

---

# Joint Semantic Mining for Weakly Supervised RGB-D Salient Object Detection

---

Jingjing Li<sup>1,\*</sup>, Wei Ji<sup>1,\*</sup> (✉), Qi Bi<sup>2</sup>, Cheng Yan<sup>3</sup>, Miao Zhang<sup>4</sup>,  
Yongri Piao<sup>4</sup>, Huchuan Lu<sup>4,5</sup>, Li Cheng<sup>1</sup>

<sup>1</sup>University of Alberta, Canada   <sup>2</sup>Wuhan University, China   <sup>3</sup>Tianjin University, China  
<sup>4</sup>Dalian University of Technology, China   <sup>5</sup>Pengcheng Lab, Shenzhen, China

## Abstract

Training saliency detection models with weak supervisions, *e.g.*, image-level tags or captions, is appealing as it removes the costly demand of per-pixel annotations. Despite the rapid progress of RGB-D saliency detection in fully-supervised setting, it however remains an unexplored territory when only weak supervision signals are available. This paper is set to tackle the problem of *weakly-supervised RGB-D salient object detection*. The key insight in this effort is the idea of maintaining per-pixel pseudo-labels with iterative refinements by reconciling the multimodal input signals in our joint semantic mining (JSM). Considering the large variations in the raw depth map and the lack of explicit pixel-level supervisions, we propose spatial semantic modeling (SSM) to capture saliency-specific depth cues from the raw depth and produce depth-refined pseudo-labels. Moreover, tags and captions are incorporated via a fill-in-the-blank training in our textual semantic modeling (TSM) to estimate the confidences of competing pseudo-labels. At test time, our model involves only a light-weight sub-network of the training pipeline, *i.e.*, it requires only an RGB image as input, thus allowing efficient inference. Extensive evaluations demonstrate the effectiveness of our approach under the weakly-supervised setting. Importantly, our method could also be adapted to work in both fully-supervised and unsupervised paradigms. In each of these scenarios, superior performance has been attained by our approach with comparing to the state-of-the-art dedicated methods. As a by-product, a *CapS* dataset is constructed by augmenting existing benchmark training set with additional image tags and captions. *Code and dataset are available at <https://github.com/jiwei0921/JSM>.*

## 1 Introduction

As a fundamental computer vision task, salient object detection (SOD) aims at locating and segmenting visually distinctive objects in a scene. It plays an important role in a variety of downstream applications including image retrieval [31, 65], medical analysis [48, 28, 25], multimodal fusion [79, 80] and video analysis [88, 81, 90]. Recent progress in supervised RGB-D SOD [5, 52, 10, 26, 89, 33] has demonstrated significant benefits of engaging depth information in saliency detection from complex scenes. The success of these fully-supervised methods, however, relies heavily on the large-scale, precise, pixel-level annotations, which are often laborious and time-consuming to acquire. On the other hand, an image usually comes with additional information in its meta-data such as tags and captions from users to describe the scene context and content, which may serve as cheap weak-supervision signals. These weak supervision signals are nonetheless noisy and have mixed qualities. Similar weak-supervision signals have been explored in RGB-image based SOD [63, 70, 38], where the noisy nature of these side information is unfortunately overlooked. To further complicate the matter, the lack of explicit pixel-level supervision brings new challenge to the RGB-D SOD task: the depth values from raw depth maps are often noisy and sometimes inconsistent. For example, in Fig. 1,

---

\*means equal contributions. Wei Ji (✉) ([wji3@ualberta.ca](mailto:wji3@ualberta.ca)) is the corresponding author.

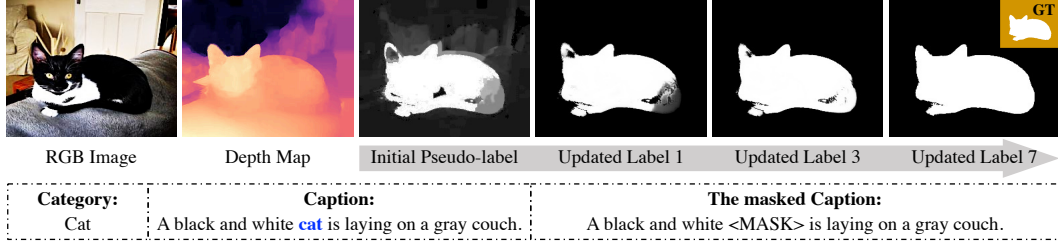


Figure 1: Illustration of weakly-supervised RGB-D salient object detection. RGB and depth images, as well as weak supervision signals such as image-level tags and captions are exploited. Initial pseudo-label is generated by the handcrafted methods, which is then iteratively updated by our joint semantic mining pipeline. GT is ground-truth label for reference.

similar depth values are shared by the cat and the underneath couch, making it difficult to discern the salient object from backgrounds. Without the explicit pixel-level supervision, existing cross-modal fusion strategies adopted by fully-supervised RGB-D methods [36, 55, 37, 27] would simply fail.

These observations motivate us to consider the new problem of *weakly-supervised RGB-D salient object detection*, which takes as input the RGB and depth images, as well as weak supervision signals such as image-level tags and captions, as illustrated in Fig. 1. By removing the demand for laborious per-pixel annotations, it also brings new challenges: 1) how to address the noisy nature of the weak supervision signals; 2) how to tackle the depth noise and inconsistency to facilitate proper separation of foreground and background regions.

This leads us to propose the use of pseudo-labels with iterative refinements in training: the pseudo-label provides internal pixel-level supervision signals, which is progressively updated by reconciling the multimodal input signals and the current information flow of the neural net, based on the previous pseudo-label. As illustrated in Fig. 2, this is realized by an interaction between two core modules, namely spatial semantic modeling (SSM) and joint semantic mining (JSM): SSM is designed to capture the saliency-specific depth semantics, to eliminate the background noises in the coarse saliency prediction, and to generate a depth-refined pseudo-label. This simple yet effective module is very generic, which could be easily plugged-in different setups, including unsupervised & fully-supervised scenarios; meanwhile, the JSM module is proposed to leverage depth semantics and weak supervision signals for attaining more reliable pseudo-labels. Specifically, a partial textual input, *i.e.*, image-level tag and caption with its salient word being masked, is fed into a dedicated textual semantic modeling or TSM to estimate the confidence scores of competing pseudo-labels, and to fill-in-the-blank. Intuitively, a semantically consistent pseudo-label should provide better context cues to reconstruct the salient word; while a closer guess of the masked word would indicate a better pseudo-label. The alternation between SSM and JSM modules is thus expected to give rise to more trustworthy pseudo-labels. At test time, it is then sufficient to take an input RGB image and activate a light-weight network to deliver its final prediction. That is, test time input involves only an RGB image, without the need of any depth map or image-level tags and captions. This drastically simplifies the input requirement and reduces the computation burden at deployment stage. Moreover, given the lack of training dataset for the weakly-supervised setting, we adapt existing RGB-D training dataset to annotate additional image-level tags and captions, which is referred to as the *CapS* dataset.

The main contributions of this paper are as follows. (1) A new problem of weakly-supervised RGB-D salient object detection is considered. In this regard, a *CapS* dataset is curated by augmenting the existing RGB-D SOD training dataset with image-level tagging and captioning annotations. (2) The key ingredient of our approach involves the production of pseudo-labels with iterative refinements, realized by iterative updates between two internal modules, SSM and JSM. Empirical experiments demonstrate the effectiveness of our approach in weakly-supervised setting. Moreover, (3) after proper adaptation of our approach in unsupervised and fully-supervised scenarios, superior performance is also observed when comparing to the respective state-of-the-art methods. (4) Our test time inference amounts to executing a light-weight saliency network: as illustrated in dotted line at Fig. 2, only an RGB image is used as its input, thus allows efficient and effective inference.

## 2 Related Work

RGB-D salient object detection [87, 93, 94] has been an active field of research in the past few years, where the incorporation of depth cues has been demonstrated [13, 34, 82, 75, 15, 35, 60, 43, 85]

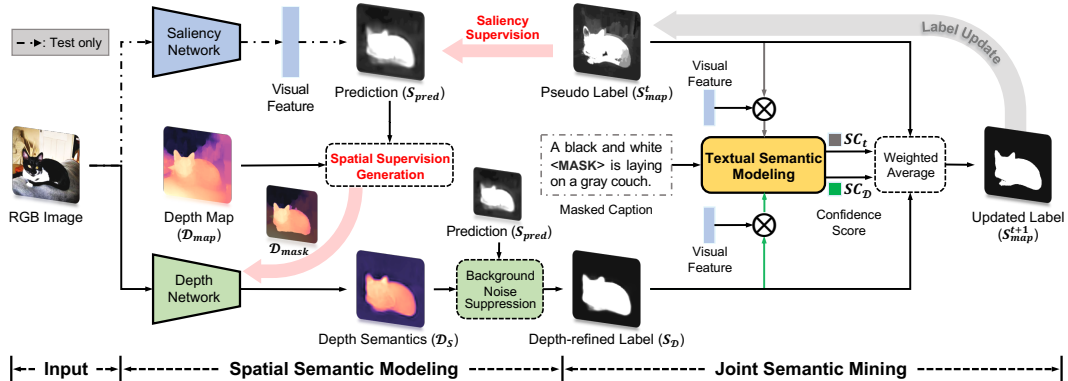


Figure 2: An overview of our approach. Its training pipeline consists of a saliency prediction network, a SSM (Sec. 3.2) to generate depth-refined pseudo-label, a TSM (Sec. 3.3) to estimate the confidences of different pseudo-labels, and a JSM (Sec. 3.4) to refine & update pseudo-label. Our testing process only involves activating a saliency network delineated in dotted lines. More details of SSM and TSM modules are illustrated in Fig. 3. The masked salient word is ‘cat’.

