

# Weakly-Supervised Saliency Detection via Salient Object Subitizing

Xiaoyang Zheng\*, Xin Tan\*, Jie Zhou, Lizhuang Ma<sup>†</sup>, and Rynson W.H. Lau<sup>†</sup>

**Abstract**—Salient object detection aims at detecting the most visually distinct objects and producing the corresponding masks. As the cost of pixel-level annotations is high, image tags are usually used as weak supervisions. However, an image tag can only be used to annotate one class of objects. In this paper, we introduce saliency subitizing as the weak supervision since it is class-agnostic. This allows the supervision to be aligned with the property of saliency detection, where the salient objects of an image could be from more than one class. To this end, we propose a model with two modules, Saliency Subitizing Module (SSM) and Saliency Updating Module (SUM). While SSM learns to generate the initial saliency masks using the subitizing information, without the need for any unsupervised methods or some random seeds, SUM helps iteratively refine the generated saliency masks. We conduct extensive experiments on five benchmark datasets. The experimental results show that our method outperforms other weakly-supervised methods and even performs comparable to some fully-supervised methods.

**Index Terms**—weak supervision, saliency detection, object subitizing

## I. INTRODUCTION

The salient object detection task aims at accurately recognizing the most distinct objects in a given image that would attract human attention. This task has received a lot of research interests in recent years, as it plays an important role in many other computer vision tasks, such as visual tracking [1], image editing/manipulation [2], [3] and image retrieval [4]. Recently, deep convolutional neural networks have achieved significant progress in saliency detection [5]–[9]. However, most of these recent methods are primarily CNN-based, which rely on a large amount of pixel-wised annotations for training. For such an image segmentation task, it is both arduous and inefficient to collect a large amount of pixel-wised saliency masks. For example, it usually takes several minutes for an experienced worker to annotate one single image. This drawback confines the amount of available training samples. In this paper, we focus on the salient object detection task with a weakly-supervised setting.

Some methods [12]–[14] tried to address salient object detection with image-tag supervisions. Li et al. [12] utilized

\* Equal Contribution.

<sup>†</sup> Corresponding Author.

X. Zheng, X. Tan, J. Zhou and L. Ma are with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, 200240, China. X. Tan and J. Zhou are also with the Department of Computer Science, City University of Hong Kong, HKSAR, China. E-mail: zxyang.reg@sjtu.edu.cn, tanxin2017@sjtu.edu.cn, lord\_liang@sjtu.edu.cn, ma-lz@cs.sjtu.edu.cn.

Rynson W.H. Lau is with the Department of Computer Science, City University of Hong Kong, HKSAR, China. E-mail: rynson.lau@cityu.edu.hk.

Manuscript received Sept. 20, 2020; revised Dec. 04, 2020.



Fig. 1. Several inconsistent cases between the given image labels and the actual salient objects. These images and tags are chosen from the Pascal VOC [10] and Microsoft COCO [11] datasets. The captions under images are the given labels. The masks are generated with our methods, which show the actually salient objects.

Class Activation Maps as coarse saliency maps. Together with results from unsupervised methods, those masks are used to supervise the segmentation stream. Wang et al [13] proposed a two-stage method, which assigns category tags to object regions and fine-tunes the network with the predicted saliency maps with the ground truth. Zeng et al. [14] proposed a unified framework to conduct saliency detection with diverse weak supervisions, including image tags, captions and pseudo labels. They achieved good performances with image labels from the Microsoft COCO [11] or Pascal VOC [10] datasets. However, their results are established on a critical assumption that *the labelled object is the most visually distinct one*. From those datasets with image tags, We observe that *this assumption is not always reliable*. As shown in Figure 1, the actual salient objects are inconsistent with the image labels. For example, the image in the second column is labelled as “fire hydrant”, it is obvious that the orange “ball” should also be a salient object. In addition, even trained on datasets with multi-class labels, these methods essentially detect object within the categories, but not salient objects. Hence, image category labels do not guarantee the property of saliency.

**Motivation.** Subitizing is the rapid enumeration of a small number of items in a given image, regardless of their semantic category. According to [15], subitizing of up to four targets is highly accurate, quick and confident. In addition, since the subitizing information may contain objects from different categories, it is class-agnostic. Inspired by the above advantages, we propose to address the saliency detection problem using only the object subitizing information as a weak supervision.

Although there exist works, e.g., [16], that use subitizing

as an auxiliary supervision, we propose to apply subitizing as the weak supervision in this work. To the best of our knowledge, we are the first to adopt subitizing as the only supervision in saliency detection task. However, the subitizing information does not indicate the position and appearance of salient objects. Therefore, we propose the Saliency Subitizing Module (SSM) to produce saliency masks. Recent works [17], [18] have proven that, even trained with classification tasks, CNNs implicitly reveal the attention regions in the given images. Trained on subitizing information, the SSM relies on the distinct regions to conduct classification. By extracting those regions, we can explicitly obtain the locations of the salient objects.

However, as pointed out in [19], in a well trained classification network, the discriminative power of each object part is different from each other, and thus lead to incomplete segmentation. In the finetune stage, we need to further enlarge the prominent regions extracted from the network. Kolesnikov et al. [20] trained their network with pseudo labels for multiple iterations and obtain integrated results, while the multi-stage training is complicated. In order to address this issue, we design the Saliency Updating Module (SUM) for refining the saliency masks produced by SSM. In each iteration, the generated saliency maps, combined with original images, are used to generate masked images. With those masked images as input to the next iteration, the network learns to recognize those related but less salient regions. During the inference phase, given an image, our model will produce the saliency maps without any iterations, and there will be no need to provide the subitizing information. Our extensive evaluations demonstrate the superiority of the proposed methods over the state-of-the-art weakly-supervised methods.

In summary, this paper has the following contributions:

- We propose to use subitizing as a weak supervision in the saliency detection task, which has the advantage of being class-agnostic.
- We propose an end-to-end multi-task saliency detection network. By introducing subitizing information, our network first generates rough saliency masks with the Saliency Subitizing Module (SSM), and then iteratively refines the saliency masks with the Saliency Updating Module (SUM).
- Our extensive experiments show that the proposed method achieves superior performance on five popular salient datasets, compared with other weak-supervised methods. It even performs comparable to some of the fully-supervised methods.

## II. RELATED WORK

Recently, the progress on salient object detection is substantial, benefiting by the development of deep neural networks. He et al. [5] proposed a convolution neural network based on super-pixel for saliency detection. Li et al. [21] utilized multi-scale features extracted from a deep CNN. Zhao et al. [22] proposed a multi-context deep learning framework for detecting salient objects with two different CNNs used to learn global and local context information. Yuan et al. [23]

proposed a saliency detection framework, which extracted the correlations among macro object contours, RGB features and low-level image features. Wang et al. [24] proposed a pyramid attention structure to enlarge the receptive field. Hou et al. [25] introduced short connections to an edge detector. Zhu et al. [8] proposed the attentional DenseASPP to exploit local and global contexts at multiple dilated convolutional layers. Hu et al. [9] proposed a spatial attenuation context network, which recurrently translated and aggregated the context features in different layers. Tu et al. [26] presented an edge-guided block to embed boundary information into saliency maps. Zhou et al. [27] proposed a multi-type self-attention network to learn more semantic details from degraded images. However, these methods rely heavily on pixel-wised supervisions. Due to the scarcity of pixel-wised data, we focus on the weakly-supervised saliency detection task.

