR curves of the proposed approach with other state-of-the-art methods using five datasets. Best viewed on screen.



**Fig. 9** F-measure curves of the proposed approach with other state-of-the-art methods using five datasets. Best viewed on screen.

439 able to detect accurate location on large objects except  
 440 our method. In the 3rd row, most methods can locate  
 441 the flowers while some details are lost. As we can see,  
 442 our approach is able to accurately find the salient ob-

ject with fewer false salient pixels detected. The 4th  
 443 row corresponds to a dog in a stadium. It is easy to  
 444 see our model segments the objects well, while other  
 445 models always detect the alphabet as salient objects.  
 446



**Fig. 10** (a) Comparison of quantitative results including average F-measure, average E-measure and average S-measure. Best viewed on screen. (b) Comparison of quantitative results including average MAE.

447 Compared with the 1st and 2nd row, a tower present-  
 448 ed in the 5th row is more difficult to segment thanks  
 449 to small objects and complex scenes. Nevertheless, our  
 450 CCFNet still highlights it very well. Different from the  
 451 last multiple examples, there are various characteris-  
 452 tics of salient objects in the 6th row of Fig. 11, such as  
 453 diverse colors and sizes. The proposed method can gener-  
 454 ate more reliable saliency maps in spite of the existing  
 455 little deficiency. Although our CCFNet erroneously seg-  
 456 ments the bottom part, it is still much better than other  
 457 methods. The 8th row shows the result of an objec-  
 458 t in cluttered backgrounds. Benefiting from EFM, our  
 459 model has more accurate edge details. Furthermore, the  
 460 9th row demonstrates that our model has good perfor-  
 461 mance with low contrast between the target and image  
 462 background. It can be observed that our model is able  
 463 to produce the complete structure of the cup whereas  
 464 previous work can not. In conclusion, our proposed ap-  
 465 proach performs better with respect to salient object  
 466 segmentation and localization, generating results that  
 467 are much closer to the ground truth in various challeng-  
 468 ing scenarios.

#### 469 4.4 Ablation studies

470 In this section, we perform a series of cases on ECSSD  
 471 and DUTS-TE datasets to assess the effectiveness of  
 472 our proposed method. All the ablation studies follow  
 473 the same implementation setup.

474 **Validity of different proposed module** We con-  
 475 duct various experiments to verify the effectiveness of  
 476 each component in CCFNet. In order to prove the valid-  
 477 ity of the proposed modules for saliency detection, we  
 478 compare our method with the other six schemes with  
 479 different proposed modules. Table 4 shows the perfor-  
 480 mance with seven schemes as well as their correspond-

481 ing saliency detection results. As seen from this table,  
 482 on the one hand, the quantitative scores of 1st ~ 3rd  
 483 lines are lower than 4th ~ 6th lines, that is both two  
 484 modules added is superior to single module, meanwhile,  
 485 the quantitative scores of 4th ~ 6th lines are lower than  
 486 7th line (ours), that is to say, three modules work to-  
 487 gether to realize the significant results. Especially, we  
 488 observe that these results perform more obviously on  
 489 complex datasets, such as DUTS-TE. Furthermore, Fig.  
 490 12 shows the results of average F-measure, average E-  
 491 measure, average S-measure and average MAE on EC-  
 492 SSD and DUTS-TE datasets. It also can be seen that  
 493 the last scheme achieves the best performance, i.e. the  
 494 scheme adopts three modules simultaneously.

495 **Validity of different loss functions** There are  
 496 three types of key loss function within the CCFNet, i.e.,  
 497 PCE loss, LSC loss and SSD loss. We design three ab-  
 498 lation experiments to evaluate the necessity of each loss  
 499 function, F-measure, MAE, E-measure and S-measure  
 500 scores are shown in Table 5. We find that LSC loss and  
 501 SSD loss can boost the performance of saliency maps  
 502 based on only using PCE loss. Especially, compared  
 503 with the first line (w/o LSC & SSD loss), our CCFNet  
 504 would promote the final performance with about 17.7%,  
 505 45.1%, 9.7% and 2.9% in F-measure, MAE, E-measure  
 506 and S-measure scores on ECSSD datasets, respective-  
 507 ly. In addition, Fig. 13 illustrates some samples of the  
 508 different loss functions, which can be seen that the pro-  
 509 posed CCFNet is well applicable to single target (1st  
 510 and 2nd lines) or multiple targets (3rd line).