to improve performance especially in complex scenes. Existing RGB-D methods aim to design effective feature fusion strategies for learning representative cross-modal features. Typically, Chen *et al.* [8] employ two-stream CNNs-based models and perform fusion by adding or concatenating paired features at shallow or deep layers. Fu *et al.* [19] utilize a Siamese network to jointly learn RGB and depth inputs for mining useful complementary features. To promote multi-modal interactions, Li *et al.* [37] design a cross-modal weighting strategy to encourage comprehensive interactions between RGB and depth information. However, those methods unfortunately rely on costly pixel-level annotations, which are tedious and time-consuming to acquire. This motivates us to consider a weakly-supervised approach. In what follows, our focus will be mainly toward related weakly-supervised methods developed for saliency detection from RGB images, where the differences of our approach from existing methods would be clarified.

Instead of using costly pixel-level annotations, some recent efforts instead explore cheap alternatives such as image-level tags (categories) [38, 63, 71, 46, 2], image captions [70], scribble labels [69, 76], and noisy pseudo-labels from handcrafted methods [45, 72, 73, 77, 68, 47, 67, 49, 66]. The pioneering work [63] leverage image-level tags or categories that could be augmented onto existing large-scale dataset at low-cost. The same scenario is also considered by Li *et al.* [38], where a composite pipeline combining graphical model with CNNs is designed. The trained network however tends to highlight the most discriminative region instead of the intended salient object out of the scene due to the sparse image-level supervisions. Image captions are examined by Zeng *et al.* [70] as supervision input; in their work image classification network and caption generation network are jointly trained to obtain pseudo-labels, which achieves descent performance. Scribble is another type of weak supervision signal, where a tiny fraction of image pixels are labeled by users as being foreground or background. Due to the annotation sparsity, object structure and details cannot be easily inferred. Zhang *et al.* [76] introduce a gated structure-aware loss as well as an auxiliary edge detection network to uncover the complete object. Meanwhile, Yu *et al.* [69] explore self-consistency among multi-scale outputs and design a local coherence loss to propagate the labels to unlabeled regions based on image features, thus enabling the detection of objects with smooth textures. However, existing weakly supervised saliency methods are solely based on RGB image. Unlike the prevalence of fully-supervised RGB-D SOD, it has never been considered the incorporation of depth cues in existing literature.

This leads us to address this problem in the presence of image tags and captions as weak supervision signals. Different from existing methods [38, 24, 63, 70] that train image classification networks or caption generation networks to delineate potential salient regions, masked salient word in captions are to be reconstructed in our work by leveraging visual saliency features. This is then used to estimate the confidence scores of pseudo-labels.

### 3 Methodology

#### 3.1 The Overall Architecture

An overview of our approach is illustrated in Fig. 2. It consists of a saliency network responsible for saliency prediction, a spatial (depth) semantic modeling (SSM) to generate depth-refined pseudo-label,

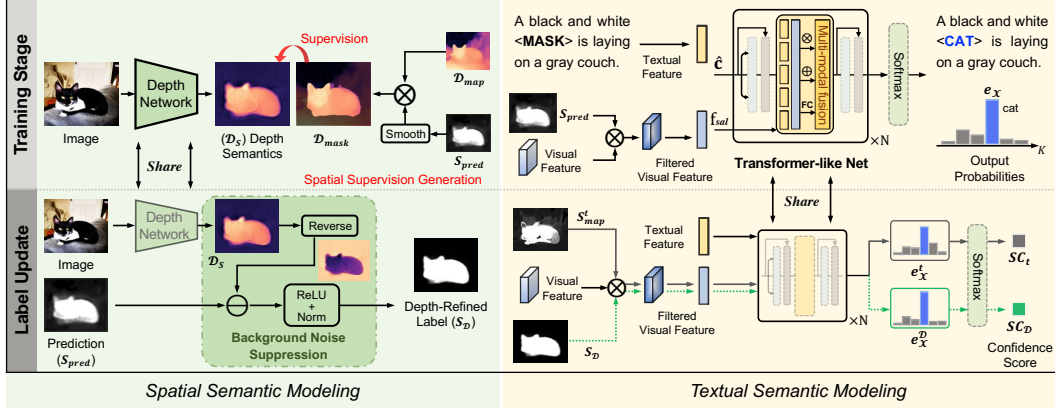


Figure 3: The detailed architecture of the proposed SSM and TSM. The upper shows their training processes, and the bottom illustrates the way of using them to perform label update.

a textual (caption) semantic modeling (TSM) to estimate the confidences of different pseudo-labels, and a joint semantic mining (JSM) strategy to refine & update pseudo-label. Overall our pipeline aims to gradually improve the quality of noisy pseudo-labels by jointly mining the useful spatial semantics from the depth map and textual semantics from the tags and captions, to produce more trustworthy supervision signals, which in turns results in better training of the saliency network.

Specifically, the popular encoder-decoder architecture [64] in SOD is adopted in both saliency network and depth network. Initial supervision signal for the saliency network is provided by traditional handcrafted methods. The predicted saliency map, together with the raw depth map, are processed to generate the saliency-guided spatial supervision signal for the depth network. Then the predicted saliency-oriented depth semantics is utilized to eliminate background noises (non-salient regions) in coarse prediction, and to generate a depth-refined pseudo-label. This is followed by our JSM strategy, which takes in the image-level tags and captions through TSM to estimate the confidence scores of pseudo-labels; updated pseudo-label is then formed based on the confidence-weighted depth-refined pseudo-label and current pseudo-label, which provides more trustworthy supervision signal for the saliency network. Note our test time inference involves only the black dashed portion in Fig. 2, which takes as input only the RGB image, thus enables efficient saliency prediction.

### 3.2 Saliency-oriented Spatial (Depth) Semantic Modeling

Initial pseudo-labels are generated by traditional saliency models, which often contain excessive noise. As illustrated in Fig. 3, our spatial semantic modeling (SSM) is to produce a more reliable depth-refined pseudo-label, achieved by explicitly capturing saliency-specific depth semantics from the depth map to eliminate possible background noise in the coarse saliency prediction.

Concretely, during training, we first generate a saliency-guided depth mask  $\mathcal{D}_{mask}$  by multiplying the rough saliency prediction  $\mathcal{S}_{pred}$  with the raw depth map  $\mathcal{D}_{map}$  in a spatial attention manner. Here, a Gaussian smooth operation is applied to smooth the predicted saliency area, to effectively perceive and capture more saliency areas from depth. The procedure is formulated as:

$$\mathcal{D}_{mask} = \Omega_{max}(\mathcal{F}_{gauss}(\mathcal{S}_{pred}, k), \mathcal{S}_{pred}) \otimes \mathcal{D}_{map}, \quad (1)$$

where  $\mathcal{F}_{gauss}(\cdot, k)$  indicates a convolution operation with Gaussian kernel  $k$  and zero bias;  $\Omega_{max}(\cdot, \cdot)$  is a maximum function to preserve the higher values between the smoothed and the original maps.  $\otimes$  is element-wise multiplication. In this paper, the size and standard deviation of kernel  $k$  are learnable through the model training procedure and are initialized with values 32 and 4, respectively.

After obtaining  $\mathcal{D}_{mask}$ , a depth network is trained to learn the saliency-specific depth semantics  $\mathcal{D}_S$ , using the mean square error (MSE) loss function. The internal inspection evidences in Fig. 4 suggest that  $\mathcal{D}_S$  (depth semantics, 6<sup>th</sup> column) is able to capture discriminative saliency cues from the raw depth map under the supervision of  $\mathcal{D}_{mask}$  (5<sup>th</sup> column).