**Weakly-Supervised Saliency detection.** There are many works using weak supervisions for the saliency detection task. For example, Li et al. [12] used the image-level labels to train the classification network and applied the coarse activation maps as saliency maps. Wang et al. [13] proposed a two-stage weakly-supervised method by designing a Foreground Inference Network (FIN) to predict foreground regions and a Global Smooth Pooling (GSP) to aggregate responses of predicted objects. Zeng et al. [14] proposed a unified network, which is trained on multi-source weak supervisions, including image tags, captions and pseudo labels. They designed an attention transfer loss to transmit signals between sub-networks with different supervisions. However, as discussed in Section I, the image-level supervisions are not always reliable. In addition, captions were used as a weak supervision in [14], combined with other supervisions. Different from those methods above, we propose to use subitizing information as the weak supervision in the saliency detection task.

**Salient object subitizing.** Zhang et al. [28] proposed a salient object subitizing dataset SOS. They firstly studied the problem of salient object subitizing and revealed the relations between subitizing and saliency detection. Lu et al. [29] formulated the subitizing task as a matching problem and exploited the self-similarity within the same class. He et al. [16] trained a subitizing network to provide additional knowledge to the pixel-level segmentation stream. Recently, Amirul et al. [30] proposed a salient object subitizing network and recognized the variability of subitizing. They also provided outputs as a distribution that reflects this variability. In this paper, our approach is motivated by these methods but we use subitizing as the weak supervision.

**Map refinement.** Li et al. [12] adopted saliency maps generated by some unsupervised methods as the initial seeds. With a graphical model, these saliency maps are used as pixel-level annotations and refined in an iterative way. However, in our proposed method, we do not utilize any unsupervised methods or initial seeds. The saliency maps are refined iteratively from the activation maps of the subitizing network. Li et al. [31] adopted a soft mask loss to an auxiliary attention stream. However, the input of [31] is updated only once while the inputs to our network are iteratively updated. In addition, there are some existing post-processing techniques used to

refine the saliency masks. In [12], [13], the authors utilized an iterative conditional random field (CRF) to enforce spatial label consistency. Zheng et al. [32] further proposed to conduct approximate inference with pair-wised Gaussian potential in CRF as a recurrent neural network. Chen et al. [33] employed the relations of the deep features to promote saliency detectors in a self-supervised way. In order to achieve better results, we adopt a refinement process, which maintains the internal structure of original images and enforces smoothness into the final saliency maps.

### III. SALIENT OBJECT DETECTION METHOD

We propose a multi-task convolutional neural network, which consists of two main modules: saliency subitizing module (SSM) and saliency updating module (SUM). SSM helps learn counting of salient objects and extract coarse saliency maps with the precise locations of the target objects. SUM helps update saliency masks produced by SSM, and extend the activation regions. Finally, we apply a refinement process to refine the object boundaries. The pipeline of our method is presented in Figure 2.

#### A. Saliency Subitizing Module

The subitizing information indicates the number of salient objects in a given image. It does not provide the location or the appearance information of the salient objects explicitly. However, when we train our network with the subitizing supervisions, the network will learn to focus on regions related to the salient objects. Hence, we design the Saliency Subitizing Module (SSM) to extract the attention regions as coarse saliency masks. We regard the saliency object subitizing task as a classification task. Training images are divided into 5 categories based on the number of salient objects: 0, 1, 2, 3 and 4+. We use ResNet-50 [34] as the backbone network, pretrained on the ImageNet [35] dataset. The original 1000-way fully-connected layer are replaced by a 5-way fully-connected layer. We use cross-entropy loss as the classification loss  $L_{cls}$ . In order to obtain denser saliency maps, the stride of the last two down-sampling layers is set as 1 in our backbone network, which produces feature maps with 1/8 of the original resolution before the classification layer. In order to enhance the representation power of the proposed network, we also apply two attention modules: channel attention module and spatial attention module, which tells the network “where” and “what” to focus, respectively. Both of them are placed in a sequential way between the ResNet blocks.

We apply the technique of Grad-CAM [18] to extract activation regions as the initial saliency maps, which contains the gradient information flowing into the last convolutional layers during the backward phase. The gradient information represents the importance of each neuron during the inference of the network. We assume that the produced features from the last convolutional layer have a channel size of  $K$ . For a given image, let  $f_k$  be the activation of unit  $k$ , where  $k \in [1, K]$ . For each class  $c$ , the gradients of the score  $y^c$  with respect

to the activation map  $f_k$  are averaged to obtain the neuron importance weight  $\alpha_k^c$  of class  $c$ :

$$\alpha_k^c = \frac{1}{N} \sum_i^m \sum_j^h \frac{\partial y^c}{\partial f_{ij}^k}, \quad (1)$$

where  $i$  and  $j$  represent the coordinates in the feature map and  $N = m \times h$ . With the neuron importance weight  $\alpha_k^c$ , we can compute the activation map  $M^c$ :

$$M^c = \text{ReLU}\left(\sum_k \alpha_k^c f^k\right). \quad (2)$$

Note that the ReLU function filters the negative gradient values, since only the positive ones contribute to the class decision, while the negative values contribute to other categories. The size of the saliency map is the same as the size of the last convolutional feature maps (1/8 of the original resolution). Since the saliency maps  $M^c$  are obtained within each inference, they become trainable during the training stage.

#### B. Saliency Updating Module

In the Saliency Subitizing Module, our proposed network learns to detect the regions that contribute to the counting of salient objects. Due to this attribute, with only the SSM module, we can obtain saliency masks with accurate locations of the target objects. However, the quality of the saliency masks may not be very high. In order to address this issue, we design a Saliency Updating Module (SUM) to fine-tune the obtained masks, with the aim to refine the activation regions. Li et al. [12] updated saliency mask with an additional CRF module. In contrast, our proposed model refines the saliency masks in an end-to-end way.

As shown in Figure 2, we fuse the origin images and the saliency maps to obtain masked images as new inputs to the next iteration. Visually, the current prominent area is eliminated from the original samples. We define  $I_0$  as the original images.  $I_i^c$  denotes the input images at the  $i$ -th iteration ( $i \geq 1$ ).  $M_i^c$  denotes the saliency maps of class  $c$  at the  $i$ -th iteration.

The fusion operation is formulated as:

$$I_i^c = I_0 - (\text{Sigmoid}(\omega \cdot (M_{i-1}^c - \sigma)) \odot I_0), \quad (3)$$

where  $\odot$  stands for element-wise multiplication,  $\sigma$  is a threshold matrix with all elements equal to  $\sigma$ . With the scale parameter  $\omega$ , the mask term gets closer to 1 when  $M_{i-1}^c > \sigma$ , and gets closer to 0 otherwise. As presented in Eq. 3, we enforce  $I_i^c$  to contain as few features from the target class  $c$  as possible.