511 **Validity of different parameters of loss func-**  
 512 **tion** Here we analyze the Validity of different paramet-  
 513 ers of the loss function in Table 6. Accordingly,  $\gamma_1 =$   
 514  $0, \gamma_2 = 0, \gamma_3 = 0$  means that our network has no ex-  
 515 tra supervision except dominant loss function  $L_{dom}$ . It  
 516 can be seen that it has the lowest scores compared with

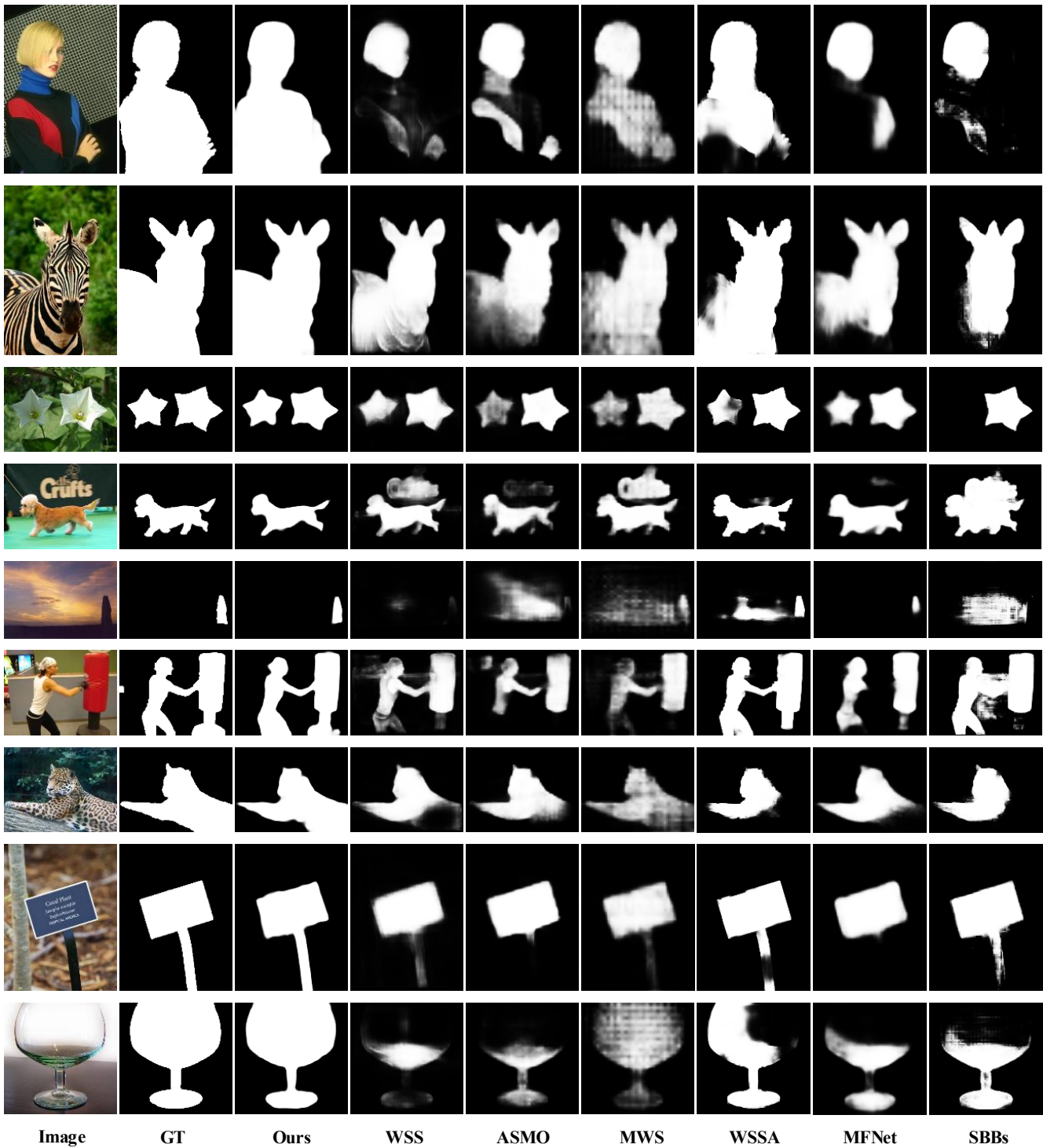


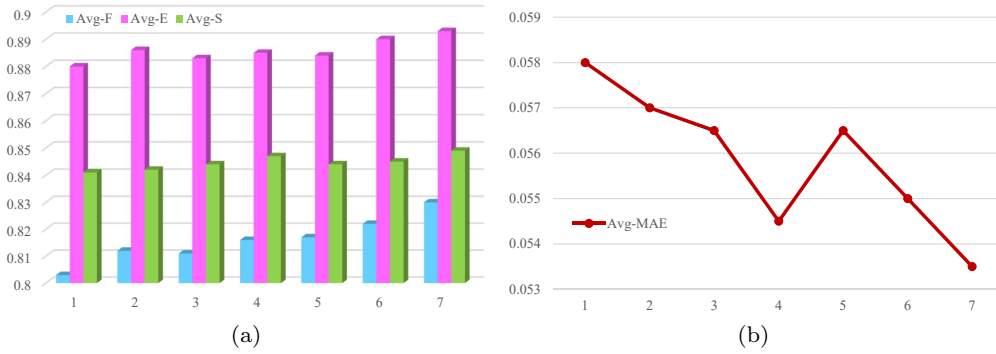
Fig. 11 Visual comparison between the proposed model and state-of-the-art methods.

517 other schemes. That means sub-stage loss is beneficial  
 518 to the network. Similarly,  $\gamma_1 = 1, \gamma_2 = 1, \gamma_3 = 1$  means  
 519 that sub-stage loss has the same weight with the dominant  
 520 loss function  $L_{dom}$ . This is not the best result  
 521 for the scribble saliency detection network, which may  
 522 be caused by the sub-stage bringing more negative in-

formation. Just because there is more and more useful  
 information from the first sub-stage to the third sub-  
 stage, hence, we balance the weight coefficients in a  