In addition, the learned  $\mathcal{D}_S$  can be further processed to generate the depth-refined pseudo-labels. This procedure is only employed when performing pseudo-label update. Specifically, we feed  $\mathcal{D}_S$  into an Background Noise Suppression block to help eliminate the background noises (non-salient

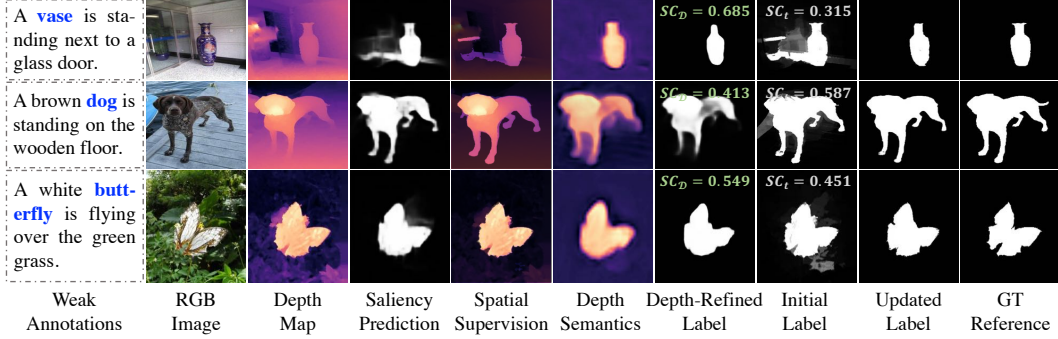


Figure 4: Step-by-step inspections of the internal processes of our approach. The GT is for reference.

regions) in the coarse prediction  $\mathcal{S}_{pred}$ . In this block, a reverse operation is first employed on  $\mathcal{D}_S$  to highlight background regions by  $1 - \mathcal{D}_S$ . This is followed by the pixel-wise subtraction to suppress the non-salient negative responses in  $\mathcal{S}_{pred}$ . Finally, a ReLU function and a normalization procedure are adopted to suppress negative numbers and normalize the result to the range of  $[0, 1]$ . This procedure of obtaining the depth-refined pseudo-label  $\mathcal{S}_D$  could be expressed by

$$\mathcal{S}_D^{i,j} = \frac{\mathcal{S}_d^{i,j} - \min(\mathcal{S}_d)}{\max(\mathcal{S}_d) - \min(\mathcal{S}_d) + \varepsilon}, \mathcal{S}_d = \text{ReLU}(\mathcal{S}_{pred} - \lambda_d(1 - \mathcal{D}_S)), i \in [1, H], j \in [1, W]. \quad (2)$$

Here  $\lambda_d \in [0, 1]$  is a constant to control the degree of the subtracted background noises and avoid negatively suppressing salient regions that have relatively low saliency scores in  $\mathcal{S}_D$ . Throughout experiments,  $\lambda_d$  is empirically set to 0.5, and  $\varepsilon$  to 1e-5.  $H$  and  $W$  are the height and width of the input image, respectively.

The pseudo-label refinement dynamics are visually inspected in Fig. 4, while Fig. 8 presents the corresponding quantitative error analysis over iterations. Empirical evidence suggests as the training proceeds, quality of the pseudo-label is significantly improved. Moreover, at this stage, we can directly utilize the depth-refined label  $\mathcal{S}_D$  to update pseudo-labels (as in Fig. 6 (b)). This can be referred as unsupervised RGB-D SOD since no weak labels are used.

### 3.3 Saliency-oriented Textual (Caption) Semantic Modeling

Previously, the mainstream use of weak labels is to train a classification network or a caption generation network, where the by-product attention maps or Class Activation Maps [92] are leveraged to determine the potential salient regions [70, 63]. It is very different in our textual semantic modeling (TSM), where the main focus is to leverage side information (*i.e.*, image-level tags and captions) to facilitate the production of reliable training signals. Inspired by the recent success of masked language models [14], captions with missing keywords are used as input, with the expectation of the complete text being reconstructed as output. In the proposed TSM, innovatively taking as input partial text with salient word being masked, as well as the saliency-filtered visual features, our TSM is to output the reconstructed text in a fill-in-the-blank manner and to estimate the confidence scores of competing pseudo-labels. The intuition is, a semantically matching pseudo-label could provide better context cues to reconstruct the masked salient word; meanwhile, a closer guess of the masked word would indicate a better pseudo-label.

Formally, for each training data, the weak labels contain caption description  $\mathbf{c} = \{\mathbf{c}_i\}_{i=1}^{n_c}$ , image-level category  $k \in \{1, \dots, K\}$ , and the position  $\mathcal{X}$  of the salient word (object) in the caption, where  $n_c$  is the word number of the caption. Let  $\mathbf{c}_i \in \mathbb{R}^{d \times 1}$  be the word embedding of the  $i$ -th word in the caption,  $K$  the total number of salient categories, and  $\mathcal{X}$  an integral number. As shown in Fig. 3, during training, the input is a masked version of caption  $\hat{\mathbf{c}} \in \mathbb{R}^{d \times n_c}$  where the salient word  $\mathbf{c}_{\mathcal{X}}$  in caption  $\mathbf{c}$  is masked with a special symbol. In order to reconstruct the masked salient word, we filter the visual feature from the saliency network by multiplying it with the learned saliency cues  $\mathcal{S}_{pred}$ . We then obtain the saliency-filtered visual feature, and transform it to a feature vector  $\mathbf{f}_{sal} \in \mathbb{R}^{d \times 1}$  for subsequent cross-modal fusion using a pooling operation and a fully-connected (FC) layer.

The center component of the TSM module is a transformer-like network. Based on original Transformer [62], we add a multi-modal fusion sub-layer into network. Three parallel operations are used to promote sufficient cross-modal feature interactions: element-wise multiplication, element-wise

addition, and concatenation followed by FC. The three outputs are concatenated and followed by a FC to change the feature dimension. Note the cross-modal fusion operation is performed word-wise. Through the textual network, we can obtain the energy vector  $\mathbf{e}_\chi \in \mathbb{R}^{K \times 1}$  of the masked salient word, which is computed over all categories by a fully-connected layer and softmax function. For each training sample, the training objective loss for the textual network is  $-\log(\mathbf{e}_\chi[k])$ , where  $\mathbf{e}_\chi[k]$  is the output probability of salient category  $k$ . *Detailed structure for our TSM can be accessed in the supplementary material.* Once trained, the TSM module is then used to estimate the confidence scores of pseudo-labels as described in Sec. 3.4.

### 3.4 Joint Semantic Mining for Label Updating

The label updating operation using joint semantic mining is iteratively conducted every  $\tau$  epochs (as a training round) during training, *i.e.*, the granularity of label updating. The choices are explored in the ablations of Sec. 4.4, where  $\tau = 5$  is shown to work best empirically. In terms of training, the saliency network, the SSM module, and the TSM module are independently trained.

At the end of each training round, our JSM strategy performs the label updating operation to generate an up-to-date pseudo-label for each training image. Now define the pseudo-label for the saliency network in current iteration as  $\mathcal{S}_{map}^t$ . As shown in the label update phase of Fig. 3, the SSM module is engaged to generate the depth-refined pseudo-label  $\mathcal{S}_D$ ; It is passed to the TSM module, where the  $\mathcal{S}_{map}^t$  and  $\mathcal{S}_D$  are taken as saliency attention maps to filter visual features, respectively. The TSM module infers the energy vectors  $\mathbf{e}_\chi^t$  and  $\mathbf{e}_\chi^D$  for  $\mathcal{S}_{map}^t$  and  $\mathcal{S}_D$ , respectively. Thus their confidence scores  $\mathcal{S}\mathcal{C}_t$  and  $\mathcal{S}\mathcal{C}_D$  can be calculated by (taking  $\mathcal{S}\mathcal{C}_t$  as an example):

$$\mathcal{S}\mathcal{C}_t = \frac{\exp(\mathbf{e}_\chi^t[k])}{\exp(\mathbf{e}_\chi^t[k]) + \exp(\mathbf{e}_\chi^D[k])}. \quad (3)$$

The updated label is the weighted average of  $\mathcal{S}_{map}^t$  and  $\mathcal{S}_D$ :  $\mathcal{S}_{map}^{t+1} = \mathcal{S}\mathcal{C}_t \times \mathcal{S}_{map}^t + \mathcal{S}\mathcal{C}_D \times \mathcal{S}_D$ . A fully-connected conditional random field [23] is applied to generate the final updated label which could provide more trustworthy supervision signal to train the saliency network.

### 3.5 The CapS Dataset

To train the weakly-supervised RGB-D SOD network, we relabel two widely-used RGB-D saliency datasets, NJUD [29] and NLPR [53], which contain a total of 2,185 training images. We provide various image-level weak annotations: categories, captions, and position of the salient word in each of the captions. The annotation process is conducted semi-automatically: the NeuralTalk2 [30] toolkit is utilized to automatically generate image captions, which are then manually checked, with unreasonable cases corrected. On average, there are 8.9 words in each caption. This is followed by categorically tagging each of the images, each corresponds to a salient category; this is different from the traditional image classification dataset ImageNet [32]. We summarize the 100 categories in the *CapS*. The positions of the salient words are first automatically identified through localizing the image categories in the captions which is subsequently followed by manual identification. *Detailed statistics and examples of our in-house CapS dataset are relegated to the supplementary material.*

## 4 Experiments

### 4.1 Datasets and Evaluation Metrics

Empirical evaluations are conducted over four large-scale RGB-D SOD benchmark datasets, including NJUD [29] with 1,985 RGB-D paired images, NLPR [53] with 1,000 samples, STERE [50] with 1,000 stereoscopic images, and DUTLF-Depth [55] with 1,200 RGB-D data. In train *vs.* test splits of the datasets, the popular setup of [18, 19, 74] is followed for a fair comparison. Training set consists of 1,485 samples from NJUD and 700 samples from NLPR. The remaining images in these datasets and other public test sets are reserved for testing purposes throughout the experiments. Four widely-used metrics are adopted for quantitative evaluation: they are E-measure ( $E_\epsilon$ ) [16], F-measure ( $F_\beta$ ) [1], weighted F-measure ( $F_\beta^w$ ) [44], and mean absolute error (MAE or  $\mathcal{M}$ ) [3].