Trained on samples without features from the current prominent area, the network learns to recognize those related but less salient regions. In other words, regions beyond the high-responding area should also include features that trigger the network to recognize the sample as class  $c$ . Similar to [31], we introduce a mask mining loss  $L_{mask}$  to extract larger activation area. This loss penalizes the prediction error for class  $c$ , as shown below,

$$L_{mask} = \frac{1}{n} \sum_c y^c (I^c), \quad (4)$$

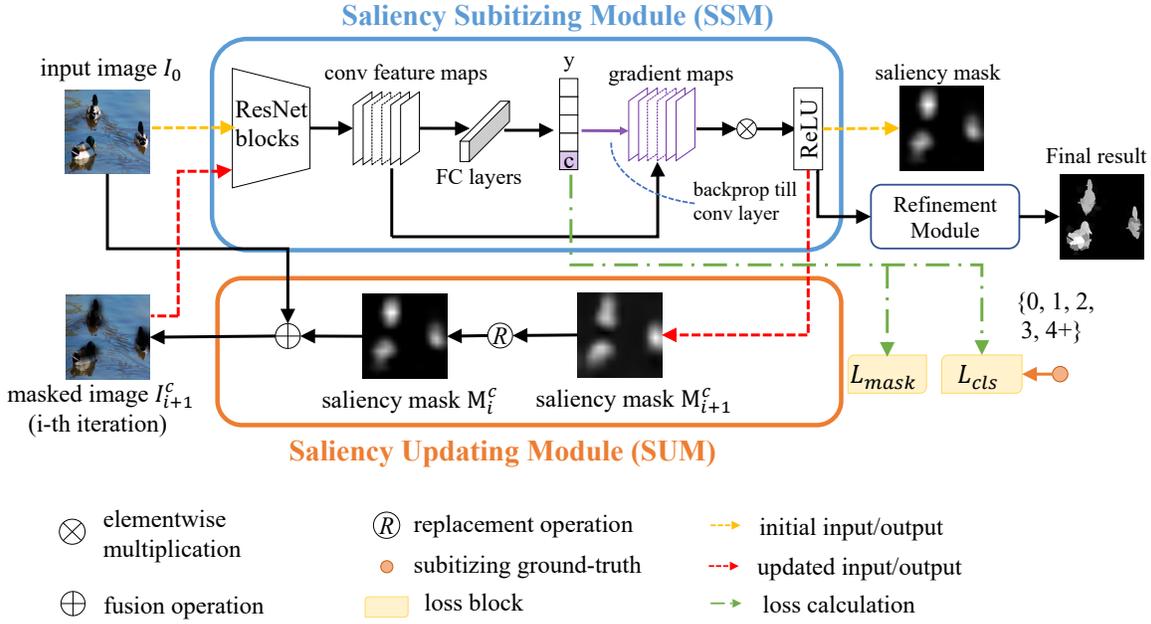


Fig. 2. The pipeline of the proposed network, with the Saliency Subitizing Module (SSM), the Saliency Updating Module (SUM) and the refinement process.

where  $n$  is the dataset size and  $y^c(I^c)$  represents the prediction score of masked images  $I^c$  for class  $c$ . With the loss perspective, the prediction scores of the right label for the masked images should be lower than those for the original images. During the training phase, the total loss  $L$  is the combination of the classification loss  $L_{cls}$  and the mask mining loss  $L_{mask}$ .

$$L = L_{cls} + \alpha L_{mask}, \quad (5)$$

where  $\alpha$  is the weighting parameter. We set  $\alpha = 1$  in all our experiments. With this loss, the network learns to extract those less salient but related parts of the target objects, while maintains the ability of recognizing subitizing information. In [31], the extracted regions for masking the input are updated only once. These regions extracted from just a single step are usually small and sparse, since CNNs tend to capture the most discriminative regions and neglect the others in the image. In contrast, our method updates the extracted regions through multiple iterations. In this way, the extracted regions in our method are more integrated.

**Discussion.** Although Wei et al. [36] also adopted an iterative strategy, our method is different from [36] in two main aspects. First, during the generation of training images for the next iteration, [36] simply replaced the internal pixels by the mean pixel values of all the training images. Instead, we use a threshold  $\sigma$  to determine the salient regions and a weighting parameter  $\omega$  to adjust the removing rate of features (as presented in Eq. 3), so that the correlations of the extracted regions and the backgrounds at different iterations would be smoothly changed. Second, [36] took the mined object region as the pixel-level label for segmentation training. Instead, our method is only trained on the given dataset with subitizing labels, avoiding training on unreliable pseudo labels.

### C. Refinement Process

To refine the object boundaries in the saliency maps, we take a graph-based optimization method. Inspired by [37], we adopt super-pixels produced with SLIC [38] as the basic representation units. Those super-pixels are organized as an adjacency graph to model the structure in both spatial and feature dimensions. A given image is segmented into a set of super-pixels  $\{x_i\}_{i=1}^N$ , where  $N = 200$ . The super-pixel graph  $A = (a_{i,j})_{N \times N}$  is defined as follows:

$$a_{i,j} = \begin{cases} K(x_i, x_j), & \text{if } x_i \text{ and } x_j \text{ are spatially adjacent;} \\ 0, & \text{otherwise,} \end{cases} \quad (6)$$

where  $K(x_i, x_j) = \exp(-\frac{1}{2}\|x_i - x_j\|_2^2)$  evaluates the feature similarity. Assume that there exist  $l$  super-pixels with initial scores. Our task is to learn a non-linear regression function  $g(x) = \sum_{j=1}^N \alpha_j K(x_i, x_j)$  for each super-pixel  $x$ . The framework is shown as:

$$\min_g \frac{1}{l} \sum_{i=1}^l (y_i - g(x_i))^2 + \theta_1 \|g\|_K^2 + \frac{\theta_2}{N^2} g^T D^{-\frac{1}{2}} L D^{-\frac{1}{2}} g, \quad (7)$$

where  $y_i$  is the weight of the  $i$ -th unit in the super-pixel graph, and  $g = (g(x_1), \dots, g(x_N))^T$ .  $\|g\|_K$  denotes the norm of  $g$  induced by the function  $K$ ;  $D$  is the diagonal matrix containing the degree value in the adjacency graph;  $L$  denotes the graph Laplacian matrix, defined as  $L = D - A$ ;  $\theta_1$  and  $\theta_2$  are two weights, set as 1 and  $1e - 6$ , respectively.