 progressive manner. It is proved that our approach is ef-  
 fective and reliable. Note that  $\gamma_1 = 0.8, \gamma_2 = 0, \gamma_3 = 0.4$   
 means that there is no first sub-stage loss function,

523  
 524  
 525  
 526  
 527  
 528





**Fig. 12** (a) Comparison of quantitative results including average F-measure, average E-measure and average S-measure. Best viewed on screen. (b) Comparison of quantitative results including average MAE.

**Table 4** Ablation study for our proposed different modules.

	FCM	EFM	GCCO	ECSSD				DUTS-TE			
				$F_\beta \uparrow$	MAE $\downarrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$	$F_\beta \uparrow$	MAE $\downarrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$
1			✓	0.869	0.054	0.905	0.874	0.737	0.062	0.854	0.807
2	✓			0.874	0.054	0.907	0.874	0.749	0.060	0.864	0.810
3		✓		0.874	0.052	0.907	0.877	0.747	0.061	0.859	0.811
4		✓	✓	0.878	0.051	0.905	0.879	0.753	0.058	0.864	0.815
5	✓		✓	0.880	0.053	0.904	0.876	0.754	0.060	0.864	0.812
6	✓	✓		0.881	0.052	0.910	0.877	0.762	0.058	0.870	0.813
7	✓	✓	✓	<b>0.890</b>	<b>0.050</b>	<b>0.912</b>	<b>0.882</b>	<b>0.770</b>	<b>0.057</b>	<b>0.873</b>	<b>0.816</b>