### 4.2 Implementation Details and Setups

The code is implemented in Pytorch toolbox on a PC with a single Tesla P40 GPU. We use decoder part [64] with ResNet-50 [22] pre-trained on ImageNet as backbone, for both saliency network and depth network. For each word in the caption, we extract word embedding with dimension  $d = 300$  using the pretrained Glove [54] word2vec network. The maximum caption length is set to 20. For

Table 1: Quantitative comparison with weakly-supervised and unsupervised saliency models. Note RGB-based methods are specifically marked by ‡.  $\mathcal{D}_S$  and  $\mathcal{T}_S$  represent the spatial semantic modeling and textual semantic modeling, respectively. ‘Un’ means unsupervised learning. ‘Cls’ is SOD with class label. ‘Cap’ represents using weak supervision signals, with both class label and image caption.

*	Sup.	DUTLF-Depth [55]				STERE [50]				NJUD [29]				NLPR [53]			
		$E_\xi$	$F_\beta^w$	$F_\beta$	$\mathcal{M}$	$E_\xi$	$F_\beta^w$	$F_\beta$	$\mathcal{M}$	$E_\xi$	$F_\beta^w$	$F_\beta$	$\mathcal{M}$	$E_\xi$	$F_\beta^w$	$F_\beta$	$\mathcal{M}$
RBD‡ [98]	Un	.733	.447	.619	.222	.730	.443	.610	.223	.684	.387	.556	.256	.765	.388	.590	.211
MST‡ [61]	Un	.678	.254	.401	.279	.681	.312	.447	.269	.670	.291	.436	.281	.762	.257	.491	.199
BSCA‡ [57]	Un	.808	.479	.682	.181	.803	.497	.676	.179	.756	.446	.623	.216	.745	.376	.554	.178
DSR‡ [39]	Un	.797	.478	.640	.164	.785	.486	.645	.165	.739	.436	.594	.196	.757	.451	.545	.120
ACSD [29]	Un	.250	.210	.188	.668	.793	.425	.661	.200	.790	.448	.696	.198	.751	.327	.547	.171
DES [11]	Un	.733	.386	.668	.280	.673	.383	.592	.297	.421	.241	.165	.448	.735	.259	.583	.301
LHM [53]	Un	.767	.350	.659	.174	.772	.360	.703	.171	.722	.311	.625	.201	.772	.320	.520	.119
GP [59]	Un	-	-	-	-	.785	.371	.710	.182	.730	.323	.666	.204	.813	.347	.670	.144
CDB [40]	Un	-	-	-	-	.808	.436	.713	.166	.752	.408	.650	.200	.810	.388	.618	.108
SE [20]	Un	.730	.339	.474	.196	.825	.546	.747	.143	.780	.518	.735	.164	.853	.578	.701	.085
DCMC [12]	Un	.712	.290	.406	.243	.832	.529	.743	.148	.796	.506	.715	.167	.684	.265	.328	.196
MB [96]	Un	.691	.464	.577	.156	.693	.455	.572	.178	.643	.369	.492	.202	.814	.574	.637	.089
CDCP [97]	Un	.794	.530	.633	.159	.797	.596	.666	.149	.751	.522	.618	.181	.785	.512	.591	.114
<b>Ours (<math>\mathcal{D}_S</math>)</b>	Un	<b>.845</b>	<b>.629</b>	<b>.741</b>	<b>.116</b>	<b>.838</b>	<b>.654</b>	<b>.738</b>	<b>.111</b>	<b>.781</b>	<b>.576</b>	<b>.689</b>	<b>.147</b>	<b>.857</b>	<b>.638</b>	<b>.698</b>	<b>.073</b>
WSS‡ [63]	Cls	.843	.634	.743	.125	.831	.664	.732	.115	.760	.566	.693	.149	.856	.651	.712	.077
MSW‡ [70]	Cap	.863	.678	.789	.105	.836	.675	.760	.103	.773	.580	.704	.141	.869	.659	.736	.071
<b>Ours (<math>\mathcal{D}_S \&amp; \mathcal{T}_S</math>)</b>	Cap	<b>.870</b>	<b>.698</b>	<b>.797</b>	<b>.093</b>	<b>.852</b>	<b>.688</b>	<b>.778</b>	<b>.095</b>	<b>.788</b>	<b>.599</b>	<b>.717</b>	<b>.133</b>	<b>.888</b>	<b>.692</b>	<b>.770</b>	<b>.060</b>

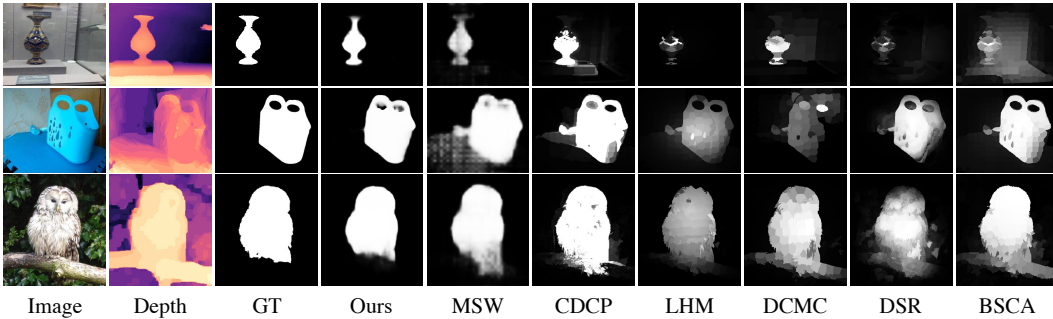


Figure 5: Visual comparison of weakly-supervised and unsupervised saliency models.

the transformer-like net in TSM, we set the hidden state to 256, the number of layers to 3, and the number of head to 4. The model is optimized by Adam with batch size of 10, and the learning rate is set to  $1 \times 10^{-4}$ . During training, we use the standard BCE loss to train the saliency network. Each image is uniformly resized to  $352 \times 352$  and is performed by randomly rotating and cropping to avoid potential overfitting. Our network is trained in an end-to-end manner and converges around 50 epochs. Initial pseudo labels are generated by the handcrafted method [97], which is freely available for annotations.

### 4.3 Model Performance

We quantitatively evaluate the performance of our approach in Table 1, with visual results shown in Fig. 5. Since our approach is the first work for weakly-supervised RGB-D saliency detection, we show the results of two state-of-the-art RGB-based weakly-supervised methods (WSS [63] and MSW [70]) for reference. To make a fair comparison, we fine-tune them on the same training set using their published code with default setups. These results show the effectiveness of our proposed method. Furthermore, the SSM and TSM in our joint semantic mining framework do not introduce any additional inference cost since they only participate in the training procedure to provide more reliable supervisory signals for saliency network. This design leads to a light-weight network that works both efficiently and effectively. As shown in Table 2, our method runs fastest among RGB-D methods, and also on par with the most efficient RGB-based methods.

Table 2: Inference time of different unsupervised and weakly-supervised saliency models. The RGB-based methods are specifically marked by ‡.

*	LHM [53]	DES [11]	GP [59]	CDCP [97]	DCMC [12]	SE [20]	ACSD [29]	CDB [40]	BSCA‡ [57]	RBD‡ [98]	MST‡ [61]	DSR‡ [39]	MSW‡ [70]	<b>Ours</b>
Inference Time (s)	2.13	7.79	12.98	5.7	1.2	1.57	0.718	0.6	2.665	0.1893	0.0302	0.3758	0.0267	0.0286

### 4.4 Empirical Analysis of Our Pipeline

Here we focus on evaluating the contribution of each component in our pipeline, and examining the internal performance of the pseudo-labels at different stages in training.

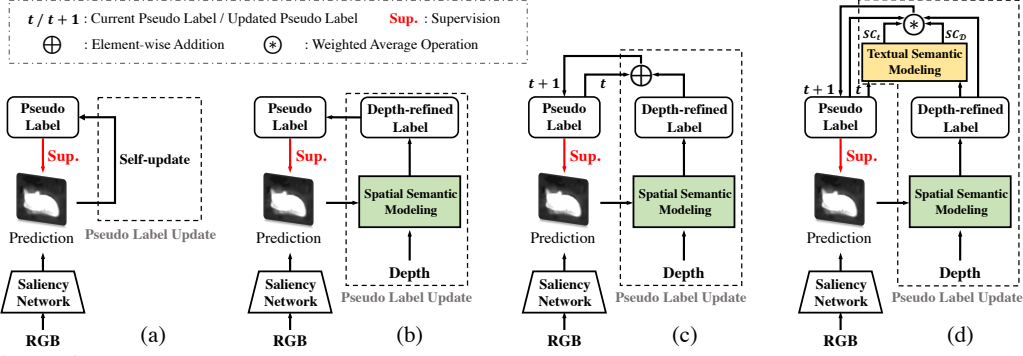


Figure 6: Diagrams of various label updating strategies used in our ablation study: (a) self-updating strategy, (b) SSM updating strategy, (c) historical moving average with equal weights, and finally (d) our JSM strategy.