In Eq. 7, the first term is the trivial square loss, and the second term aims at normalizing the desired regression function. However, unlike [37], we introduce matrix  $D$  in the third term to enforce constraints between units of the super-pixel graph, and normalize the optimized results. Since we introduce constraints between different graph units, our method can help strengthen the connections and smoothness among different

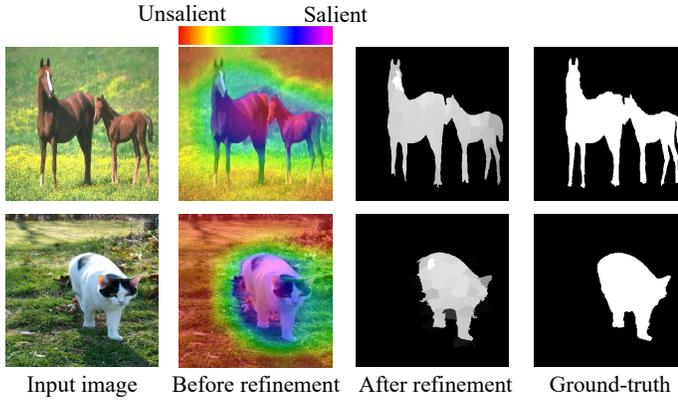


Fig. 3. Comparison between coarse saliency maps and refined saliency maps. The color code above the second column indicates the degree of saliency.

graph units. The optimization objective function is transformed into a matrix form as:

$$\min_{\alpha} \frac{1}{l} \sum_{i=1}^l \|y - JK\alpha\|_2^2 + \theta_1 \alpha^T K\alpha + \frac{\theta_2}{N^2} \alpha^T KD^{-\frac{1}{2}} LD^{-\frac{1}{2}} K\alpha, \quad (8)$$

where  $J$  is a diagonal matrix with the first  $l$  elements set to 1, while the other elements are set to 0;  $\alpha = (\alpha_1, \dots, \alpha_N)^T$ ;  $K$  is the kernel gram matrix. The solution to the above optimization problem is formulated as:

$$\alpha^* = (JK + \theta_1 lI + \frac{\theta_2 l}{N^2} D^{-\frac{1}{2}} LD^{-\frac{1}{2}} K)^{-1} y, \quad (9)$$

where  $I$  is an identity matrix. With the optimized  $\alpha^*$ , we can calculate the saliency score  $g(x)$  for each super-pixel. As presented in Figure 3, the refinement process optimizes the boundaries of the salient maps.

## IV. EXPERIMENT RESULTS

### A. Implementation Details

In this paper, we utilize ResNet-101 [34] as the backbone and modify it to meet our requirement. Those unmodified layers are initialized with weights pretrained on ImageNet [35], while the rest are randomly initialized. Our method is implemented based on the PyTorch framework. We use Stochastic Gradient Descent for parameter updating. The learning rate is initially set up as 1e-3, and will go down as the training progresses. The weight decay and the momentum is set as 5e-4 and 0.9, respectively.

It has been widely proved that inputs with various scales helps the accurate localization of target objects. Hence, the input scales are set as  $\{0.5, 0.75, 1.0\}$  of the original size. Saliency maps from three replicate networks with shared weights are summed up and regularized as  $[0, 1]$ . The proposed method is trained and evaluated on a PC with i9 3.3GHz CPU, an Nvidia 1080Ti GPU and 64GB RAM. Given an image of  $400 \times 400$ , the network takes about 0.05s to produce a single saliency map and the refinement procedure takes 0.03s per image. We apply random horizontal flipping, color scale and random rotations ( $\pm 30^\circ$ ) to augment the training datasets.

### B. Datasets and Evaluation Metrics

The ESOS dataset [15] is a saliency detection dataset, annotated with subitizing labels. It contains 17,000 images, which are selected from four datasets: MS COCO [11], Pascal VOC 2007 [10], ImageNet [35] and SUN [39]. Each image in the dataset is re-labeled by [15] with 0, 1, 2, 3 and 4+ salient objects (5 classes). We randomly choose 80% of the whole dataset (around 13,000 images) as the training set and use the rest 20% as the validation set for model selection. All images are scaled to  $256 \times 256$  during training. To compare with other saliency detection methods, we evaluate our proposed method on five benchmarks: Pascal-S [40], ECSSD [41], HKU-IS [42], DUT-OMRON [43] and MSRA-B [44]. These five datasets are all commonly used in the saliency detection task. All of them provide images and pixel-level masks.

In this paper, we adopt four metrics to measure the performance: precision-recall (PR) curve,  $F_\beta$ , S-measure [45] and mean absolute error (MAE). The continuous saliency maps are binarized with different threshold values. The PR curve is computed by comparing a series of binary masks with the ground truth masks. The second metric is defined as:

$$F_\beta = \frac{(1 + \beta^2) \cdot \text{precision} \cdot \text{recall}}{\beta^2 \cdot \text{precision} + \text{recall}}, \quad (10)$$

where  $\beta^2$  is set as 0.3 to balance between precision and recall [38]. The maximum F-measure is selected among all precision-recall pairs. The structural measure, or S-measure [45], is used to evaluate the structural similarity of non-binary saliency maps. It is defined as:

$$S = \frac{S_o + S_r}{2}, \quad (11)$$

$S_o$  is used to assess the object structure similarity against the ground truth, while  $S_r$  measures the global similarity of the target objects. Please refer to the original paper [45] for the definition of these two terms. The metric of MAE measures the average pixel-wised absolute difference between the binary ground-truth and the saliency maps. It is calculated as:

$$MAE(S, GT) = \frac{1}{W \times H} \sum_{x=1}^W \sum_{y=1}^H \|S(x, y) - GT(x, y)\|, \quad (12)$$

where  $S$  and  $GT$  are generated saliency maps and the ground truth of size  $W \times H$ .

### C. Visualized Results

We present some visualized results of our proposed method in Figure 5, which is close to the ground-truth.

### D. Comparison with Other Methods

We conduct the saliency detection task in a weakly-supervised way. There exist several other weakly-supervised methods. The difference on the settings is presented in Table II. Our method requires less extra information than those existing weakly-supervised methods. In addition, we apply the refinement process to optimize the saliency maps, instead of the commonly used CRF.

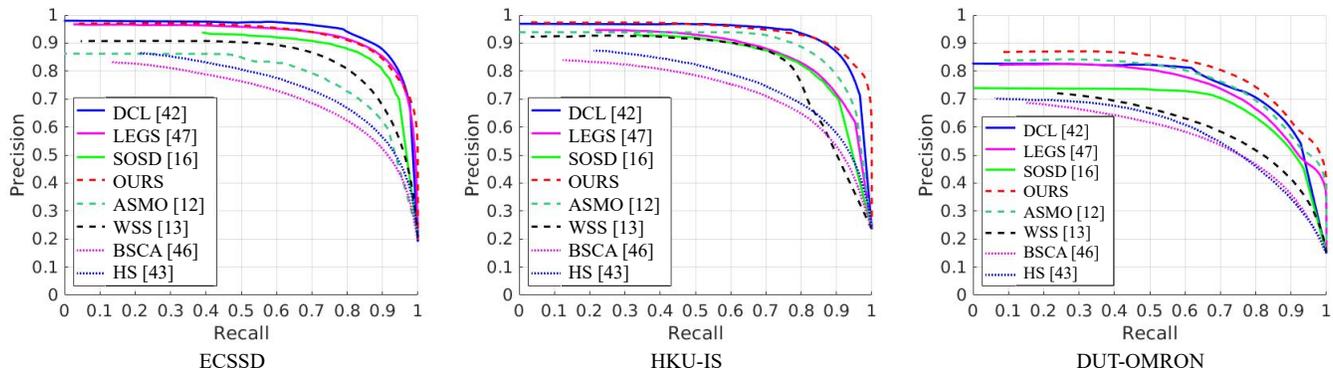


Fig. 4. Precision-recall curves of our proposed method and other methods on three benchmark datasets. Our proposed method consistently outperforms unsupervised and weakly-supervised methods. It is comparable to some fully-supervised methods.