**Table 5** Ablation study for different loss functions.

	PCE	LSC	SSD	ECSSD				DUTS-TE			
				$F_\beta \uparrow$	MAE $\downarrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$	$F_\beta \uparrow$	MAE $\downarrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$
1	✓			0.756	0.091	0.831	0.799	0.575	0.105	0.725	0.711
2	✓		✓	0.763	0.087	0.839	0.809	0.589	0.095	0.742	0.727
3	✓	✓		0.875	0.053	0.900	0.875	0.751	0.052	0.857	0.812
4	✓	✓	✓	<b>0.890</b>	<b>0.050</b>	<b>0.912</b>	<b>0.882</b>	<b>0.770</b>	<b>0.057</b>	<b>0.873</b>	<b>0.816</b>

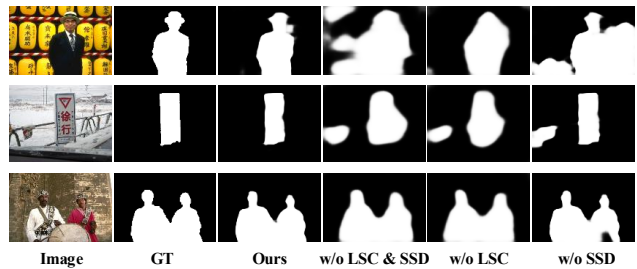
**Table 6** Ablation study for different parameters of loss function.

		ECSSD				DUTS-TE			
		$F_\beta \uparrow$	MAE $\downarrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$	$F_\beta \uparrow$	MAE $\downarrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$
1	$\gamma_1 = 0, \gamma_2 = 0, \gamma_3 = 0$	0.866	0.055	0.896	0.872	0.732	0.064	0.852	0.804
2	$\gamma_1 = 0, \gamma_2 = 0.6, \gamma_3 = 0.4$	0.869	0.054	0.901	0.874	0.732	0.064	0.848	0.807
3	$\gamma_1 = 0.8, \gamma_2 = 0, \gamma_3 = 0.4$	0.878	0.051	0.911	0.878	0.758	0.057	0.873	0.815
4	$\gamma_1 = 0.8, \gamma_2 = 0.6, \gamma_3 = 0$	0.874	0.051	0.906	0.879	0.734	0.065	0.848	0.807
5	$\gamma_1 = 1, \gamma_2 = 1, \gamma_3 = 1$	0.875	0.052	0.900	0.876	0.747	0.061	0.856	0.811
6	$\gamma_1 = 0.8, \gamma_2 = 0.6, \gamma_3 = 0.4$	<b>0.890</b>	<b>0.050</b>	<b>0.912</b>	<b>0.882</b>	<b>0.770</b>	<b>0.057</b>	<b>0.873</b>	<b>0.816</b>

529 which can be seen that the result is poor on two dataset-  
 530 s. The reason for this issue may be that the first sub-  
 531 stage supervision has more semantics, which plays a  
 532 decisive role in the subsequent prediction.

## 533 5 Conclusion

534 In this paper, we proposed a novel and effective comple-  
 535 mentary characteristics fusion network (CCFNet) for  
 536 salient object detection with scribble annotations. First,  
 537 a global context guiding operation and edge fusion mod-



**Fig. 13** Ablation study of different loss function.

ule are introduced, which are used to obtain global semantics and learn salient edge information. It is proved that they can better understand global high level information and edge information. Next, to exploit complete salient regions with different level features, this paper proposes the feature correlation module for saliency detection. In order to better distinguish foreground and background information for a given image, self-supervised saliency detection loss is illustrated. Finally, to demonstrate the performance of our proposed method, we conduct experiment results on five well-known datasets. Extensive experimental results demonstrate that our approach outperforms state-of-the-art weakly supervised methods and ablation studies prove the effectiveness of each component as well. Future work will focus on developing a more lightweight weakly supervised model as well as investigating how to deploy SOD algorithms in mobile devices to strengthen practicality.