Table 3: Ablation study of our pipeline.  $\uparrow$  ( $\downarrow$ ) denote performance gains (relative to backbone).

Model Setups	DUTL-Depth [55]			STERE [50]			NJUD [29]			NLPR [53]		
	$F_{\beta}^w$	$F_{\beta}$	$\mathcal{M}$	$F_{\beta}^w$	$F_{\beta}$	$\mathcal{M}$	$F_{\beta}^w$	$F_{\beta}$	$\mathcal{M}$	$F_{\beta}^w$	$F_{\beta}$	$\mathcal{M}$
Backbone trained on Pseudo Labels ( <i>i.e.</i> , ‘B’)	.512	.644	.167	.555	.666	.158	.510	.627	.186	.479	.570	.126
‘B’ trained on Pseudo Labels with CRF	.568	.670	.140	.601	.684	.135	.546	.642	.165	.585	.624	.094
‘B’ with Self-updating Strategy	.616	.697	.130	.643	.708	.123	.571	.673	.154	.607	.651	.087
‘B’ with Spatial Semantic Modeling	.629	.741	.116	.654	.738	.111	.576	.689	.147	.638	.698	.073
‘B’ with SSM and Historical Moving Average	.644	.750	.113	.674	.753	.105	.588	.698	.141	.664	.722	.068
‘B’ with Joint Semantic Mining ( <i>i.e.</i> , Ours)	<b>.698</b>	<b>.797</b>	<b>.093</b>	<b>.688</b>	<b>.778</b>	<b>.095</b>	<b>.599</b>	<b>.717</b>	<b>.133</b>	<b>.692</b>	<b>.770</b>	<b>.060</b>
	$\uparrow 36\%$	$\uparrow 24\%$	$\downarrow 44\%$	$\uparrow 24\%$	$\uparrow 17\%$	$\downarrow 40\%$	$\uparrow 17\%$	$\uparrow 14\%$	$\downarrow 28\%$	$\uparrow 44\%$	$\uparrow 35\%$	$\downarrow 52\%$

**Ablation analysis.** We present in Table 3 the ablation results of our pipeline on four benchmarks. To start with, we consider the backbone as the saliency network trained with initial pseudo labels. As our proposed SSM and TSM are gradually incorporated into the backbone to generate the depth-refined pseudo-labels and estimate their confidence scores for label updating, noticeable performance gains are consistently achieved in all datasets. The SSM significantly reduces the MAE metric and increases the F-measure score by 31% and 14.5% on average in four datasets, respectively. The TSM further boosts the saliency detection performance to a higher level where a significant amount of performance gains with 41% and 22.5% are finally achieved on MAE and F-measure metrics.

To further demonstrate the effectiveness of our SSM in exploiting the depth semantics to refine pseudo-labels, we retrain the saliency network using the self-updating strategy as in Fig. 6 (a). In this strategy, the saliency prediction with CRF are directly utilized to update pseudo-labels, at the end of each training round. As shown in Table 3 (3<sup>rd</sup> row vs. 4<sup>th</sup> row), when excluding saliency-oriented depth semantics captured by the SSM, the performance of model degrades greatly. This indicates our SSM can effectively suppress background noise and providing reliable training labels. Furthermore, we replace TSM with a heuristic historical moving average strategy where the pseudo-labels are assigned with equal weights as in Fig. 6 (c). Table 3 shows it achieves better performance than the backbone using SSM due to the consideration of historical information, but it is consistently inferior compared to our pipeline with the TSM module. These results suggest that our TSM can better estimate the confidence or quality of pseudo-labels and generate the trustworthy supervision signals. Meanwhile, we also provide the internal inspections of our approach in Fig. 4, in terms of the generation of pseudo-labels and their corresponding confidence scores, for better understanding.

In addition, we further discuss the effect of different update intervals in our pipeline when performing label updating. As listed in Table 4, the larger or smaller update interval leads to inferior performance due to the insufficient or excessive learning of models.

Table 4: Parameter analysis of the update interval  $\tau$  (epoch) in the JSM.

Interval ( $\tau$ )	STERE [50]		NJUD [29]		NLPR [53]	
	$F_{\beta}$	$\mathcal{M}$	$F_{\beta}$	$\mathcal{M}$	$F_{\beta}$	$\mathcal{M}$
$\tau = 1$	.769	.098	.703	.142	.766	.063
$\tau = 5$	.778	.095	.717	.133	.770	.060
$\tau = 10$	.758	.101	.695	.138	.763	.065

**Analysis of pseudo-labels.** We analyze the evolutions of pseudo-labels in our training process. Here the quality of pseudo-labels is measured using the ground-truth labels of the training set (for evaluation only). Visual evidences in Fig. 7 show that the quality of pseudo-labels is gradually improved as our JSM is performed. We can see that the initial pseudo-labels unfortunately tend to miss important salient parts as well as fine-grained details. By adopting our proposed joint semantic mining for label updating, the missing parts could be gradually retrieved, with the object silhouette



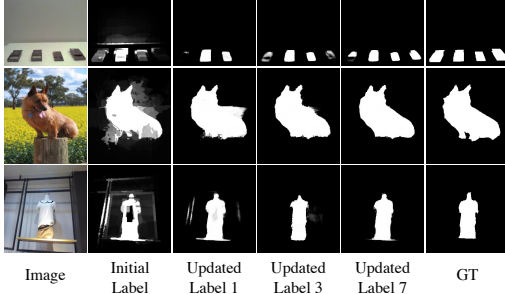


Figure 7: Visualization of updated labels.

Table 5: Ablation analysis of the supervised variant of our approach, where the human annotations (*i.e.*, ground truths) are used to train SOD models.

Model Setups (fully)	STERE [50]			NJUD [29]			NLPR [53]		
	$F_{\beta}^w$	$F_{\beta}$	$\mathcal{M}$	$F_{\beta}^w$	$F_{\beta}$	$\mathcal{M}$	$F_{\beta}^w$	$F_{\beta}$	$\mathcal{M}$
Backbone ( <i>i.e.</i> , 'B')	.858	.869	.045	.864	.871	.046	.865	.863	.028
'B' + SSM (Ours)	.876	.896	.039	.885	.906	.038	.892	.905	.022

also being refined. The final pseudo-label is closest to the true label, which could provide more reliable guiding signal for training the saliency network. Moreover, we present in Fig. 8 the error reduction curves of the updated pseudo-labels with our full pipeline (blue line) and SSM only (red line), respectively. It can be seen that only SSM is able to improve the quality of pseudo-labels by exploiting the useful depth semantics to refine pseudo-labels. Our JSM further boosts the performance by leveraging the textual semantics to integrate reliable pseudo-labels.

#### 4.5 Generalization Analysis

**Adapting to unsupervised setting.** Our approach can be adapted to unsupervised setting by using the architecture illustrated in Fig. 6 (b), *i.e.*, SSM, where only the depth semantics are mined to refine pseudo-labels without weak labels. The quantitative results in Table 1 show the effectiveness of our SSM in unsupervised setting.

**Adapting to fully-supervised setting.** Apart from the adaption to unsupervised setting, a variant of our approach can also be applied to fully-supervised RGB-D SOD scenario. This is achieved by modifying the generation of saliency-guided depth mask as  $\mathcal{D}_{mask} = \mathcal{S}_{GT} \otimes \mathcal{D}_{map}$  in Eq. 1, with  $\mathcal{S}_{GT}$  being the ground-truth annotations. The saliency network and depth network are trained by  $\mathcal{S}_{GT}$  and  $\mathcal{D}_{mask}$ , respectively. Then the background noise suppression block in SSM module is applied to obtain the final saliency. As ablated in Table 5, our fully-supervised variant achieves consistent performance improvement compared to the backbone trained on  $\mathcal{S}_{GT}$ . In addition, our results compare favorably with those of 21 fully-supervised RGB-D saliency models, as quantitatively shown in Table 6, and qualitatively illustrated in Fig. 9. Notice that, different from existing fully-supervised RGB-D SOD methods in designing the complicated cross-modal feature interaction strategy, our method directly exploits the learned depth semantics to promote saliency accuracy, which is simple yet effective and brand new for this field.

## 5 Failure Cases

Due to the sparsity of weak annotations, the network is usually difficult to identify the fine-grained object boundaries. As depicted in Fig. 10, although these models can effectively detect the salient objects, fine-grained details are still missing. A doable solution is to introduce auxiliary edge constraint during training. For example, the edge detection loss can be employed on the low-level features of the model, which could force the model to produce better features highlighting the object details [76, 41]. The edge maps can be generated by classical Canny operator [4] in an unsupervised manner.

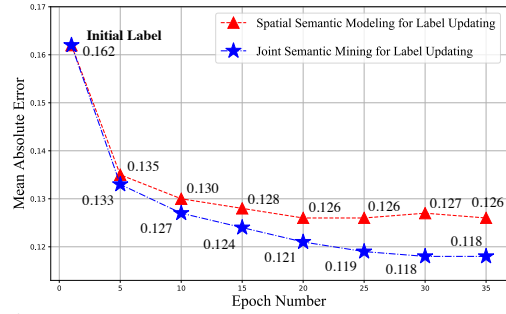


Figure 8: Error reduction plot of the updated pseudo-labels over iterations.