TABLE I

QUANTITATIVE RESULTS ON  $F_\beta$  AND MAE. THE RED ONES REFER TO THE BEST RESULTS AND THE BLUE ONES REFER TO THE SECOND BEST RESULTS.

Dataset	Metric	Unsupervised		Weakly-supervised		Fully-supervised				
		BSCA [46]	HS [43]	ASMO [12]	WSS [13]	Ours	SOSD [16]	LEGS [47]	DCL [42]	MSWS [14]
ECSSD	$F_\beta \uparrow$	0.705	0.727	0.837	0.823	<b>0.858</b>	0.832	0.827	0.859	<b>0.878</b>
	MAE $\downarrow$	0.183	0.228	0.110	0.120	0.108	<b>0.105</b>	0.118	0.106	<b>0.096</b>
MSRA-B	$F_\beta \uparrow$	0.830	0.813	0.881	0.845	<b>0.897</b>	0.875	0.870	<b>0.905</b>	0.890
	MAE $\downarrow$	0.131	0.161	0.095	0.112	0.082	0.104	<b>0.081</b>	0.072	<b>0.071</b>
DUT-OMRON	$F_\beta \uparrow$	0.500	0.616	0.722	0.657	<b>0.778</b>	0.665	0.669	<b>0.733</b>	0.718
	MAE $\downarrow$	0.196	0.227	0.110	0.150	<b>0.083</b>	0.198	0.133	<b>0.094</b>	0.114
Pascal-S	$F_\beta \uparrow$	0.597	0.641	0.752	0.720	<b>0.803</b>	0.794	0.752	<b>0.815</b>	0.790
	MAE $\downarrow$	0.225	0.264	0.152	0.145	<b>0.131</b>	0.114	0.157	<b>0.113</b>	0.134
HKU-IS	$F_\beta \uparrow$	0.654	0.710	0.846	0.821	<b>0.882</b>	0.860	0.770	<b>0.892</b>	/
	MAE $\downarrow$	0.174	0.213	0.086	0.093	<b>0.080</b>	0.129	0.118	<b>0.074</b>	/

TABLE II

COMPARISON OF OUR SETTING WITH OTHER METHODS. *Seed* MEANS USING UNSUPERVISED SALIENCY MAPS AS THE INITIAL SEEDS. *Pixel* MEANS APPLYING PIXEL-WISED SUPERVISION. *CRF* MEANS USING CONDITIONAL RANDOM FIELDS AS THE POST-PROCESSING STEP. *Pseudo* MEANS ADOPTING PSEUDO LABELS.

Setting	ASMO [12]	WSS [13]	SOSD [16]	Ours
Label	image tag	image tag	subitizing+pixel	subitizing
Seed	✓	×	×	×
Pixel	×	×	✓	×
Pseudo	✓	✓	×	×
CRF	✓	✓	✓	×

We compare our results with eight state-of-the-art methods, including two unsupervised ones: BSCA [46] and HS [43]; two weakly-supervised ones using image-label supervisions: ASMO [12] and WSS [13]; four fully-supervised ones: SOSD [16], LEGS [13], DCL [42] and MSWS [14].

As shown in Table I, our proposed method outperforms existing unsupervised methods with a considerable margin. Compared to weakly-supervised methods with image-label supervisions, our method achieves better performance on all benchmarks. It proves that the subitizing supervision helps boost the saliency detection task. In addition, our method compares favorably against some fully-supervised counterparts. Note that on the DUT-OMRON dataset, our method obtains more precise results than the fully-supervised methods. Since the masks of the DUT-OMRON dataset are complex in appearance, sometimes with holes, it reveals that our method is capable of handling difficult situations. Compared to SOSD [16], which utilized additional subitizing information, our

method extracts more valid information from the subitizing supervision. Moreover, our method achieve comparable results with MSWS [14], which applied multi-source weak supervisions, including subitizing, image labels and captioning.

TABLE III

COMPARISON OF OUR RESULTS WITH TWO WEAKLY-SUPERVISED METHODS (ASMO [12] AND WSS [13] USING IMAGE-TAG SUPERVISION) AND A FULLY-SUPERVISED METHOD SOSD [16] ADDITIONALLY USING SUBITIZING INFORMATION IN TERMS OF S-MEASURE (LARGER IS BETTER).

Methods	ASMO [12]	WSS [13]	SOSD [16]	Ours
ECSSD	0.827	0.829	0.837	<b>0.860</b>
DUT	0.736	0.803	0.816	<b>0.832</b>
Pascal-S	0.702	0.815	0.742	<b>0.854</b>

The PR curves on the ECSSD, HKU-IS and DUT-OMRON datasets are presented in Figure 4. Our method consistently outperforms other unsupervised methods and weakly-supervised methods. It is also better than some fully-supervised methods like SOSD [16], LEGS [13] and DCL [42], except on the ECSSD dataset where ours is very close to the result of DCL. We also evaluate those methods on *S-measure*. The results are shown in Table III. It reveals that our method generates saliency maps of higher structural similarity compared with the ground-truth masks. The qualitative result is shown in Figure 6. The first two rows show that our method provides clear separation between multiple objects. The next four rows present that our results maintain the complete appearance of the salient objects. Moreover, we generate saliency maps with clear boundaries, as shown in the last three rows.



Fig. 5. Visualized results of our proposed method. From top to bottom, there are five groups, organized as: images, our results and the ground-truth.