## Declaration of Interests

We declare that we have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgements

This work was supported by National Natural Science Foundation of China (No. 61973066,61471110) , Major Science and Technology Projects of Liaoning Province (No.2021JH1/10400049), Fundation of Key Laboratory of Aerospace System Simulation(No.6142002200301), Open Research Projects of Zhejiang Lab (No.2019KD0A-D01/006) and Major Science and technology innovation engineering projects of Shandong Province (No.2019JZ-ZY010128).

## References

1. Araslanov, N., Roth, S.: Single-stage semantic segmentation from image labels. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4253–4262 (2020)
2. Borji, A., Itti, L.: Exploiting local and global patch rarities for saliency detection. In: 2012 IEEE conference on computer vision and pattern recognition, pp. 478–485. IEEE (2012)
3. Chen, S., Tan, X., Wang, B., Hu, X.: Reverse attention for salient object detection. In: Proceedings of the European Conference on Computer Vision (ECCV), pp. 234–250 (2018)
4. Cheng, M.M., Mitra, N.J., Huang, X., Torr, P.H., Hu, S.M.: Global contrast based salient region detection. IEEE transactions on pattern analysis and machine intelligence **37**(3), 569–582 (2014)
5. Cheng, M.M., Zhang, F.L., Mitra, N.J., Huang, X., Hu, S.M.: Repfinder: finding approximately repeated scene elements for image editing. ACM Transactions on Graphics (TOG) **29**(4), 1–8 (2010)
6. Craye, C., Filliat, D., Goudou, J.F.: Environment exploration for object-based visual saliency learning. In: 2016 IEEE International Conference on Robotics and Automation (ICRA), pp. 2303–2309. IEEE (2016)
7. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. International journal of computer vision **88**(2), 303–338 (2010)
8. Fan, D.P., Cheng, M.M., Liu, Y., Li, T., Borji, A.: Structure-measure: A new way to evaluate foreground maps. In: Proceedings of the IEEE international conference on computer vision, pp. 4548–4557 (2017)
9. Fan, D.P., Gong, C., Cao, Y., Ren, B., Cheng, M.M., Borji, A.: Enhanced-alignment measure for binary foreground map evaluation. arXiv preprint arXiv:1805.10421 (2018)
10. Goferman, S., Zelnik-Manor, L., Tal, A.: Context-aware saliency detection. IEEE transactions on pattern analysis and machine intelligence **34**(10), 1915–1926 (2011)
11. Hong, S., You, T., Kwak, S., Han, B.: Online tracking by learning discriminative saliency map with convolutional neural network. In: International conference on machine learning, pp. 597–606. PMLR (2015)
12. Hou, Q., Cheng, M.M., Hu, X., Borji, A., Tu, Z., Torr, P.: Deeply supervised salient object detection with short connections. IEEE Transactions on Pattern Analysis and Machine Intelligence **41**(4), 815–828 (2019)
13. Hou, Q., Cheng, M.M., Hu, X., Borji, A., Tu, Z., Torr, P.H.: Deeply supervised salient object detection with short connections. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 3203–3212 (2017)
14. Jiang, H., Wang, J., Yuan, Z., Wu, Y., Zheng, N., Li, S.: Salient object detection: A discriminative regional feature integration approach. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2083–2090 (2013)
15. Li, G., Xie, Y., Lin, L.: Weakly supervised salient object detection using image labels. In: Thirty-second AAAI conference on artificial intelligence (2018)
16. Li, G., Yu, Y.: Visual saliency based on multiscale deep features. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 5455–5463 (2015)
17. Li, G., Yu, Y.: Visual saliency detection based on multiscale deep cnn features. IEEE transactions on image processing **25**(11), 5012–5024 (2016)
18. Li, X., Yang, F., Cheng, H., Liu, W., Shen, D.: Contour knowledge transfer for salient object detection. In: Proceedings of the European Conference on Computer Vision (ECCV), pp. 355–370 (2018)
19. Liu, J.J., Hou, Q., Cheng, M.M., Feng, J., Jiang, J.: A simple pooling-based design for real-time salient object detection. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3917–3926 (2019)
20. Liu, N., Zhang, N., Wan, K., Han, J., Shao, L.: Visual saliency transformer. arXiv preprint arXiv:2104.12099 (2021)
21. Liu, Y., Wang, P., Cao, Y., Liang, Z., Lau, R.W.: Weakly-supervised salient object detection with saliency bounding boxes. IEEE Transactions on Image Processing **30**, 4423–4435 (2021)