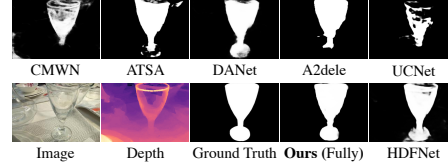


Figure 9: Visual comparison of fully supervised RGB-D SOD methods.

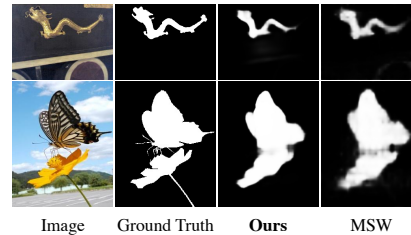


Figure 10: Failure cases of the existing weakly-supervised saliency methods.

Table 6: Quantitative comparison of existing fully-supervised RGB-D SOD methods. Notice that, in fully-supervised scenario, when evaluating the newly released DUTLF-Depth dataset, the specific setup used by [55] is adopted to make a fair comparison, *i.e.*, using a total of 2,985 training samples that contain 1,485 from NJUD, 700 from NLPR and 800 from DUTLF-Depth.

Method	DUTLF-Depth [55]				STERE [50]				NJUD [29]				NLPR [53]			
	$E_{\xi}$	$F_{\beta}^w$	$F_{\beta}$	$\mathcal{M}$	$E_{\xi}$	$F_{\beta}^w$	$F_{\beta}$	$\mathcal{M}$	$E_{\xi}$	$F_{\beta}^w$	$F_{\beta}$	$\mathcal{M}$	$E_{\xi}$	$F_{\beta}^w$	$F_{\beta}$	$\mathcal{M}$
CTMF <sup>TCyB'17</sup> [21]	.884	.690	.792	.097	.841	.747	.771	.086	.864	.732	.788	.085	.869	.691	.723	.056
DF <sup>TIP'17</sup> [58]	.842	.542	.748	.145	.691	.596	.742	.141	.818	.552	.744	.151	.838	.524	.682	.099
PCA <sup>CVPR'18</sup> [6]	.858	.696	.760	.100	.887	.801	.826	.064	.896	.811	.844	.059	.916	.772	.794	.044
TANet <sup>TIP'19</sup> [7]	.866	.712	.779	.093	.893	.804	.835	.060	.893	.812	.844	.061	.916	.789	.795	.041
PDNet <sup>ICME'19</sup> [95]	.861	.650	.757	.112	.880	.799	.813	.071	.890	.798	.832	.062	.876	.659	.740	.064
MMCI <sup>PR'19</sup> [8]	.855	.636	.753	.113	.873	.757	.829	.068	.878	.749	.813	.079	.871	.688	.729	.059
CPFP <sup>CVPR'19</sup> [86]	.814	.644	.736	.099	.912	.808	.830	.051	.895	.837	.850	.053	.924	.820	.822	.036
DMRA <sup>ICCV'19</sup> [55]	.927	.858	.883	.048	.923	.841	.876	.049	.908	.853	.872	.051	.942	.845	.855	.031
SSF <sup>CVPR'20</sup> [83]	.946	.894	.914	.034	.921	.850	.867	.046	.913	.871	.886	.043	.949	.874	.875	.026
A2dele <sup>CVPR'20</sup> [56]	.924	.864	.890	.043	.915	.855	.874	.044	.897	.851	.874	.051	.945	.867	.878	.028
JL-DCF <sup>CVPR'20</sup> [19]	.931	.863	.883	.043	.919	.857	.869	.040	-	-	-	-	.954	.882	.878	<b>.022</b>
S2MA <sup>CVPR'20</sup> [42]	.921	.861	.866	.044	.907	.825	.855	.051	-	-	-	-	.938	.852	.853	.030
UCNet <sup>CVPR'20</sup> [74]	.903	.821	.856	.056	.922	.867	.885	<b>.039</b>	-	-	-	-	.953	.878	.890	.025
PGAR <sup>ECCV'20</sup> [9]	.944	.889	.914	.035	.919	.856	.880	.041	.915	.871	.893	.042	.955	.881	.885	.024
D3Net <sup>NNLS'20</sup> [17]	.847	.668	.756	.097	.920	.845	.855	.046	.913	.860	.863	.047	.943	.854	.857	.030
CMWN <sup>ECCV'20</sup> [37]	.916	.831	.866	.056	.917	.847	.869	.043	.910	.855	.878	.047	.940	.856	.859	.029
BBSNet <sup>ECCV'20</sup> [18]	.833	.663	.774	.120	.925	.858	.885	.041	.924	.884	.902	<b>.035</b>	.952	.879	.882	.023
DANet <sup>ECCV'20</sup> [91]	.925	.847	.884	.047	.914	.830	.858	.047	-	-	-	-	.949	.858	.871	.028
FRDT <sup>ACMM'20</sup> [84]	.941	.878	.902	.039	.925	.858	.872	.042	.917	.862	.879	.048	.946	.863	.868	.029
ATSA <sup>ECCV'20</sup> [78]	.947	.901	.918	.032	.919	.866	.874	.040	.921	.883	.893	.040	.945	.867	.876	.028
HDFNet <sup>ECCV'20</sup> [51]	.934	.865	.892	.040	.925	.863	.879	.040	.915	.879	.893	.038	.948	.869	.878	.027
<b>Ours (Fully Sup.)</b>	<b>.949</b>	<b>.908</b>	<b>.934</b>	<b>.030</b>	<b>.929</b>	<b>.876</b>	<b>.896</b>	<b>.039</b>	<b>.926</b>	<b>.885</b>	<b>.906</b>	.038	<b>.959</b>	<b>.892</b>	<b>.905</b>	<b>.022</b>

## 6 Conclusion

To tackle the new problem of weakly-supervised RGB-D salient object detection, we propose in this paper an end-to-end approach based on iterative updates of the internal pseudo-labels. This allows us to leverage depth information in eliminating non-salient background noises and generating reliable depth-refined pseudo-labels. Textual semantics is incorporated in the fill-in-the-blank fashion, which is used to estimate the confidence scores of pseudo-labels. Extensive experiments demonstrate the effectiveness and efficiency of our approach. In addition, our method is very generic and can be easily adapted to fully-supervised and unsupervised paradigms. In these scenarios, our variants also obtain superior performance over the existing state-of-the-art dedicated methods.

**Acknowledgement.** This research was supported by the University of Alberta Start-up Grant, UAHJIC Grants, and NSERC Discovery Grants (No. RGPIN-2019-04575). The authors are grateful to the anonymous reviewers for their valuable suggestions in improving the quality of the paper.

## References

- [1] Radhakrishna Achanta, Sheila Hemami, Francisco Estrada, and Sabine Susstrunk. Frequency-tuned salient region detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1597–1604, 2009.
- [2] Qi Bi, Shuang Yu, Wei Ji, Cheng Bian, Lijun Gong, Hanruo Liu, Kai Ma, and Yefeng Zheng. Local-global dual perception based deep multiple instance learning for retinal disease classification. In *MICCAI*, pages 55–64. Springer, 2021.
- [3] Ali Borji, Ming-Ming Cheng, Huaizu Jiang, and Jia Li. Salient object detection: A benchmark. *IEEE Transactions on Image Processing*, 24(12):5706–5722, 2015.
- [4] John Canny. A computational approach to edge detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 8(6):679–698, 1986.
- [5] Chenglizhao Chen, Jipeng Wei, Chong Peng, and Hong Qin. Depth-quality-aware salient object detection. *IEEE Transactions on Image Processing*, 30:2350–2363, 2021.
- [6] Hao Chen and Youfu Li. Progressively complementarity-aware fusion network for RGB-D salient object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3051–3060, 2018.
- [7] Hao Chen and Youfu Li. Three-stream attention-aware network for RGB-D salient object detection. *IEEE Transactions on Image Processing*, 28(6):2825–2835, 2019.