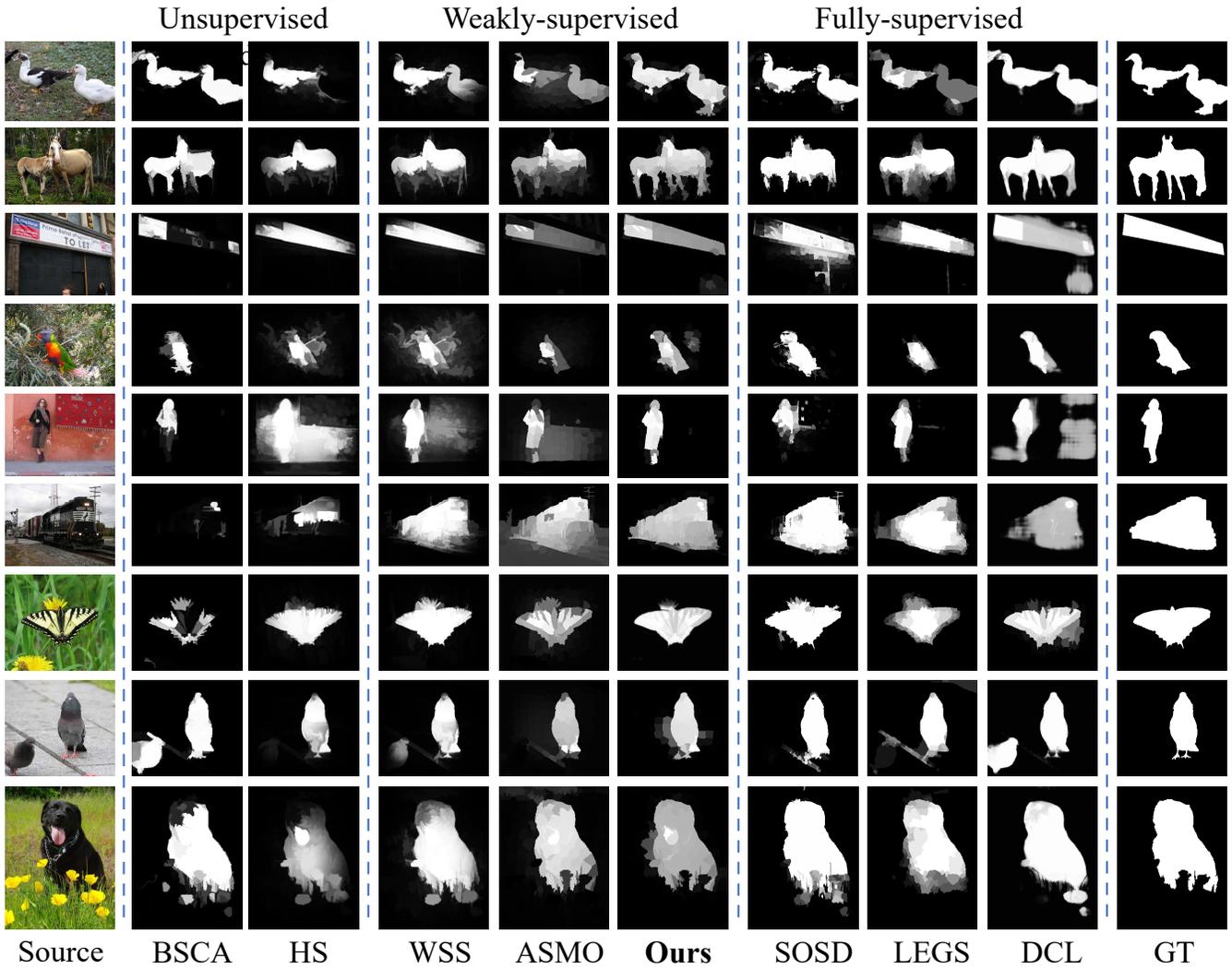


Fig. 6. Qualitative comparison with other saliency detection methods. Unsupervised methods, weakly-supervised methods and fully-supervised methods are placed from left to right. Among those weakly-supervised methods, our proposed method produces saliency maps closest to the ground-truth masks.

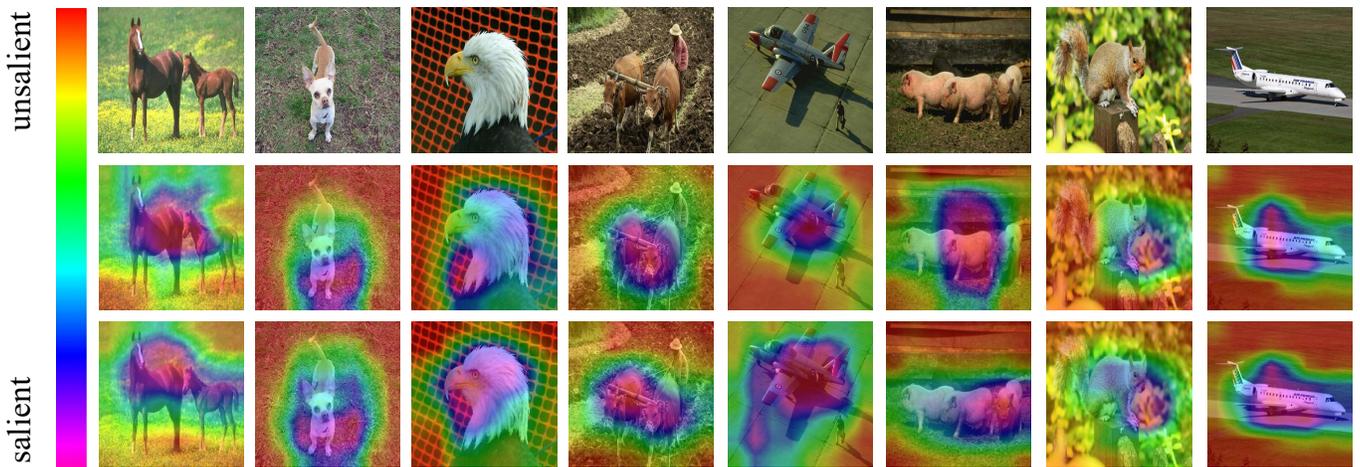


Fig. 7. Comparison between different iterations. From top to bottom, they are original images, saliency maps generated with Grad-CAM, and saliency maps after 50 iterations. Saliency maps after 50 iterations cover larger activation area belonging to the salient objects. The color code on the left represents the degree of saliency.

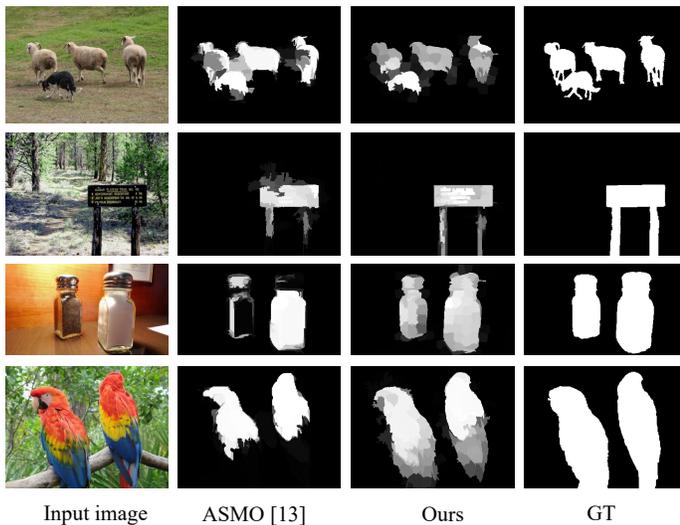


Fig. 8. Comparison between our method and ASMO [12], a weakly-supervised method using image-tag supervision. The results reveal the superiority of subitizing supervision.