22. Liu, Y., Zhang, Y., Liu, S., Coleman, S., Wang, Z., Qiu, F.: Salient object detection by aggregating contextual information. *Pattern Recognition Letters* **153**, 190–199 (2022)
23. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3431–3440 (2015)
24. Luo, Z., Mishra, A., Achkar, A., Eichel, J., Li, S., Jodoin, P.M.: Non-local deep features for salient object detection. In: *Proceedings of the IEEE Conference on computer vision and pattern recognition*, pp. 6609–6617 (2017)
25. Nguyen, D.T., Dax, M., Mummadi, C.K., Ngo, T.P.N., Nguyen, T.H.P., Lou, Z., Brox, T.: Deepusps: Deep robust unsupervised saliency prediction with self-supervision. *arXiv preprint arXiv:1909.13055* (2019)
26. Noori, M., Mohammadi, S., Majelan, S.G., Bahri, A., Havaei, M.: Dfnet: Discriminative feature extraction and integration network for salient object detection. *Engineering Applications of Artificial Intelligence* **89**, 103419 (2020)
27. Pan, Z., Jiang, P., Tu, C.: Scribble-supervised semantic segmentation by random walk on neural representation and self-supervision on neural eigenspace. *arXiv preprint arXiv:2011.05621* (2020)
28. Pang, Y., Zhao, X., Zhang, L., Lu, H.: Multi-scale interactive network for salient object detection. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9413–9422 (2020)
29. Piao, Y., Wang, J., Zhang, M., Lu, H.: Mfnet: Multi-filter directive network for weakly supervised salient object detection. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4136–4145 (2021)
30. Qin, X., Zhang, Z., Huang, C., Dehghan, M., Zaiane, O.R., Jagersand, M.: U2-net: Going deeper with nested u-structure for salient object detection. *Pattern Recognition* **106**, 107404 (2020)
31. Siva, P., Russell, C., Xiang, T., Agapito, L.: Looking beyond the image: Unsupervised learning for object saliency and detection. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3238–3245 (2013)
32. Tang, M., Djelouah, A., Perazzi, F., Boykov, Y., Schroers, C.: Normalized cut loss for weakly-supervised cnn segmentation. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1818–1827 (2018)
33. Ullah, I., Jian, M., Hussain, S., Lian, L., Ali, Z., Qureshi, I., Guo, J., Yin, Y.: Global context-aware multi-scale features aggregative network for salient object detection. *Neurocomputing* **455**, 139–153 (2021)
34. Wang, L., Lu, H., Wang, Y., Feng, M., Wang, D., Yin, B., Ruan, X.: Learning to detect salient objects with image-level supervision. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 136–145 (2017)
35. Wang, T., Borji, A., Zhang, L., Zhang, P., Lu, H.: A stagewise refinement model for detecting salient objects in images. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4019–4028 (2017)
36. Wang, W., Shen, J., Dong, X., Borji, A.: Salient object detection driven by fixation prediction. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1711–1720 (2018)
37. Wei, J., Wang, S., Huang, Q.: F<sup>3</sup>net: Fusion, feedback and focus for salient object detection. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 12321–12328 (2020)
38. Wei, J., Wang, S., Wu, Z., Su, C., Huang, Q., Tian, Q.: Label decoupling framework for salient object detection. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13025–13034 (2020)
39. Wei, Y., Liang, X., Chen, Y., Shen, X., Cheng, M.M., Feng, J., Zhao, Y., Yan, S.: Stc: A simple to complex framework for weakly-supervised semantic segmentation. *IEEE transactions on pattern analysis and machine intelligence* **39**(11), 2314–2320 (2016)