- [8] Hao Chen, Youfu Li, and Dan Su. Multi-modal fusion network with multi-scale multi-path and cross-modal interactions for RGB-D salient object detection. *Pattern Recognition*, 86:376–385, 2019.
- [9] Shuhan Chen and Yun Fu. Progressively guided alternate refinement network for RGB-D salient object detection. In *Proceedings of the European Conference on Computer Vision*, 2020.
- [10] Zuyao Chen, Runmin Cong, Qianqian Xu, and Qingming Huang. DPANet: Depth potentiality-aware gated attention network for RGB-D salient object detection. *IEEE Transactions on Image Processing*, 2020.
- [11] Yupeng Cheng, Huazhu Fu, Xingxing Wei, Jiangjian Xiao, and Xiaochun Cao. Depth enhanced saliency detection method. In *Proceedings of International Conference on Internet Multimedia Computing and Service*, pages 23–27, 2014.
- [12] Runmin Cong, Jianjun Lei, Changqing Zhang, Qingming Huang, Xiaochun Cao, and Chunping Hou. Saliency detection for stereoscopic images based on depth confidence analysis and multiple cues fusion. *IEEE Signal Processing Letters*, 23(6):819–823, 2016.
- [13] Karthik Desingh, K Madhava Krishna, Deepu Rajan, and CV Jawahar. Depth really matters: Improving visual salient region detection with depth. In *British Machine Vision Conference*, 2013.
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina N. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4171–4186, 2018.
- [15] Deng-Ping Fan, Ming-Ming Cheng, Yun Liu, Tao Li, and Ali Borji. Structure-measure: A new way to evaluate foreground maps. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4548–4557, 2017.
- [16] Deng-Ping Fan, Cheng Gong, Yang Cao, Bo Ren, Ming-Ming Cheng, and Ali Borji. Enhanced-alignment measure for binary foreground map evaluation. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 698–704, 2018.
- [17] Deng-Ping Fan, Zheng Lin, Zhao Zhang, Menglong Zhu, and Ming-Ming Cheng. Rethinking RGB-D salient object detection: Models, data sets, and large-scale benchmarks. *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [18] Deng-Ping Fan, Yingjie Zhai, Ali Borji, Jufeng Yang, and Ling Shao. BBS-Net: RGB-D salient object detection with a bifurcated backbone strategy network. In *Proceedings of the European Conference on Computer Vision*, 2020.
- [19] Keren Fu, Deng-Ping Fan, Ge-Peng Ji, and Qijun Zhao. JL-DCF: Joint learning and densely-cooperative fusion framework for RGB-D salient object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3052–3062, 2020.
- [20] Jingfan Guo, Tongwei Ren, and Jia Bei. Salient object detection for RGB-D image via saliency evolution. In *IEEE International Conference on Multimedia and Expo*, pages 1–6, 2016.
- [21] Junwei Han, Hao Chen, Nian Liu, Chenggang Yan, and Xuelong Li. CNNs-based RGB-D saliency detection via cross-view transfer and multiview fusion. *IEEE Transactions on Cybernetics*, 48(11):3171–3183, 2017.
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [23] Qibin Hou, Ming-Ming Cheng, Xiaowei Hu, Ali Borji, Zhuowen Tu, and Philip HS Torr. Deeply supervised salient object detection with short connections. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3203–3212, 2017.
- [24] Kuang-Jui Hsu, Yen-Yu Lin, and Yung-Yu Chuang. Weakly supervised salient object detection by learning a classifier-driven map generator. *IEEE Transactions on Image Processing*, 28(11):5435–5449, 2019.
- [25] Wei Ji, Wenting Chen, Shuang Yu, Kai Ma, Li Cheng, Linlin Shen, and Yefeng Zheng. Uncertainty quantification for medical image segmentation using dynamic label factor allocation among multiple raters. In *MICCAI on QUBIQ Workshop*, 2020.
- [26] Wei Ji, Jingjing Li, Shuang Yu, Miao Zhang, Yongri Piao, Shunyu Yao, Qi Bi, Kai Ma, Yefeng Zheng, Huchuan Lu, and Li Cheng. Calibrated RGB-D salient object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9471–9481, June 2021.

- [27] Wei Ji, Jingjing Li, Miao Zhang, Yongri Piao, and Huchuan Lu. Accurate RGB-D salient object detection via collaborative learning. In *Proceedings of the European Conference on Computer Vision*, pages 52–69, 2020.
- [28] Wei Ji, Shuang Yu, Junde Wu, Kai Ma, Cheng Bian, Qi Bi, Jingjing Li, Hanruo Liu, Li Cheng, and Yefeng Zheng. Learning calibrated medical image segmentation via multi-rater agreement modeling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12341–12351, June 2021.
- [29] Ran Ju, Ling Ge, Wenjing Geng, Tongwei Ren, and Gangshan Wu. Depth saliency based on anisotropic center-surround difference. In *IEEE International Conference on Image Processing*, pages 1115–1119, 2014.
- [30] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3128–3137, 2015.
- [31] ByoungChul Ko, Soo Yeong Kwak, and Hyeran Byun. SVM-based salient region(s) extraction method for image retrieval. In *Proceedings of the 17th International Conference on Pattern Recognition*, volume 2, pages 977–980, 2004.
- [32] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105, 2012.
- [33] Chongyi Li, Seed Anwar, Junhui Hou, Runmin Cong, Chunle Guo, and Wenqi Ren. Underwater image enhancement via medium transmission-guided multi-color space embedding. *IEEE Transactions on Image Processing*, 2021.
- [34] Chongyi Li, Runmin Cong, Junhui Hou, Sanyi Zhang, Yue Qian, and Sam Kwong. Nested network with two-stream pyramid for salient object detection in optical remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 57(11):9156–9166, 2019.
- [35] Chongyi Li, Runmin Cong, Sam Kwong, Junhui Hou, Huazhu Fu, Guopu Zhu, Dingwen Zhang, and Qingming Huang. ASIF-Net: Attention steered interweave fusion network for RGB-D salient object detection. *IEEE Transactions on Cybernetics*, PP(99):1–13, 2020.
- [36] Chongyi Li, Runmin Cong, Yongri Piao, Qianqian Xu, and Chen Change Loy. RGB-D salient object detection with cross-modality modulation and selection. In *Proceedings of the European Conference on Computer Vision*, 2020.
- [37] Gongyang Li, Zhi Liu, Linwei Ye, Yang Wang, and Haibin Ling. Cross-modal weighting network for RGB-D salient object detection. In *Proceedings of the European Conference on Computer Vision*, 2020.
- [38] Guanbin Li, Yuan Xie, and Liang Lin. Weakly supervised salient object detection using image labels. In *The AAAI Conference on Artificial Intelligence*, 2018.
- [39] Xiaohui Li, Huchuan Lu, Lihe Zhang, Xiang Ruan, and Ming-Hsuan Yang. Saliency detection via dense and sparse reconstruction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2976–2983, 2013.
- [40] Fangfang Liang, Lijuan Duan, Wei Ma, Yuanhua Qiao, Zhi Cai, and Laiyun Qing. Stereoscopic saliency model using contrast and depth-guided-background prior. *Neurocomputing*, 275:2227–2238, 2018.
- [41] Jiang-Jiang Liu, Qibin Hou, Ming-Ming Cheng, Jiashi Feng, and Jianmin Jiang. A simple pooling-based design for real-time salient object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3917–3926, 2019.
- [42] Nian Liu, Ni Zhang, and Junwei Han. Learning selective self-mutual attention for RGB-D saliency detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 13756–13765, 2020.
- [43] Yunqiu Lv, Jing Zhang, Yuchao Dai, Aixuan Li, Bowen Liu, Nick Barnes, and Deng-Ping Fan. Simultaneously localize, segment and rank the camouflaged objects. In *CVPR*, 2021.
- [44] Ran Margolin, Lihi Zelnik-Manor, and Ayellet Tal. How to evaluate foreground maps? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2014.
- [45] Tam Nguyen, Maximilian Dax, Chaithanya Kumar Mummadi, Nhung Ngo, Thi Hoai Phuong Nguyen, Zhongyu Lou, and Thomas Brox. DeepUSPS: Deep robust unsupervised saliency prediction via self-supervision. In *Advances in Neural Information Processing Systems*, pages 204–214, 2019.

- [46] Munan Ning, Cheng Bian, Donghuan Lu, Hong-Yu Zhou, Shuang Yu, Chenglang Yuan, Yang Guo, Yaohua Wang, Kai Ma, and Yefeng Zheng. A macro-micro weakly-supervised framework for as-oc tissue segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 725–734. Springer, 2020.
- [47] Munan Ning, Cheng Bian, Dong Wei, Shuang Yu, Chenglang Yuan, Yaohua Wang, Yang Guo, Kai Ma, and Yefeng Zheng. A new bidirectional unsupervised domain adaptation segmentation framework. In *International Conference on Information Processing in Medical Imaging*, pages 492–503. Springer, 2021.
- [48] Munan Ning, Cheng Bian, Chenglang Yuan, Kai Ma, and Yefeng Zheng. Ensembled resnet for anatomical brain barriers segmentation. *Segmentation, Classification, and Registration of Multi-modality Medical Imaging Data*, 12587:27, 2021.
- [49] Munan Ning, Donghuan Lu, Dong Wei, Cheng Bian, Chenglang Yuan, Shuang Yu, Kai Ma, and Yefeng Zheng. Multi-anchor active domain adaptation for semantic segmentation. *arXiv preprint arXiv:2108.08012*, 2021.
- [50] Yuzhen Niu, Yujie Geng, Xueqing Li, and Feng Liu. Leveraging stereopsis for saliency analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 454–461, 2012.
- [51] Youwei Pang, Lihe Zhang, Xiaoqi Zhao, and Huchuan Lu. Hierarchical dynamic filtering network for RGB-D salient object detection. In *Proceedings of the European Conference on Computer Vision*, 2020.
- [52] Youwei Pang, Xiaoqi Zhao, Lihe Zhang, and Huchuan Lu. Multi-scale interactive network for salient object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9413–9422, 2020.
- [53] Houwen Peng, Bing Li, Weihua Xiong, Weiming Hu, and Rongrong Ji. RGBD salient object detection: a benchmark and algorithms. In *Proceedings of the European Conference on Computer Vision*, pages 92–109, 2014.
- [54] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1532–1543, 2014.
- [55] Yongri Piao, Wei Ji, Jingjing Li, Miao Zhang, and Huchuan Lu. Depth-induced multi-scale recurrent attention network for saliency detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7254–7263, 2019.
- [56] Yongri Piao, Zhengkun Rong, Miao Zhang, Weisong Ren, and Huchuan Lu. A2dele: Adaptive and attentive depth distiller for efficient RGB-D salient object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9060–9069, 2020.
- [57] Yao Qin, Huchuan Lu, Yiqun Xu, and He Wang. Saliency detection via cellular automata. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 110–119, 2015.
- [58] Liangqiong Qu, Shengfeng He, Jiawei Zhang, Jiandong Tian, Yandong Tang, and Qingxiong Yang. RGBD salient object detection via deep fusion. *IEEE Transactions on Image Processing*, 26(5):2274–2285, 2017.
- [59] Jianqiang Ren, Xiaojin Gong, Lu Yu, Wenhui Zhou, and Michael Ying Yang. Exploiting global priors for RGB-D saliency detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 25–32, 2015.
- [60] Riku Shigematsu, David Feng, Shaodi You, and Nick Barnes. Learning RGB-D salient object detection using background enclosure, depth contrast, and top-down features. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 2749–2757, 2017.
- [61] Wei-Chih Tu, Shengfeng He, Qingxiong Yang, and Shao-Yi Chien. Real-time salient object detection with a minimum spanning tree. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2334–2342, 2016.
- [62] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, pages 5998–6008, 2017.
- [63] Lijun Wang, Huchuan Lu, Yifan Wang, Mengyang Feng, Dong Wang, Baocai Yin, and Xiang Ruan. Learning to detect salient objects with image-level supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 136–145, 2017.

- [64] Zhe Wu, Li Su, and Qingming Huang. Cascaded partial decoder for fast and accurate salient object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3907–3916, 2019.
- [65] Cheng Yan, Guansong Pang, Lei Wang, Jile Jiao, Xuetao Feng, Chunhua Shen, and Jingjing Li. Bv-person: A large-scale dataset for bird-view person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10943–10952, October 2021.
- [66] Qingsong Yao, Zecheng He, Hu Han, and S Kevin Zhou. Miss the point: Targeted adversarial attack on multiple landmark detection. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 692–702. Springer, 2020.
- [67] Qingsong Yao, Quan Quan, Li Xiao, and S Kevin Zhou. One-shot medical landmark detection. *arXiv preprint arXiv:2103.04527*, 2021.
- [68] Qingsong Yao, Li Xiao, Peihang Liu, and S Kevin Zhou. Label-free segmentation of covid-19 lesions in lung ct. *IEEE Transactions on Medical Imaging*, 2021.
- [69] Siyue Yu, Bingfeng Zhang, Jimin Xiao, and Eng Gee Lim. Structure-consistent weakly supervised salient object detection with local saliency coherence. In *The AAAI Conference on Artificial Intelligence*, 2021.
- [70] Yu Zeng, Yunzhi Zhuge, Huchuan Lu, Lihe Zhang, Mingyang Qian, and Yizhou Yu. Multi-source weak supervision for saliency detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6074–6083, 2019.
- [71] Dingwen Zhang, Junwei Han, Gong Cheng, and Ming-Hsuan Yang. Weakly supervised object localization and detection: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [72] Dingwen Zhang, Junwei Han, and Yu Zhang. Supervision by fusion: Towards unsupervised learning of deep salient object detector. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4048–4056, 2017.
- [73] Dingwen Zhang, Junwei Han, Yu Zhang, and Dong Xu. Synthesizing supervision for learning deep saliency network without human annotation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(7):1755–1769, 2019.
- [74] Jing Zhang, Deng-Ping Fan, Yuchao Dai, Saeed Anwar, Fatemeh Sadat Saleh, Tong Zhang, and Nick Barnes. UC-Net: Uncertainty inspired RGB-D saliency detection via conditional variational autoencoders. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8582–8591, 2020.
- [75] Jing Zhang, Nick Barnes Jianwen Xie, and Ping Li. Learning generative vision transformer with energy-based latent space for saliency prediction. In *NeurIPS*, 2021.
- [76] Jing Zhang, Xin Yu, Aixuan Li, Peipei Song, Bowen Liu, and Yuchao Dai. Weakly-supervised salient object detection via scribble annotations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12546–12555, 2020.
- [77] Jing Zhang, Tong Zhang, Yuchao Dai, Mehrtash Harandi, and Richard Hartley. Deep unsupervised saliency detection: A multiple noisy labeling perspective. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9029–9038, 2018.
- [78] Miao Zhang, Sun Xiao Fei, Jie Liu, Shuang Xu, Yongri Piao, and Huchuan Lu. Asymmetric two-stream architecture for accurate RGB-D saliency detection. In *Proceedings of the European Conference on Computer Vision*, 2020.
- [79] Miao Zhang, Wei Ji, Yongri Piao, Jingjing Li, Yu Zhang, Shuang Xu, and Huchuan Lu. LFNet: Light field fusion network for salient object detection. *IEEE Transactions on Image Processing*, 29:6276–6287, 2020.
- [80] Miao Zhang, Jingjing Li, Wei Ji, Yongri Piao, and Huchuan Lu. Memory-oriented decoder for light field salient object detection. In *Advances in Neural Information Processing Systems*, pages 896–906, 2019.
- [81] Miao Zhang, Jie Liu, Yifei Wang, Yongri Piao, Shunyu Yao, Wei Ji, Jingjing Li, Huchuan Lu, and Zhongxuan Luo. Dynamic context-sensitive filtering network for video salient object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, 2021.
- [82] Miao Zhang, Tingwei Liu, Yongri Piao, Shunyu Yao, and Huchuan Lu. Auto-msfnet: Search multi-scale fusion network for salient object detection. In *ACM Multimedia Conference 2021*, 2021.

- [83] Miao Zhang, Weisong Ren, Yongri Piao, Zhengkun Rong, and Huchuan Lu. Select, supplement and focus for RGB-D saliency detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3472–3481, 2020.
- [84] Miao Zhang, Yu Zhang, Yongri Piao, Beiqi Hu, and Huchuan Lu. Feature reintegration over differential treatment: A top-down and adaptive fusion network for RGB-D salient object detection. In *ACM International Conference on Multimedia*, 2020.
- [85] Qijian Zhang, Runmin Cong, Junhui Hou, Chongyi Li, and Yao Zhao. Coadnet: Collaborative aggregation-and-distribution networks for co-salient object detection. *NeurIPS*, 2020.
- [86] Jia-Xing Zhao, Yang Cao, Deng-Ping Fan, Ming-Ming Cheng, Xuan-Yi Li, and Le Zhang. Contrast prior and fluid pyramid integration for RGBD salient object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3927–3936, 2019.
- [87] Jiawei Zhao, Yifan Zhao, Jia Li, and Xiaowu Chen. Is depth really necessary for salient object detection? In *ACM International Conference on Multimedia*, 2020.
- [88] Xiaoqi Zhao, Youwei Pang, Jiaying Yang, Lihe Zhang, and Huchuan Lu. Multi-source fusion and automatic predictor selection for zero-shot video object segmentation. *arXiv preprint arXiv:2108.05076*, 2021.
- [89] Xiaoqi Zhao, Youwei Pang, Lihe Zhang, Huchuan Lu, and Lei Zhang. Suppress and balance: A simple gated network for salient object detection. In *Proceedings of the European Conference on Computer Vision*, 2020.
- [90] Xiaoqi Zhao, Lihe Zhang, and Huchuan Lu. Automatic polyp segmentation via multi-scale subtraction network. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 120–130. Springer, 2021.
- [91] Xiaoqi Zhao, Lihe Zhang, Youwei Pang, Huchuan Lu, and Lei Zhang. A single stream network for robust and real-time RGB-D salient object detection. In *Proceedings of the European Conference on Computer Vision*, 2020.
- [92] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2921–2929, 2016.
- [93] Tao Zhou, Deng-Ping Fan, Ming-Ming Cheng, Jianbing Shen, and Ling Shao. RGB-D salient object detection: A survey. *Computational Visual Media*, pages 1–33, 2021.
- [94] Tao Zhou, Huazhu Fu, Geng Chen, Yi Zhou, Deng-Ping Fan, and Ling Shao. Specificity-preserving rgb-d saliency detection. *Proceedings of the IEEE International Conference on Computer Vision*, 2021.
- [95] Chunbiao Zhu, Xing Cai, Kan Huang, Thomas H Li, and Ge Li. PDNet: Prior-model guided depth-enhanced network for salient object detection. In *IEEE International Conference on Multimedia and Expo*, pages 199–204, 2019.
- [96] Chunbiao Zhu, Ge Li, Xiaoqiang Guo, Wenmin Wang, and Ronggang Wang. A multilayer backpropagation saliency detection algorithm based on depth mining. In *International Conference on Computer Analysis of Images and Patterns*, pages 14–23, 2017.
- [97] Chunbiao Zhu, Ge Li, Wenmin Wang, and Ronggang Wang. An innovative salient object detection using center-dark channel prior. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 1509–1515, 2017.
- [98] Wangjiang Zhu, Shuang Liang, Yichen Wei, and Jian Sun. Saliency optimization from robust background detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2814–2821, 2014.