The superior performance of our proposed method confirms that object subitizing generalizes better to the saliency detection task than image-level supervision. We also conduct extensive experiments to validate the performance of each component in our method.

### V. ABLATION STUDY

In this section, we discuss the advantage of subitizing supervisions over image-level supervisions. We also evaluate the effectiveness of the Saliency Updating Module and the refinement process.

#### A. The Advantage of Subitizing Supervisions

In this subsection, we aim to reveal the advantage of subitizing supervisions over image-tag supervisions. The generated saliency maps from our method are compared against those from ASMO [12] with image-label supervisions, as presented in Figure 8. The first two rows reveal that the subitizing supervision helps recognize the border between multiple objects. The last two rows indicate that our method captures the whole regions of the salient objects, while those methods supervised with image tags leave out some parts of the salient objects. In addition, we train our network with the image-tag supervisions. The results are also process with the refinement module. The performance with different training data is presented in Figure IV. The subitizing-supervised framework performs better than the image-tag-supervised framework, which reveals the advantage of subitizing supervision.

#### B. The Effect of Saliency Updating Module

In this subsection, we aim to evaluate the effect of the Saliency Updating Module. As shown in Figure 7, the results after updating are more complete in appearance, while those without any updating only focus on a limited but notable region due to the property of neural networks. With the

TABLE IV  
THE PERFORMANCE OF OUR FRAMEWORK TRAINED WITH IMAGE TAGS AND SUBITIZING, RESPECTIVELY. THE BETTER ONES ARE MARKED IN BOLD.

Dataset	Metric	w/ image-tag	w/ subitizing
ECSSD	$F_\beta \uparrow$	0.825	<b>0.858</b>
	MAE $\downarrow$	0.110	<b>0.108</b>
DUT-OMRON	$F_\beta \uparrow$	0.745	<b>0.778</b>
	MAE $\downarrow$	0.103	<b>0.083</b>

TABLE V  
THE PERFORMANCE WITH/WITHOUT THE REFINEMENT PROCESS. THE SALIENCY RESULTS WITH CRF IS ALSO PRESENTED. THE BEST PERFORMANCE IS MARKED IN BOLD.

Dataset	Metric	w/o post-pro.	w/ CRF	w/ ref.
ECSSD	$F_\beta \uparrow$	0.707	0.721	<b>0.858</b>
	MAE $\downarrow$	0.197	0.185	<b>0.108</b>
MSRA-B	$F_\beta \uparrow$	0.731	0.782	<b>0.897</b>
	MAE $\downarrow$	0.167	0.152	<b>0.082</b>
Pascal-S	$F_\beta \uparrow$	0.644	0.680	<b>0.803</b>
	MAE $\downarrow$	0.206	0.191	<b>0.131</b>

SUM module, the network captures more parts within the semantic affinity. In addition, on the DUT-OMRON dataset, the  $F_\beta$  and MAE measures of saliency predictions after 10 and 50 iterations with our SUM module are 0.638/0.252 and 0.704/0.139, respectively. It reveals that the SUM module helps boost the performance of saliency detection.

#### C. The Effect of Refinement Process

In order to obtain promising results, we adopt the refinement process to optimize the boundaries of saliency maps. As CRF is the most popular technique to refine segmentation results, we apply CRF on coarse maps and evaluate the outputs as well. The results on the ECSSD, MSRA-B and Pascal-S datasets are presented in Table V. The refinement process contributes a lot to the recognition of salient objects. It reveals that our refinement process achieves better optimization results than CRF. In addition, to evaluate the effectiveness of the refinement module on other methods, the refinement process is conducted on unsupervised results from BSCA [46]. As shown in Table VI, the refinement module helps improve the performance, but the processed results are still worse than our results.

### VI. CONCLUSION

In this paper, we propose a novel method for the salient object detection task with the subitizing supervision. We design a model with the Saliency Subitizing Module and the Saliency Updating Module, which generates the initial

TABLE VI  
THE PERFORMANCE OF UNSUPERVISED RESULTS (BSCA [46]) PROCESSED BY OUR REFINEMENT MODULE.

Dataset	Metric	BSCA [46]	BSCA w/ ref.	ours w/ ref.
ECSSD	$F_\beta \uparrow$	0.705	0.756	<b>0.858</b>
	MAE $\downarrow$	0.183	0.140	<b>0.108</b>
DUT-O	$F_\beta \uparrow$	0.500	0.618	<b>0.778</b>
	MAE $\downarrow$	0.196	0.134	<b>0.083</b>

masks using subitizing information and iteratively refines the generated saliency masks, respectively. Without any seeds from unsupervised methods, our method outperforms other weakly-supervised methods and even performs comparable to some fully-supervised methods.

#### ACKNOWLEDGMENT

We thank for the support from National Natural Science Foundation of China(61972157, 61902129), Shanghai Pujiang Talent Program (19PJ1403100), Economy and Information Commission of Shanghai (XX-RGZN-01-19-6348), National Key Research and Development Program of China (No. 2019YFC1521104), Science and Technology Commission of Shanghai Municipality Program (No. 18D1205903). Xin Tan is also supported by the Postgraduate Studentship (Mainland Schemes) from City University of Hong Kong.

#### REFERENCES

- [1] X. Qin, S. He, Z. V. Zhang, M. Dehghan, and M. Jägersand, “Real-time salient closed boundary tracking using perceptual grouping and shape priors.” in *BMVC*, 2017.
- [2] S. Goferman, L. Zelnik-Manor, and A. Tal, “Context-aware saliency detection,” *IEEE TPAMI*, vol. 34, no. 10, pp. 1915–1926, 2011.
- [3] M.-M. Cheng, F.-L. Zhang, N. J. Mitra, X. Huang, and S.-M. Hu, “Repfinder: finding approximately repeated scene elements for image editing,” *TOG*, vol. 29, no. 4, p. 83, 2010.
- [4] J. He, J. Feng, X. Liu, T. Cheng, T.-H. Lin, H. Chung, and S.-F. Chang, “Mobile product search with bag of hash bits and boundary reranking,” in *CVPR*, 2012, pp. 3005–3012.
- [5] S. He, R. W. Lau, W. Liu, Z. Huang, and Q. Yang, “Supercnn: A superpixelwise convolutional neural network for salient object detection,” *IJCV*, vol. 115, no. 3, pp. 330–344, 2015.
- [6] N. Liu and J. Han, “Dhsnet: Deep hierarchical saliency network for salient object detection,” in *CVPR*, 2016, pp. 678–686.
- [7] Y. Zhuge, Y. Zeng, and H. Lu, “Deep embedding features for salient object detection,” in *AAAI*, vol. 33, 2019, pp. 9340–9347.
- [8] L. Zhu, J. Chen, X. Hu, C.-W. Fu, X. Xu, J. Qin, and P.-A. Heng, “Aggregating attentional dilated features for salient object detection,” *IEEE TCSVT*, 2019.
- [9] X. Hu, C.-W. Fu, L. Zhu, T. Wang, and P.-A. Heng, “Sac-net: spatial attenuation context for salient object detection,” *IEEE TCSVT*, 2020.
- [10] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes (voc) challenge,” *IJCV*, vol. 88, no. 2, pp. 303–338, 2010.
- [11] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *ECCV*, 2014, pp. 740–755.
- [12] G. Li, Y. Xie, and L. Lin, “Weakly supervised salient object detection using image labels,” in *AAAI*, 2018.
- [13] L. Wang, H. Lu, Y. Wang, M. Feng, D. Wang, B. Yin, and X. Ruan, “Learning to detect salient objects with image-level supervision,” in *CVPR*, 2017, pp. 136–145.
- [14] Y. Zeng, Y. Zhuge, H. Lu, L. Zhang, M. Qian, and Y. Yu, “Multi-source weak supervision for saliency detection,” in *CVPR*, 2019, pp. 6074–6083.
- [15] J. Zhang, S. Ma, M. Sameki, S. Sclaroff, M. Betke, Z. Lin, X. Shen, B. Price, and R. Mech, “Salient object subitizing,” *IJCV*, 2017.
- [16] S. He, J. Jiao, X. Zhang, G. Han, and R. W. Lau, “Delving into salient object subitizing and detection,” in *ICCV*, 2017, pp. 1059–1067.
- [17] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, “Learning deep features for discriminative localization,” in *CVPR*, 2016, pp. 2921–2929.
- [18] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-cam: Visual explanations from deep networks via gradient-based localization,” in *ICCV*, 2017, pp. 618–626.
- [19] S. Woo, J. Park, J.-Y. Lee, and I. So Kweon, “Cbam: Convolutional block attention module,” in *ECCV*, 2018, pp. 3–19.
- [20] A. Kolesnikov and C. H. Lampert, “Seed, expand and constrain: Three principles for weakly-supervised image segmentation,” in *ECCV*, 2016, pp. 695–711.
- [21] G. Li and Y. Yu, “Visual saliency based on multiscale deep features,” in *CVPR*, 2015, pp. 5455–5463.
- [22] R. Zhao, W. Ouyang, H. Li, and X. Wang, “Saliency detection by multi-context deep learning,” in *CVPR*, 2015, pp. 1265–1274.
- [23] Y. Yuan, C. Li, J. Kim, W. Cai, and D. D. Feng, “Dense and sparse labeling with multidimensional features for saliency detection,” *IEEE TCSVT*, 2016.
- [24] W. Wang, S. Zhao, J. Shen, S. C. Hoi, and A. Borji, “Salient object detection with pyramid attention and salient edges,” in *CVPR*, 2019, pp. 1448–1457.
- [25] H. Qibin, C. Ming-Ming, H. Xiaowei, B. Ali, T. Zhuowen, and H. S. T. Philip, “Deeply supervised salient object detection with short connections,” *IEEE TPAMI*, vol. 41, no. 4, pp. 815–828, 2019.
- [26] Z. Tu, Y. Ma, C. Li, J. Tang, and B. Luo, “Edge-guided non-local fully convolutional network for salient object detection,” *IEEE TCSVT*, 2020.
- [27] Z. Zhou, Z. Wang, H. Lu, S. Wang, and M. Sun, “Multi-type self-attention guided degraded saliency detection,” in *AAAI*, 2020, pp. 13 082–13 089.
- [28] J. Zhang, S. Ma, M. Sameki, S. Sclaroff, M. Betke, Z. Lin, X. Shen, B. Price, and R. Mech, “Salient object subitizing,” in *CVPR*, 2015, pp. 4045–4054.
- [29] E. Lu, W. Xie, and A. Zisserman, “Class-agnostic counting,” in *ACCV*, 2018, pp. 669–684.
- [30] M. Amirul Islam, M. Kalash, and N. D. Bruce, “Revisiting salient object detection: Simultaneous detection, ranking, and subitizing of multiple salient objects,” in *CVPR*, 2018, pp. 7142–7150.
- [31] K. Li, Z. Wu, K.-C. Peng, J. Ernst, and Y. Fu, “Tell me where to look: Guided attention inference network,” in *CVPR*, 2018, pp. 9215–9223.
- [32] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. Torr, “Conditional random fields as recurrent neural networks,” in *ICCV*, 2015, pp. 1529–1537.
- [33] C. Chen, X. Sun, Y. Hua, J. Dong, and H. Xv, “Learning deep relations to promote saliency detection,” in *AAAI*, 2020, pp. 10 510–10 517.
- [34] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016, pp. 770–778.
- [35] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *CVPR*, 2009, pp. 248–255.
- [36] Y. Wei, J. Feng, X. Liang, M.-M. Cheng, Y. Zhao, and S. Yan, “Object region mining with adversarial erasing: A simple classification to semantic segmentation approach,” in *CVPR*, 2017, pp. 1568–1576.
- [37] X. Li, L. Zhao, L. Wei, M.-H. Yang, F. Wu, Y. Zhuang, H. Ling, and J. Wang, “Deepsaliency: Multi-task deep neural network model for salient object detection,” *IEEE TIP*, vol. 25, no. 8, pp. 3919–3930, 2016.
- [38] R. Achanta, S. Hemami, F. Estrada, and S. Süsstrunk, “Frequency-tuned salient region detection,” in *CVPR*, 2009, pp. 1597–1604.
- [39] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba, “Sun database: Large-scale scene recognition from abbey to zoo,” in *CVPR*, 2010, pp. 3485–3492.
- [40] Y. Li, X. Hou, C. Koch, J. M. Rehg, and A. L. Yuille, “The secrets of salient object segmentation,” in *CVPR*, 2014, pp. 280–287.
- [41] J. Shi, Q. Yan, L. Xu, and J. Jia, “Hierarchical image saliency detection on extended cssd,” *IEEE TPAMI*, vol. 38, no. 4, pp. 717–729, 2015.
- [42] G. Li and Y. Yu, “Deep contrast learning for salient object detection,” in *CVPR*, 2016, pp. 478–487.
- [43] C. Yang, L. Zhang, H. Lu, X. Ruan, and M.-H. Yang, “Saliency detection via graph-based manifold ranking,” in *CVPR*, 2013, pp. 3166–3173.
- [44] H. Jiang, J. Wang, Z. Yuan, Y. Wu, N. Zheng, and S. Li, “Salient object detection: A discriminative regional feature integration approach,” in *CVPR*, 2013, pp. 2083–2090.
- [45] D.-P. Fan, M.-M. Cheng, Y. Liu, T. Li, and A. Borji, “Structure-measure: A new way to evaluate foreground maps,” in *ICCV*, 2017, pp. 4548–4557.
- [46] Y. Qin, H. Lu, Y. Xu, and H. Wang, “Saliency detection via cellular automata,” in *CVPR*, 2015, pp. 110–119.
- [47] L. Wang, H. Lu, X. Ruan, and M.-H. Yang, “Deep networks for saliency detection via local estimation and global search,” in *CVPR*, 2015, pp. 3183–3192.