40. Wu, Z., Su, L., Huang, Q.: Cascaded partial decoder for fast and accurate salient object detection. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3907–3916 (2019)
41. Wu, Z., Su, L., Huang, Q.: Stacked cross refinement network for edge-aware salient object detection. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7264–7273 (2019)
42. Yan, Q., Xu, L., Shi, J., Jia, J.: Hierarchical saliency detection. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1155–1162 (2013)
43. Yang, C., Zhang, L., Lu, H., Ruan, X., Yang, M.H.: Saliency detection via graph-based manifold ranking. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3166–3173 (2013)
44. Yu, S., Zhang, B., Xiao, J., Lim, E.G.: Structure-consistent weakly supervised salient object detection with local saliency coherence. In: *AAAI Conf. Art. Intell* (2021)
45. Zeng, Y., Zhuge, Y., Lu, H., Zhang, L., Qian, M., Yu, Y.: Multi-source weak supervision for saliency detection. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6074–6083 (2019)
46. Zhai, Y., Shah, M.: Visual attention detection in video sequences using spatiotemporal cues. In: *Proceedings of the 14th ACM international conference on Multimedia*, pp. 815–824 (2006)
47. Zhang, D., Han, J., Zhang, Y.: Supervision by fusion: Towards unsupervised learning of deep salient object detector. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4048–4056 (2017)
48. Zhang, J., Yu, X., Li, A., Song, P., Liu, B., Dai, Y.: Weakly-supervised salient object detection via scribble annotations. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12546–12555 (2020)
49. Zhang, J., Zhang, T., Dai, Y., Harandi, M., Hartley, R.: Deep unsupervised saliency detection: A multiple noisy labeling perspective. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9029–9038 (2018)
50. Zhang, L., Dai, J., Lu, H., He, Y., Wang, G.: A bi-directional message passing model for salient object detection. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1741–1750 (2018)
51. Zhang, P., Wang, D., Lu, H., Wang, H., Ruan, X.: A-mulet: Aggregating multi-level convolutional features for salient object detection. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 202–211 (2017)
52. Zhang, P., Wang, D., Lu, H., Wang, H., Yin, B.: Learning uncertain convolutional features for accurate saliency detection. In: *Proceedings of the IEEE International Conference on computer vision*, pp. 212–221 (2017)

- 
- 788 53. Zhang, P., Zhuo, T., Huang, W., Chen, K., Kankanhalli,  
789 M.: Online object tracking based on cnn with spatial-  
790 temporal saliency guided sampling. *Neurocomputing*  
791 **257**, 115–127 (2017)
- 792 54. Zhang, X., Wang, T., Qi, J., Lu, H., Wang, G.: Progres-  
793 sive attention guided recurrent network for salient object  
794 detection. In: *Proceedings of the IEEE Conference on*  
795 *Computer Vision and Pattern Recognition*, pp. 714–722  
796 (2018)
- 797 55. Zhao, J.X., Liu, J.J., Fan, D.P., Cao, Y., Yang, J., Cheng,  
798 M.M.: Egnet: Edge guidance network for salient object  
799 detection. In: *Proceedings of the IEEE/CVF Interna-*  
800 *tional Conference on Computer Vision*, pp. 8779–8788  
801 (2019)
- 802 56. Zhao, T., Wu, X.: Pyramid feature attention network for  
803 saliency detection. In: *Proceedings of the IEEE/CVF*  
804 *Conference on Computer Vision and Pattern Recognition*  
805 (2019)

**Declaration of interests**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: