



# Cross-Modal Weighting Network for RGB-D Salient Object Detection

Gongyang Li<sup>1</sup> , Zhi Liu<sup>1</sup> , Linwei Ye<sup>2</sup> , Yang Wang<sup>2,4</sup> ,  
and Haibin Ling<sup>3</sup> 

<sup>1</sup> Shanghai University, Shanghai, China  
ligongyang@shu.edu.cn, liuzhi@staff.shu.edu.cn

<sup>2</sup> University of Manitoba, Winnipeg, Canada  
{ye13,ywang}@cs.umanitoba.ca

<sup>3</sup> Stony Brook University, Stony Brook, NY, USA  
hling@cs.stonybrook.edu

<sup>4</sup> Huawei Technologies Canada, Markham, Canada  
<https://github.com/MathLee/CMWNet>

**Abstract.** Depth maps contain geometric clues for assisting Salient Object Detection (SOD). In this paper, we propose a novel Cross-Modal Weighting (CMW) strategy to encourage comprehensive interactions between RGB and depth channels for RGB-D SOD. Specifically, three RGB-depth interaction modules, named CMW-L, CMW-M and CMW-H, are developed to deal with respectively low-, middle- and high-level cross-modal information fusion. These modules use Depth-to-RGB Weighing (DW) and RGB-to-RGB Weighting (RW) to allow rich cross-modal and cross-scale interactions among feature layers generated by different network blocks. To effectively train the proposed Cross-Modal Weighting Network (CMWNet), we design a composite loss function that summarizes the errors between intermediate predictions and ground truth over different scales. With all these novel components working together, CMWNet effectively fuses information from RGB and depth channels, and meanwhile explores object localization and details across scales. Thorough evaluations demonstrate CMWNet consistently outperforms 15 state-of-the-art RGB-D SOD methods on seven popular benchmarks.

**Keywords:** RGB-D salient object detection · Cross-Modal Weighting · Depth-to-RGB weighting · RGB-to-RGB weighting

## 1 Introduction

Salient object detection (SOD) aims to pick the regions/objects in an image that are most attractive to human visual attention. It has a wide range of applications as summarized in recent surveys [1, 2, 39]. Most existing SOD solutions

take as input an RGB image (or video), which is convenient in many application scenarios, but may suffer from challenges such as low contrast and disturbing background. Alternatively, one can seek help from the depth information typically provided with an RGB-D input. In fact, with the popularity of depth sensors/devices, RGB-D SOD has received extensive attention recently, and numerous approaches [4, 8, 13, 16, 20, 25, 28, 31, 36, 42] have been proposed to extract salient objects from paired RGB images and depth maps.

Starting with the first stereoscopic image SOD dataset STEREO [28], traditional RGB-D SOD methods mainly apply contrast cue [14, 15, 35], fusion framework [19, 29, 36, 37] and measure strategy [9, 16, 25] to extract the complementary information in depth maps. These well-designed hand-crafted features-based methods, which are influenced by RGB SOD solutions, have achieved remarkable results. However, salient objects in the generated saliency maps are sometimes blocky because of inappropriate over-segmentation, while salient objects may be confused by complex scenes.

Recently, with the rapid development of deep learning, convolutional neural networks (CNNs) have shown strong dominance in many computer vision problems. Many CNN-based RGB-D SOD methods have been proposed and greatly outperformed traditional ones. Early CNN-based methods [31, 33] feed the superpixel-based hand-crafted features of RGB-D pairs into CNNs, but their results are still patch-based. Subsequent methods instead assign saliency values for each pixel based on the RGB image and depth map in an end-to-end manner. Among these methods, the two-stream architecture [4, 6, 10, 20, 38, 44] fuses cross-modal features/saliency maps in the middle/late stage, while the single-stream architecture [13, 26] directly handles RGB-D pairs. These methods, despite achieving great performance gain, do not take full advantage of rich interactive information between different modalities and scales of CNN blocks.

Motivated by the above observation, in this paper, we propose a novel *Cross-Modal Weighting Network* (CMWNet) that significantly improves RGB-depth interactions, and hence boosts RGB-D SOD performances as demonstrated in our thorough experiments. Our key idea is to jointly explore the information carried by both RGB and depth channels, and to encourage cross-modal and cross-scale RGB-depth interactions among different CNN feature blocks. This way, our algorithm can capture both microscopic details carried by shallow blocks and macroscopic object location information carried by deep blocks. CMWNet adopts a three-level representation, capturing low-, middle-, and high-level information respectively; and multiple blocks at different scales are allowed to be within a level. The cross-modal cross-scale interactions are modeled through the novel Cross-Modal Weighting (CMW) modules to highlight salient objects.

In particular, we propose three CMW modules, CMW-L, CMW-M and CMW-H. For low- and middle-level parts, CMW-L and CMW-M are used to enhance salient object details in a cross-scale manner. For high-level part, CMW-H is used to enhance salient object localization, which plays a crucial role in subsequent prediction of salient objects. The key components in these CMW modules are the proposed Depth-to-RGB Weighting (DW) and RGB-to-RGB

Weighting (RW) operations that enhance RGB features in each channel based on corresponding response maps. In addition to the encoder, a three-level decoder is designed to connect the three-level enhanced features to predict the final salient objects. In this way, the proposed CMWNet effectively exploits the properties of CNN features and strengthens the cross-modal and cross-scale interactions, resulting in excellent performance.

Our major contributions are summarized as follows:

- We explore the complex complementarity between RGB image and depth map in a three-level encoder-decoder structure, and propose a novel *Cross-Modal Weighting Network* (CMWNet) to encourage the cross-modal and cross-scale interactions, boosting the performance of RGB-D SOD.
- We propose three novel RGB-depth interaction modules to effectively enhance both salient object details (CMW-L and CMW-M) and salient object localization (CMW-H).
- Extensive experiments on seven popular public datasets under six commonly used evaluation metrics show that the proposed method achieves the best performance compared with 15 state-of-the-art RGB-D SOD methods.

## 2 Related Work

**Traditional RGB-D SOD.** Starting from the first work for saliency detection [21], the contrast-based approaches are the mainstream for saliency detection. This trend has spread to traditional RGB-D SOD. Numerous contrast-based RGB-D SOD methods have been proposed, such as disparity contrast [28], depth contrast [7, 14, 15, 35], and multi-contextual contrast [29]. Song *et al.* [36] employed the multi-scale fusion to merge saliency maps to obtain the final RGB-D saliency map, which is similar to methods based on two-stream saliency fusion [29] and multiple-cues fusion [19, 37]. By adopting the objectness measure [25], depth confidence measure [9] and salient structure measure [16], the performance gets clear improvement. Besides, other methods (*i.e.*, global prior [32], cellular automata [18], transformation strategy [8]) have been proposed for RGB-D SOD. However, these traditional methods are often based on superpixels, regions and patches, which cause saliency maps to appear blocky and saliency values to be scattered.

**CNN-Based RGB-D SOD.** In recent years, numerous CNN-based RGB-D SOD methods [4–6, 10, 13, 20, 26, 31, 33, 38, 42, 44] have been proposed. As pioneering work based on CNNs, Qu *et al.* [31] fed the superpixel-based RGB-D saliency features into a five-layer CNN. Shigematsu *et al.* [33] sent ten superpixel-based depth features to a network. Being patch-based methods, these methods sometimes generate results that appear blocky. To overcome the limitation, Han *et al.* [20] proposed a transfer and fusion based network to predict pixel-level saliency values. The single-stream architecture [13, 26] adopts a straightforward way to handle the four-channel RGB-D pair. This architecture does not effectively capture the cross-modal interactions between the RGB image and the

depth map, so the performance depends largely on the network structure rather than the cross-modal interactions. The two-stream architecture employs two separate networks to extract features [4, 6, 20, 44] and saliency maps [10, 38], and then fuse them with various strategies. Some works [10, 20, 38, 44] only fuse saliency maps and high-level features. As a result, they do not capture more complex cross-modal interactions at other levels of the network. Some other works [4, 5] consider cross-modal CNN features, but the same module is used to process cross-modal CNN features at different blocks. Consequently, these methods ignore the different properties of CNN features at different blocks and cannot provide specific enhancements to object details and object localization.

In this work, we propose a novel three-level CMWNet to encourage interactions between RGB and depth channels and propose several modules to treat differently detail features and localization features carried in CNN feature blocks at various scales. Moreover, we process CNN features in a cross-modal and cross-scale manner to effectively capture the interactions across modalities and scales. Thus, our network can accurately enhance the details and localization of salient objects in the encoder part and precisely infer salient objects in the three-level decoder.

### 3 Proposed Method

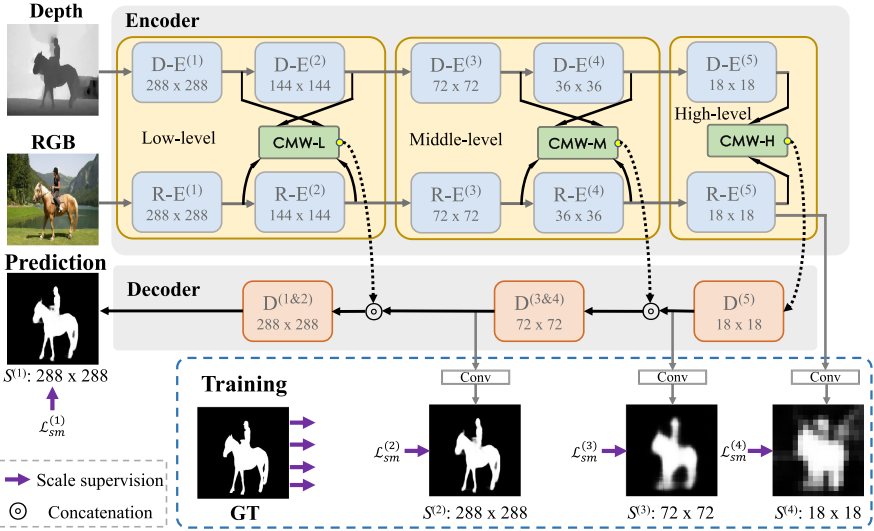
In this section, we start with the overview of Cross-Modal Weighting Network (CMWNet) (Sect. 3.1). In Sect. 3.2, we provide the details of low- and middle-level cross-modal weighting modules, *i.e.* CMW-L and CMW-M, and then in Sect. 3.3 we introduce the high-level cross-modal weighting module CMW-H. Finally, we describe the implementation details in Sect. 3.4.

#### 3.1 Network Overview and Motivation

The proposed CMWNet follows a three-level Siamese encoder-decoder structure, as summarized in Fig. 1.

**Three-Level Encoder.** We adopt the VGG16 [34] as the backbone. The depth map branch and the RGB image branch share the same weights. In the Siamese encoder part, five CNN blocks of the depth map and the RGB image are denoted as D-E<sup>(*l*)</sup> and R-E<sup>(*l*)</sup> (*l* ∈ {1, 2, 3, 4, 5}) is the block index), respectively. Considering the unique properties of features, we divide the first and second CNN blocks into low-level part, the third and fourth CNN blocks into middle-level part, and the last CNN block into high-level part.

**Low- and Middle-Level Cross-Modal Weighting Modules.** The weighting mechanism [43] is an extended version of attention mechanism, and it aims to modulate features of each channel according to particular response maps. The abundant geometric knowledge of depth maps is helpful to provide object details and object localization for SOD. We novelly extend the weighting mechanism with cross-modal information (*i.e.* RGB image and depth map), and propose

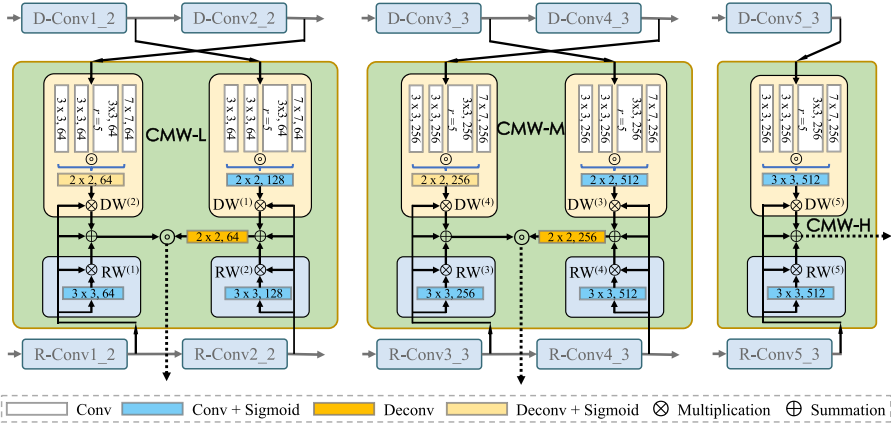


**Fig. 1.** Illustration of the proposed CMWNet. For both RGB and depth channel, the Siamese encoder network is employed to extract feature blocks organized in three levels. Three Cross-Modal Weighting (CMW) modules, CMW-L, CMW-M and CMW-H, are proposed to capture the interactions at corresponding level, and provide inputs for the decoder. The decoder progressively aggregates all the cross-modal cross-scale information for the final prediction. For training, multi-scale pixel-level supervision for intermediate predictions are utilized.

cross-modal RGB-depth interaction modules, which adopt weighting mechanism to reweight the RGB features based on depth response maps and RGB response maps to focus on salient objects.

Considering that the low- and middle-level parts carry abundant information about object details, we treat the features in these two levels as responsible for object details enhancement, and propose CMW-L and CMW-M. Each of the low- and middle-level contains the two adjacent CNN blocks, one contains relatively macroscopic information while the other relatively microscopic. To balance these two types of information within a level, we use the higher *depth* block to enhance the lower *RGB* block, and use the lower *depth* block to enhance the higher *RGB* block, namely *cross-scale Depth-to-RGB weighting*. It is an important component of CMW-L and CMW-M. Concretely, the higher depth response maps are in charge of modulating the lower RGB features, and the lower depth response maps are responsible to modulate the higher RGB features. Such a cross-scale way captures cross-scale complementarity of cross-modal features. Besides, the Depth-to-RGB weighting is executed between two adjacent blocks, which can capture the continuity of features.

On the other hand, RGB features have the ability to modulate themselves. For this purpose, we introduce the *RGB-to-RGB weighting* to CMW-L and CMW-M. RGB features are enhanced by RGB response maps, which are



**Fig. 2.** Details of the three proposed RGB-depth interaction modules: CMW-L, CMW-M and CMW-H. All modules consist of Depth-to-RGB Weighting (DW) and RGB-to-RGB Weighting (RW) as key operations. Notably, the DW in CMW-L and CMW-M is performed in the cross-scale manner between two adjacent blocks, which effectively captures the feature continuity and activates cross-modal cross-scale interactions.

generated from RGB features. This allows our weighting modules to learn and adjust salient parts in an adaptive manner. Depth-to-RGB weighting and RGB-to-RGB weighting complement each other, improving the stability and robustness of our inference. Thus, for example, in CMW-M, R-E<sup>(3)</sup> is enhanced by D-E<sup>(4)</sup> and R-E<sup>(3)</sup>, and R-E<sup>(4)</sup> is enhanced by D-E<sup>(3)</sup> and R-E<sup>(4)</sup>. Notably, CMW-L and CMW-M perform the same cross-scale scheme, but with different resolutions. The multi-resolution enhanced object details of features benefit SOD.

**High-Level Cross-Modal Weighting Module.** The high-level part is distinct from the other two parts, and it contains rich global information. Thus, in the high-level part, we adopt CNN features of the highest blocks (*i.e.*, D-E<sup>(5)</sup> and R-E<sup>(5)</sup>) to accurately locate salient objects. We propose the CMW-H module, which is the modified variant of CMW-L and CMW-M, to enhance the macroscopic localization of salient objects. The RGB features of R-E<sup>(5)</sup> are enhanced by the depth response maps generated from D-E<sup>(5)</sup> and the RGB response maps generated from R-E<sup>(5)</sup>.

**Three-Level Decoder.** The decoder can make good use of features from the encoder with skip-connections. To fuse all the enhanced features for effective inference, the decoder part also consists of three levels, *i.e.*, D<sup>(5)</sup>, D<sup>(3&4)</sup> and D<sup>(1&2)</sup> as shown in Fig. 1, corresponding to high-level, middle-level and low-level encoder parts. Between the two adjacent levels, we employ the deconvolutional layer for  $4\times$  upsampling. Specifically, to effectively train the proposed CMWNet, we adopt the deep scale supervision [40] behind each level to force features of the decoder network to focus on salient objects.

### 3.2 Low- and Middle-Level Cross-Modal Weighting

We perform the enhancement on RGB features at low-level and middle-level parts with the CMW-L and CMW-M, respectively. The details of **CMW-L** and **CMW-M** are shown in Fig. 2. There are two types of weighting in each module, *i.e.*, cross-scale Depth-to-RGB weighting (DW) and RGB-to-RGB weighting (RW). For each encoder block, those two types of weighting only apply to the last feature layer. So we denote the last layer of features in D-E<sup>(l)</sup> and R-E<sup>(l)</sup> as  $\mathbf{f}_d^{(l)}$  and  $\mathbf{f}_r^{(l)}$ , respectively.

We provide a simplified version to explain the principle of DW and RW. For a feature map  $\mathbf{F} \in \mathbb{R}^{C_0 \times W_0 \times H_0}$ , there are two groups of weighting response maps  $\mathbf{r}_1 \in [0, 1]^{C_0 \times W_0 \times H_0}$  and  $\mathbf{r}_2 \in [0, 1]^{C_0 \times W_0 \times H_0}$  to modulate it at pixel level, and then the two types of weighting can be formulated as:

$$\mathbf{EF} = \mathbf{F} + \mathbf{r}_1 \otimes \mathbf{F} + \mathbf{r}_2 \otimes \mathbf{F}, \quad (1)$$

where  $\mathbf{EF} \in \mathbb{R}^{C_0 \times W_0 \times H_0}$  is the enhanced feature,  $\otimes$  is the element-wise multiplication, and  $+$  is the element-wise summation.  $\mathbf{r}_1 \otimes \mathbf{F}$  and  $\mathbf{r}_2 \otimes \mathbf{F}$  can be regarded as the DW operation and RW operation, respectively. If  $\mathbf{r}_1$  and  $\mathbf{r}_2$  have good responses to salient objects (*i.e.*, the pixel value is close to 1 on salient objects and close to 0 on background),  $\mathbf{F}$  will be accurately modulated to focus on the desired salient parts and  $\mathbf{EF}$  will have a stronger representation for salient objects. Thus, we apply the two types of weighting to RGB features and depth features for enhancement of salient object details in the CMW-L and CMW-M modules.

**Depth-to-RGB Weighting.** In our network, the Depth-to-RGB weighting is the most important operation to mine the complementarity of depth maps. It works in a cross-modal and cross-scale manner. To expand the receptive field and increase the feature diversity, we design a comprehensive structure of filters to generate the local and global features  $\mathbf{f}_{lg}^{(l)}$ .

Concretely, we adopt two convolutional layers with  $3 \times 3$  kernel as local filters. We also adopt a dilated convolution [41] with  $3 \times 3$  kernel and *rate* = 5 and a convolutional layers with  $7 \times 7$  kernel as global filters, as shown in DW<sup>(l)</sup> of Fig. 2. The global filters in the comprehensive structure expand the receptive field of convolution operations. The obtained global features can capture macro-level information of depth features, which are complementary to the local features. For each  $\mathbf{f}_d^{(l)}$ ,  $\mathbf{f}_{lg}^{(l)}$  can be computed as:

$$\mathbf{f}_{lg}^{(l)} = \text{concat}(C(\mathbf{f}_d^{(l)}; \mathbf{W}_{loc}^{(l_1)}), C(\mathbf{f}_d^{(l)}; \mathbf{W}_{loc}^{(l_2)}), C(\mathbf{f}_d^{(l)}; \mathbf{W}_{glo}^{(l_1)}), C(\mathbf{f}_d^{(l)}; \mathbf{W}_{glo}^{(l_2)})), \quad (2)$$

where  $\text{concat}(\cdot)$  denotes the cross-channel concatenation,  $C(*; \mathbf{W}_{loc}^{(l_i)})$  is a convolutional layer with parameters  $\mathbf{W}_{loc}^{(l_i)}$  (*i.e.*,  $\mathbf{W}_{loc}^{(l_1)}$  and  $\mathbf{W}_{loc}^{(l_2)}$  are  $3 \times 3$  kernel) for producing local features,  $C(*; \mathbf{W}_{glo}^{(l_1)})$  is a convolutional layer with parameters  $\mathbf{W}_{glo}^{(l_1)}$  (*i.e.*,  $\mathbf{W}_{glo}^{(l_1)}$  is  $7 \times 7$  kernel), and  $C(*; \mathbf{W}_{glo}^{(l_2)})$  is the dilated convolution with parameters  $\mathbf{W}_{glo}^{(l_2)}$  (*i.e.*,  $\mathbf{W}_{glo}^{(l_2)}$  is  $3 \times 3$  kernel with *rate* = 5).

Then, the multi-scale features in  $\mathbf{f}_{lg}^{(l)}$  are fused to generate the depth response maps  $\mathbf{r}_{dw}^{(l)}$ . Specifically, to make  $\mathbf{r}_{dw}^{(l)}$  have the same resolution as the corresponding cross-scale RGB features, the fusion operation is a convolutional layer with stride 2 for  $2\times$  downsampling for  $DW^{(1)}$  and  $DW^{(3)}$ , while a deconvolutional layer is used for  $DW^{(2)}$  and  $DW^{(4)}$  for  $2\times$  upsampling. Thus,  $\mathbf{r}_{dw}^{(l)}$  can be computed as:

$$\mathbf{r}_{dw}^{(l)} = \begin{cases} \sigma(C(\mathbf{f}_{lg}^{(l)}; \mathbf{W}_{dw}^{(l)})), l = 1, 3 \\ \sigma(De(\mathbf{f}_{lg}^{(l)}; \mathbf{W}_{dw}^{(l)})), l = 2, 4 \end{cases}, \quad (3)$$

where  $\sigma(\cdot)$  is the sigmoid function, and  $De(*; \mathbf{W}_{dw}^{(l)})$  is the deconvolutional layer with parameters  $\mathbf{W}_{dw}^{(l)}$ , which are  $2 \times 2$  kernel with stride 2.

Finally,  $\mathbf{r}_{dw}^{(l)}$  is used to enhance the cross-scale RGB features as follows:

$$\mathbf{f}_{dw}^{(l)} = \begin{cases} \mathbf{r}_{dw}^{(l+1)} \otimes \mathbf{f}_r^{(l)}, l = 1, 3 \\ \mathbf{r}_{dw}^{(l-1)} \otimes \mathbf{f}_r^{(l)}, l = 2, 4 \end{cases}. \quad (4)$$

**RGB-to-RGB Weighting.** Considering that the RGB features of low- and middle-level parts also contain rich information about details of salient objects, we propose the RGB-to-RGB weighting to adaptively enhance RGB features with the RGB response maps  $\mathbf{r}_{rw}^{(l)}$ , which are generated from  $\mathbf{f}_r^{(l)}$  as follows:

$$\mathbf{r}_{rw}^{(l)} = \sigma(C(\mathbf{f}_r^{(l)}; \mathbf{W}_{rw}^{(l)})). \quad (5)$$

The details of filters in  $RW^{(l)}$  are also presented in Fig. 2. Then, similar as depth response maps,  $\mathbf{r}_{rw}^{(l)}$  can enhance  $\mathbf{f}_r^{(l)}$  as follows:

$$\mathbf{f}_{rw}^{(l)} = \mathbf{r}_{rw}^{(l)} \otimes \mathbf{f}_r^{(l)}. \quad (6)$$

**Aggregation of Double Weighting Features.** After the DW and RW operations, the RGB features are enhanced twice and can capture the details of salient objects. To preserve the original color information, we add RGB features to these two groups of enhanced features to produce the double enhanced features  $\mathbf{f}_{de}^{(l)}$ . The aggregation of double weighting features is defined as:

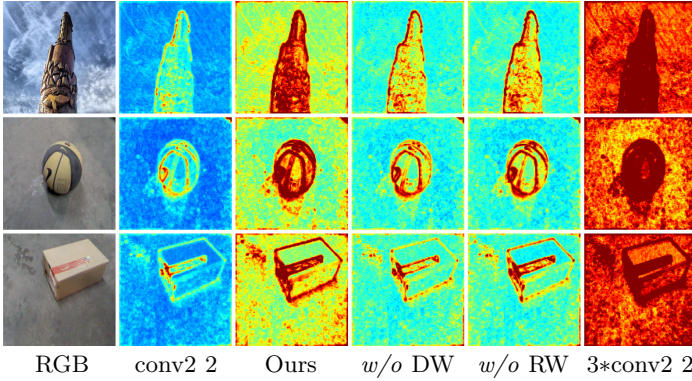
$$\mathbf{f}_{de}^{(l)} = \mathbf{f}_r^{(l)} + \mathbf{f}_{dw}^{(l)} + \mathbf{f}_{rw}^{(l)}. \quad (7)$$

Combining Eqs. 4, 6 and 7, we find that Eq. 7 is similar to Eq. 1. Thus, for the CMW-L and CMW-M, the output features  $\mathbf{f}_{cmw}^{(k)}$  can be computed as:

$$\mathbf{f}_{cmw}^{(k)} = Cat(\mathbf{f}_{de}^{(2k-1)}, De(\mathbf{f}_{de}^{(2k)}; \mathbf{W}_{cmw}^{(2k)})), \quad k = 1, 2. \quad (8)$$

Then,  $\mathbf{f}_{cmw}^{(2)}$  and  $\mathbf{f}_{cmw}^{(1)}$  boost the salient object inference in  $D^{(3\&4)}$  and  $D^{(1\&2)}$ , respectively, through skip-connections (*i.e.*, the dashed line in Fig. 1). More detailed parameters of CMW-L and CMW-M are shown in Fig. 2.





**Fig. 3.** Visualizing features of RGB conv2.2 in CMW-L. *w/o* DW: without adding DW features; *w/o* RW: without adding RW features; 3\*conv2.2: features with triple linear enhancement.

In Fig. 3, we visualize features of RGB conv2.2 in **CMW-L** to verify the effectiveness of the double weighting enhancement. By comparing “conv2.2” and “Ours”, salient objects are highlighted more clearly in “Ours”. If we delete  $\mathbf{f}_{dw}^{(l)}$  (*w/o* DW) or  $\mathbf{f}_{rw}^{(l)}$  (*w/o* RW), salient objects are more indistinct than “Ours”. We also show the features with triple linear enhancement (“3\*conv2.2”) to demonstrate that the double weighting enhancement is more effective than the conventional linear enhancement.

### 3.3 High-Level Cross-Modal Weighting

As for the high-level part, we modify the cross-scale CMW-L to the same-scale manner to effectively utilize the macroscopic semantic information. We propose the CMW-H for object localization enhancement. For DW operation in CMW-H, the RGB features are enhanced by depth response maps of the same layer. So, the DW operation in Eq. 4 is modified as follows:

$$\mathbf{f}_{dw}^{(l)} = \sigma(C(\mathbf{f}_{lg}^{(l)}; \mathbf{W}_{dw}^{(l)})) \otimes \mathbf{f}_r^{(l)}, \quad l = 5. \quad (9)$$

Other operations, such as the RW and features aggregation, are the same as those in CMW-L. Notably, the output of CMW-H is  $\mathbf{f}_{de}^{(5)}$ , which is directly fed to the decoder part. It leads the inference process of SOD, as shown in Fig. 1. The detailed structure of CMW-H is present in Fig. 2.

### 3.4 Implementation Details

**Loss Function.** As shown in Fig. 1, we add a convolutional layer after R-E<sup>(5)</sup>, D<sup>(5)</sup> and D<sup>(3&4)</sup> to generate intermediate predictions  $S^{(4)}$ ,  $S^{(3)}$  and  $S^{(2)}$  at the

training phase. Then, we utilize different scales of ground truth (GT) to supervise them and the final prediction  $S^{(1)}$  with the softmax loss. The total loss  $\mathbb{L}$  can be defined as:

$$\mathbb{L} = \sum_{t=1}^4 \alpha_t \cdot \mathcal{L}_{sm}^{(t)}(S^{(t)}, G^{(t)}), \quad (10)$$

where  $\mathcal{L}_{sm}^{(t)}(\cdot, \cdot)$  is the softmax loss,  $\alpha_t$  is the loss weight and set to 1, and  $G^{(t)}$  is a GT of the same resolution as  $S^{(t)}$ .

**Network Training Protocol.** Our CMWNet is implemented in Caffe [22] with an NVIDIA Titan X GPU. The parameters of the Siamese encoder part are initialized by the VGG16 model [34], except that the conv1\_1 of the depth stream is initialized by the Gaussian distribution with a standard deviation of 0.01. Other newly added layers are initialized using the Xavier initialization [17]. Following [13, 20], the training set consists of 1,400 triplets from NJU2K [23] and 650 triplets from NLPR [29]. We resize all training triplets to  $288 \times 288$ , and then we adopt the rotation ( $90^\circ$ ,  $180^\circ$  and  $270^\circ$ ) and mirror reflection for augmentation, resulting 10.25K training triplets. We employ the SGD [3] to train the network 22.5K iterations. The learning rate, batch size, iteration size, momentum and weight decay are set to  $10^{-7}$ , 1, 8, 0.9 and 0.0001, respectively. The learning rate will be divided by 10 12.5K iterations.

## 4 Experiments

### 4.1 Datasets and Evaluation Metrics

**Datasets.** We evaluate the proposed method and all compared methods on seven public benchmark datasets, including STEREO [28], NJU2K [23], LFSD [25], DES [7], NLPR [29], SSD [24] and SIP [13].

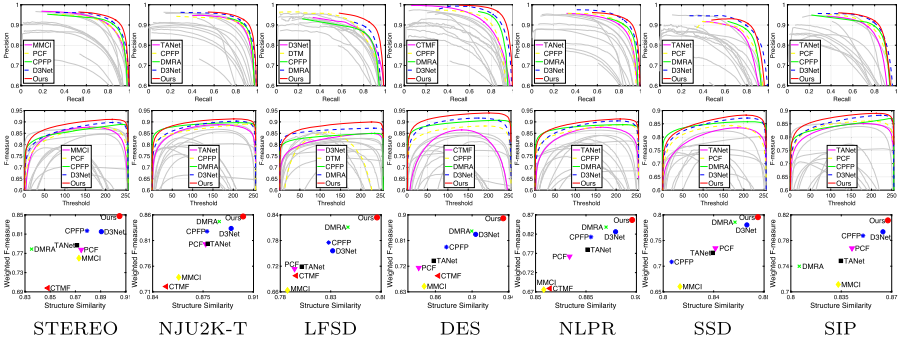
**Evaluation Metrics.** We evaluate the performance of our method and other methods using six widely used evaluation metrics including maximum F-measure ( $\mathcal{F}_\beta$ ,  $\beta^2 = 0.3$ ), weighted F-measure ( $\mathcal{F}_\beta^w$ ,  $\beta^2 = 1$ ) [27], mean absolute error (MAE,  $\mathcal{M}$ ), precision-recall (PR) curve, S-measure ( $\mathcal{S}_\lambda$ ,  $\lambda = 0.5$ ) [11], and maximum E-measure ( $\mathcal{E}_\xi$ ) [12].

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

**Comparison Methods.** We compare the proposed CMWNet with 6 state-of-the-art traditional methods, which are LBE [16], DCMC [9], SE [18], CDCP [45], MDSF [36] and DTM [8], and 9 state-of-the-art CNN-based methods, which are DF [31], CTMF [20], PCF [4], AFNet [38], MMCI [6], TANet [5], CFPF [42], DMRA [30] and D3Net [13]. The saliency maps of all compared methods are provided by authors or obtained by running their released codes. Notably, we retest DMRA [30] on STEREO dataset with 1,000 images, which results in different performances from the original DMRA paper.

**Table 1.** Quantitative results of 15 state-of-the-art methods on 7 datasets.  $\uparrow$  and  $\downarrow$  stand for larger and smaller is better, respectively. The best two results are in **bold** and *italics*. The *corner note* of each method is the publication year.

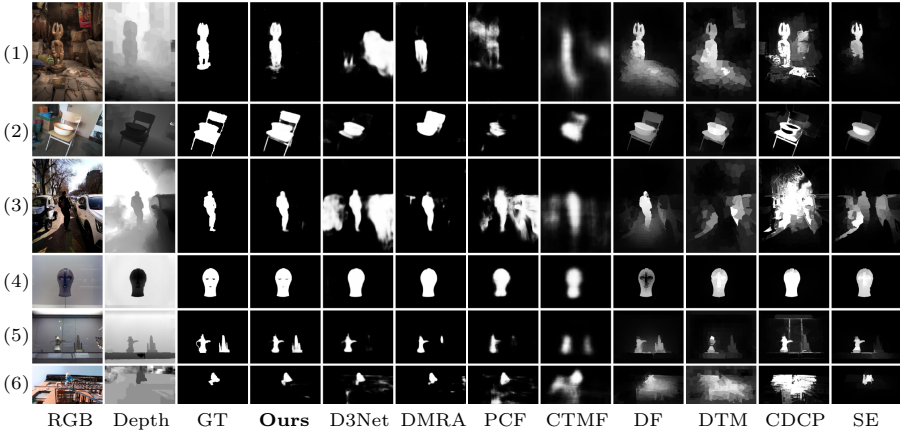
Models	STEREO [28]				NJU2K-T [23]				LFSD [25]				DES [7]				NLPR-T [29]				SSD [24]				SIP [13]			
	$S_\lambda$	$\mathcal{F}_\beta$	$\mathcal{E}_\xi$	$\mathcal{M}$	$S_\lambda$	$\mathcal{F}_\beta$	$\mathcal{E}_\xi$	$\mathcal{M}$	$S_\lambda$	$\mathcal{F}_\beta$	$\mathcal{E}_\xi$	$\mathcal{M}$	$S_\lambda$	$\mathcal{F}_\beta$	$\mathcal{E}_\xi$	$\mathcal{M}$	$S_\lambda$	$\mathcal{F}_\beta$	$\mathcal{E}_\xi$	$\mathcal{M}$	$S_\lambda$	$\mathcal{F}_\beta$	$\mathcal{E}_\xi$	$\mathcal{M}$	$S_\lambda$	$\mathcal{F}_\beta$	$\mathcal{E}_\xi$	$\mathcal{M}$
<b>LB</b> E <sub>16</sub> [16]	0.60	0.33	0.787	0.250	0.695	0.748	0.803	0.153	0.736	0.796	0.804	0.208	0.703	0.788	0.890	0.208	0.702	0.745	0.855	0.081	0.621	0.619	0.736	0.278	0.727	0.751	0.853	0.200
<b>DC</b> MC <sub>16</sub> [9]	0.731	0.740	0.819	0.148	0.686	0.715	0.799	0.172	0.753	0.817	0.856	0.155	0.707	0.666	0.773	0.111	0.724	0.648	0.793	0.117	0.704	0.711	0.786	0.169	0.683	0.618	0.743	0.186
<b>SE</b> <sub>16</sub> [18]	0.708	0.755	0.846	0.143	0.664	0.748	0.813	0.169	0.698	0.791	0.840	0.167	0.741	0.741	0.856	0.090	0.756	0.713	0.847	0.091	0.675	0.710	0.800	0.165	0.628	0.661	0.771	0.164
<b>CDCP</b> <sub>17</sub> [45]	0.713	0.664	0.786	0.149	0.669	0.621	0.741	0.180	0.717	0.703	0.786	0.167	0.709	0.631	0.811	0.115	0.727	0.645	0.820	0.112	0.603	0.535	0.700	0.214	0.595	0.505	0.721	0.224
<b>MDSF</b> <sub>17</sub> [36]	0.728	0.719	0.809	0.176	0.748	0.775	0.838	0.157	0.700	0.783	0.826	0.190	0.741	0.746	0.851	0.122	0.805	0.793	0.885	0.095	0.673	0.703	0.779	0.192	0.717	0.698	0.798	0.167
<b>DTM</b> <sub>18</sub> [8]	0.747	0.743	0.837	0.168	0.706	0.716	0.799	0.190	0.783	0.825	0.853	0.160	0.752	0.697	0.858	0.123	0.733	0.677	0.833	0.145	0.677	0.651	0.773	0.199	0.690	0.659	0.778	0.203
<b>DF</b> <sub>17</sub> [31]	0.757	0.757	0.847	0.141	0.763	0.804	0.864	0.141	0.791	0.816	0.865	0.138	0.752	0.766	0.870	0.093	0.802	0.778	0.880	0.085	0.747	0.735	0.828	0.142	0.653	0.657	0.759	0.185
<b>CTMF</b> <sub>18</sub> [20]	0.848	0.831	0.912	0.086	0.849	0.845	0.913	0.085	0.796	0.791	0.865	0.119	0.863	0.844	0.932	0.055	0.860	0.825	0.929	0.056	0.776	0.729	0.805	0.099	0.716	0.694	0.829	0.139
<b>PCF</b> <sub>18</sub> [4]	0.875	0.860	0.925	0.064	0.877	0.872	0.924	0.059	0.794	0.779	0.835	0.112	0.842	0.804	0.893	0.049	0.874	0.841	0.925	0.044	0.841	0.807	0.894	0.062	0.842	0.838	0.901	0.071
<b>AFNet</b> <sub>19</sub> [38]	0.825	0.823	0.887	0.075	0.772	0.775	0.853	0.100	0.738	0.744	0.815	0.133	0.770	0.728	0.881	0.068	0.799	0.771	0.879	0.058	0.714	0.687	0.807	0.118	0.720	0.712	0.819	0.118
<b>MMCI</b> <sub>19</sub> [6]	0.873	0.863	0.927	0.068	0.858	0.852	0.915	0.079	0.787	0.771	0.839	0.132	0.848	0.822	0.928	0.065	0.856	0.815	0.913	0.059	0.813	0.781	0.882	0.082	0.833	0.818	0.897	0.086
<b>TANet</b> <sub>19</sub> [5]	0.871	0.861	0.923	0.060	0.878	0.874	0.925	0.060	0.801	0.796	0.847	0.111	0.858	0.827	0.910	0.046	0.886	0.863	0.941	0.041	0.839	0.810	0.897	0.063	0.835	0.830	0.895	0.075
<b>CPPP</b> <sub>19</sub> [42]	0.879	0.874	0.925	0.051	0.878	0.877	0.923	0.053	0.828	0.826	0.872	0.088	0.872	0.846	0.923	0.038	0.888	0.867	0.932	0.036	0.807	0.766	0.852	0.082	0.850	0.851	0.903	0.064
<b>DMRA</b> <sub>19</sub> [30]	0.835	0.847	0.911	0.066	0.886	0.886	0.927	0.051	0.847	0.856	0.900	0.075	0.900	0.888	0.943	0.030	0.899	0.879	0.947	0.031	0.857	0.844	0.906	0.058	0.806	0.821	0.875	0.085
<b>D3Net</b> <sub>19</sub> [13]	0.891	0.881	0.930	0.054	0.895	0.889	0.932	0.051	0.832	0.819	0.864	0.099	0.904	0.885	0.946	0.030	0.906	0.885	0.946	0.034	0.866	0.847	0.910	0.058	0.864	0.862	0.903	0.062
<b>Ours</b>	0.905	0.901	0.944	0.043	0.903	0.902	0.936	0.046	0.876	0.883	0.912	0.066	0.934	0.930	0.969	0.022	0.917	0.903	0.951	0.029	0.875	0.871	0.930	0.051	0.867	0.874	0.913	0.062



**Fig. 4.** Quantitative comparisons on PR curve, F-measure curve and  $S_\lambda - \mathcal{F}_\beta^w$  coordinates. The top 5 methods on PR and F-measure curves are shown in color. For  $S_\lambda - \mathcal{F}_\beta^w$  coordinates, we only compare with several representative methods. (Color figure online)

**Quantitative Comparison.** We evaluate our method and the other 15 state-of-the-art methods under four quantitative metrics, including S-measure  $S_\lambda$ , max F-measure  $\mathcal{F}_\beta$ , max E-measure  $\mathcal{E}_\xi$  and MAE  $\mathcal{M}$ . As shown in Table 1, our method favorably outperforms all compared methods under these four metrics, and the recently proposed CNN-based methods [5, 13, 30, 42] and our method are superior to traditional methods by a large margin. Comparing to the second best results in Table 1, the performance of our method on the largest and challenging dataset STEREO [28] is improved by 1.6% and 2.0% in  $S_\lambda$  and  $\mathcal{F}_\beta$ , respectively. The improvement on the relatively small dataset DES [7] is remarkable, with an increase of 3.0% and 4.2% in  $S_\lambda$  and  $\mathcal{F}_\beta$ , respectively. For the salient person detection, our method improves the performance by 1.2% in  $\mathcal{F}_\beta$  on SIP [13].

In addition, we also present the PR curves, F-measure curves and  $S_\lambda$  (X-axis) -  $\mathcal{F}_\beta^w$  (Y-axis) coordinates in Fig. 4. The performance under these metrics is consistent with that in Table 1. The superiority of our method is more visible



**Fig. 5.** Visual comparisons with eight representative methods, including five CNN-based methods and three traditional methods.

on STEREO [28], LFS [25] and DES [7]. Both Table 1 and Fig. 4 demonstrate that our method is consistently better than all compared methods in terms of different evaluation metrics.

**Visual Comparison.** We show visual comparisons with 8 representative methods in Fig. 5. Each row in Fig. 5 represents a challenging scenario for SOD, including low contrast (1<sup>st</sup> row), disturbing background (2<sup>nd</sup> row), salient person detection (3<sup>rd</sup> row), object with fine structures (4<sup>th</sup> row), multiple objects (5<sup>th</sup> row) and small object (6<sup>th</sup> row). Regardless of different scene, our method can accurately highlight salient objects with fine details. Notably, in the 4<sup>th</sup> row, the mask has a fine structure with three holes on it. Thanks to the DW and RW operations in our method, our method successfully highlights the mask with three holes, while other methods fail.

### 4.3 Ablation Studies

We conduct detailed ablation studies of our CMWNet on a big dataset, NJU2K [23], and a small but challenging dataset, SSD [24]. Specifically, we assess (1) the rationality of enhancing RGB features with depth features; (2) the individual contributions of CMW-L&M and CMW-H; (3) the importance of weighting; (4) the rationality of cross-scale weighting of CMW-L&M; and (5) the necessity of deep scale supervision of decoder part. We change one component at a time to evaluate individual contributions.

**Rationality of Enhancing RGB Features with Depth Features.** In our method, we use depth features to enhance RGB features (“DeR”). To study the rationality of this enhancement manner, we explore another baseline variant: adopting RGB features to enhance depth features (“ReD”). From Table 2, we

**Table 2.** Ablation studies on *NJU2K* [23] and *SSD* [24]. The best result of each column is **bold**. Details are introduced in Sect. 4.3.

Models	NJU2K-T [23]				SSD [24]			
	$\mathcal{S}_\lambda \uparrow$	$\mathcal{F}_\beta \uparrow$	$\mathcal{E}_\xi \uparrow$	$\mathcal{M} \downarrow$	$\mathcal{S}_\lambda \uparrow$	$\mathcal{F}_\beta \uparrow$	$\mathcal{E}_\xi \uparrow$	$\mathcal{M} \downarrow$
<b>Ours (DeR)</b>	<b>.903</b>	<b>.902</b>	<b>.936</b>	<b>.046</b>	<b>.875</b>	<b>.871</b>	<b>.930</b>	<b>.051</b>
ReD	.889	.887	.927	.056	.864	.850	.909	.063
<i>w/o</i> depth ( <i>w/o</i> DW)	.886	.886	.924	.056	.855	.842	.915	.064
<i>w/o</i> CMW-L&M	.891	.886	.932	.053	.849	.839	.909	.066
<i>w/o</i> CMW-H	.896	.894	.929	.051	.853	.845	.908	.063
<i>w/o</i> RW	.901	.899	.933	.046	.868	.861	.919	.056
<i>w/o</i> Wei	.900	.898	.933	.048	.858	.839	.899	.061
DW <i>w/o</i> GF	.900	.898	.933	.048	.868	.858	.923	.054
RW <i>w/</i> GF	.901	.900	.934	.046	.870	.867	.924	.052
<i>w/o</i> CS	.901	.898	.932	.047	.864	.861	.922	.060
C2S	.900	.899	.933	.049	.864	.847	.906	.060
<i>w/o</i> DS	.898	.898	.933	.049	.866	.862	.923	.055

observe that the performance on both datasets has dropped (*e.g.*  $\mathcal{M}$ : 0.046  $\rightarrow$  0.056 on NJU2K and 0.051  $\rightarrow$  0.063 on SSD). This confirms that using depth features to enhance RGB features is more reasonable than the other direction for extracting the cross-modal complementarity.

In addition, we remove the depth map input in our network to evaluate the power of depth map (*w/o* depth). This variant is for RGB SOD. The performance of *w/o* depth drops sharply (*e.g.*  $\mathcal{M}$ : 0.046  $\rightarrow$  0.056 on NJU2K). This confirms that the way of exploring complementary distance information of depth map in our network is effective.

**Individual Contributions of CMW-L&M and CMW-H.** The proposed three RGB-depth interaction modules can be divided into two types. CMW-L and CMW-M (CMW-L&M) are responsible for object details enhancement, and CMW-H for object localization enhancement. Thus, we provide two variants of our network: removing the CMW-L&M (*w/o* CMW-L&M) and removing the CMW-H (*w/o* CMW-H). From Table 2, we observe a significant performance degradation (*e.g.*  $\mathcal{S}_\lambda$ : 0.903  $\rightarrow$  0.891 on NJU2K and 0.875  $\rightarrow$  0.849 on SSD) of *w/o* CMW-L&M. This confirms that the proposed CMW-L&M are momentous to our network, and they enhance the details of salient object in low- and middle-level features clearly. Some enhanced features in **CMW-L** are shown in the third column of Fig. 3. The performance drop (*e.g.*  $\mathcal{F}_\beta$ : 0.902  $\rightarrow$  0.894 on NJU2K and 0.871  $\rightarrow$  0.845 on SSD) of *w/o* CMW-H means that the proposed CMW-H is also important to our network and it enhances the salient object localization in high-level features accurately.

**Importance of Weighting.** To study the importance of two types of weighting, we derive three variants: removing the Depth-to-RGB weighting (*w/o* DW), removing the RGB-to-RGB weighting (*w/o* RW), and using concatenation of depth features and RGB features instead of two types of weighting (*w/o* Wei). Specifically, *w/o* DW is the same as *w/o* depth, *i.e.*, the depth map does not participate in SOD. The depth map is still utilized in *w/o* Wei, which is not equal to *w/o* (DW + RW). It focuses on evaluating the impact of weighting mechanism. According to the statistics in Table 2, we observe the performance of these three variants is worse than our complete CMWNet. This demonstrates that the two types of weighting can help our CMWNet to better highlight salient objects with effective feature enhancement. We also provide two variants, *i.e.* DW *w/o* GF and RW *w/o* GF, to confirm the rationality of specific global filters (GF) in DW.

In addition, the visualization of features *w/o* DW and *w/o* RW is shown in Fig. 3, in which the boundaries of salient objects *w/o* RW are much clearer than *w/o* DW. This demonstrates that the depth map does assist the RGB-D SOD, and the enhancement effect of DW is more effective than RW.

**Rationality of Cross-Scale Weighting in CMW-L&M.** To study the rationality of cross-scale weighting in CMW-L&M, we modify the cross-scale DW to the same-scale manner (*w/o* CS), which is the same as CMW-H. The double enhanced features of each scale in CMW-L&M are concatenated for inference. By comparing *w/o* CS and Ours in Table 2, we find that the performance of *w/o* CS decreases on both NJU2K and SSD. This demonstrates that the DW of CMW-L&M in the cross-scale manner is rational, and this manner can enhance interactions between different scales to further boost performance.

Besides, to study the rationality of performing cross-scale weighting between adjacent CNN blocks, we provide a variant which performs cross-scale weighting between nonadjacent CNN blocks, *i.e.*, R-E<sup>(1)</sup> is enhanced by D-E<sup>(3)</sup>, R-E<sup>(2)</sup> is enhanced by D-E<sup>(4)</sup>, R-E<sup>(3)</sup> is enhanced by D-E<sup>(1)</sup> and R-E<sup>(4)</sup> is enhanced by D-E<sup>(2)</sup> (C2S). As the results presented in Table 2, we observe that the results of C2S are worse than Ours. The reason behind this is that the weighting performed across two scales (*i.e.* C2S) may lose the continuity of features, causing the depth response maps to fail to highlight the salient objects of RGB features. In contrast, the weighting performed between two adjacent CNN blocks (*i.e.* Ours) can capture the continuity of features and precisely increase cross-scale interactions.

**Necessity of Deep Scale Supervision in Decoder.** To study the necessity of deep scale supervision, we provide a baseline with only one supervision of the final prediction  $S^{(1)}$  (*w/o* DS). As shown in Table 2, we observe that network training with additional supervision is better than the single supervision. This verifies that the multiple scale supervision during network training can improve the testing performance. Besides, the intermediate predictions  $S^{(4)}$ ,  $S^{(3)}$  and  $S^{(2)}$  are also shown in Fig. 1. We can observe that the refinement process of prediction from coarse ( $S^{(4)}$ ) to fine ( $S^{(1)}$ ) benefits from the deep scale supervision in decoder part, which visually confirms the necessity of deep scale supervision.

## 5 Conclusion

In this paper, we propose a novel Cross-Modal Weighting Network (CMWNet) for RGB-D SOD. In particular, three novel cross-modal cross-scale weighting modules (CWM-L, CMW-M and CMW-H) are designed to encourage feature interactions for improving SOD performance. Based on these improvements, a three-level decoder progressively refines salient objects. Extensive experiments are conducted to validate our CMWNet, which achieves the best performance on seven public RGB-D SOD benchmarks in comparison with 15 state-of-the-arts.

**Acknowledgments.** This work was supported by the National Natural Science Foundation of China under Grant 61771301. Linwei Ye and Yang Wang were supported by NSERC.

## References

1. Borji, A., Cheng, M.-M., Hou, Q., Jiang, H., Li, J.: Salient object detection: a survey. *Comput. Vis. Media* **5**(2), 117–150 (2019)
2. Borji, A., Cheng, M.M., Jiang, H., Li, J.: Salient object detection: a benchmark. *IEEE TIP* **24**(12), 5706–5722 (2015)
3. Bottou, L.: Large-scale machine learning with stochastic gradient descent. In: *COMPSTAT* (2010)
4. Chen, H., Li, Y.: Progressively complementarity-aware fusion network for RGB-D salient object detection. In: *IEEE CVPR* (2018)
5. Chen, H., Li, Y.: Three-stream attention-aware network for RGB-D salient object detection. *IEEE TIP* **28**(6), 2825–2835 (2019)
6. Chen, H., Li, Y., Su, D.: Multi-modal fusion network with multi-scale multi-path and cross-modal interactions for RGB-D salient object detection. *Pattern Recogn.* **86**, 376–385 (2019)
7. Cheng, Y., Fu, H., Wei, X., Xiao, J., Cao, X.: Depth enhanced saliency detection method. In: *ACM ICIMCS* (2014)
8. Cong, R., Lei, J., Fu, H., Hou, J., Huang, Q., Kwong, S.: Going from RGB to RGBD saliency: a depth-guided transformation model. *IEEE TCYB* **50**, 3627–3639 (2019). <https://doi.org/10.1109/TCYB.2019.2932005>
9. Cong, R., Lei, J., Zhang, C., Huang, Q., Cao, X., Hou, C.: Saliency detection for stereoscopic images based on depth confidence analysis and multiple cues fusion. *IEEE SPL* **23**(6), 819–823 (2016)
10. Ding, Y., Liu, Z., Huang, M., Shi, R., Wang, X.: Depth-aware saliency detection using convolutional neural networks. *J. Vis. Commun. Image Represent.* **61**, 1–9 (2019)
11. Fan, D.P., Cheng, M.M., Liu, Y., Li, T., Borji, A.: Structure-measure: a new way to evaluate foreground maps. In: *IEEE ICCV* (2017)
12. Fan, D.P., Gong, C., Cao, Y., Ren, B., Cheng, M.M., Borji, A.: Enhanced-alignment measure for binary foreground map evaluation. In: *IJCAI* (2018)
13. Fan, D.P., et al.: Rethinking RGB-D salient object detection: models, datasets, and large-scale benchmarks. *arXiv preprint arXiv:1907.06781* (2019)
14. Fan, X., Liu, Z., Sun, G.: Salient region detection for stereoscopic images. In: *IEEE DSP* (2014)

15. Fang, Y., Wang, J., Narwaria, M., Callet, P.L., Lin, W.: Saliency detection for stereoscopic images. *IEEE TIP* **23**(6), 2625–2636 (2014)
16. Feng, D., Barnes, N., You, S., McCarthy, C.: Local background enclosure for RGB-D salient object detection. In: *IEEE CVPR* (2016)
17. Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: *AISTATS* (2010)
18. Guo, J., Ren, T., Bei, J.: Salient object detection for RGB-D image via saliency evolution. In: *IEEE ICME* (2016)
19. Guo, J., Ren, T., Jia, B., Zhu, Y.: Salient object detection in RGB-D image based on saliency fusion and propagation. In: *ACM ICIMCS* (2015)
20. Han, J., Chen, H., Liu, N., Yan, C., Li, X.: CNNs-based RGB-D saliency detection via cross-view transfer and multiview fusion. *IEEE TCYB* **48**(11), 3171–3183 (2018)
21. Itti, L., Koch, C., Niebur, E.: A model of saliency-based visual attention for rapid scene analysis. *IEEE TPAMI* **20**(11), 1254–1259 (1998)
22. Jia, Y., et al.: Caffe: convolutional architecture for fast feature embedding. In: *ACM MM* (2014)
23. Ju, R., Ge, L., Geng, W., Ren, T., Wu, G.: Depth saliency based on anisotropic center-surround difference. In: *IEEE ICIP* (2014)
24. Li, G., Zhu, C.: A three-pathway psychobiological framework of salient object detection using stereoscopic technology. In: *IEEE ICCVW* (2017)
25. Li, N., Ye, J., Ji, Y., Ling, H., Yu, J.: Saliency detection on light field. In: *IEEE CVPR* (2014)
26. Liu, Z., Shi, S., Duan, Q., Zhang, W., Zhao, P.: Salient object detection for RGB-D image by single stream recurrent convolution neural network. *Neurocomputing* **363**, 46–57 (2019)
27. Margolin, R., Zelnik-Manor, L., Tal, A.: How to evaluate foreground maps. In: *IEEE CVPR* (2014)
28. Niu, Y., Geng, Y., Li, X., Liu, F.: Leveraging stereopsis for saliency analysis. In: *IEEE CVPR* (2012)
29. Peng, H., Li, B., Xiong, W., Hu, W., Ji, R.: RGBD salient object detection: a benchmark and algorithms. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) *ECCV 2014*. LNCS, vol. 8691, pp. 92–109. Springer, Cham (2014). [https://doi.org/10.1007/978-3-319-10578-9\\_7](https://doi.org/10.1007/978-3-319-10578-9_7)
30. Piao, Y., Ji, W., Li, J., Zhang, M., Lu, H.: Depth-induced multi-scale recurrent attention network for saliency detection. In: *IEEE ICCV* (2019)
31. Qu, L., He, S., Zhang, J., Tian, J., Tang, Y., Yang, Q.: RGBD salient object detection via deep fusion. *IEEE TIP* **26**(5), 2274–2285 (2017)
32. Ren, J., Gong, X., Yu, L., Zhou, W., Yang, M.Y.: Exploiting global priors for RGB-D saliency detection. In: *IEEE CVPRW* (2015)
33. Shigematsu, R., Feng, D., You, S., Barnes, N.: Learning RGB-D salient object detection using background enclosure, depth contrast, and top-down features. In: *IEEE ICCVW* (2017)
34. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: *ICLR* (2015)
35. Song, H., Liu, Z., Du, H., Sun, G., Bai, C.: Saliency detection for RGBD images. In: *ACM ICIMCS* (2015)
36. Song, H., Liu, Z., Du, H., Sun, G., Olivier, L.M., Ren, T.: Depth-aware salient object detection and segmentation via multiscale discriminative saliency fusion and bootstrap learning. *IEEE TIP* **26**(9), 4204–4216 (2017)



37. Wang, A., Wang, M.: RGB-D salient object detection via minimum barrier distance transform and saliency fusion. *IEEE SPL* **24**(5), 663–667 (2017)
38. Wang, N., Gong, X.: Adaptive fusion for RGB-D salient object detection. *IEEE Access* **7**, 55277–55284 (2019)
39. Wang, W., Lai, Q., Fu, H., Shen, J., Ling, H.: Salient object detection in the deep learning era: an in-depth survey. *arXiv preprint [arXiv:1904.09146](https://arxiv.org/abs/1904.09146)* (2019)
40. Xie, S., Tu, Z.: Holistically-nested edge detection. In: *IEEE ICCV* (2015)
41. Yu, F., Koltun, V.: Multi-scale context aggregation by dilated convolutions. In: *ICLR* (2016)
42. Zhao, J.X., Cao, Y., Fan, D.P., Cheng, M.M., Li, X.Y., Zhang, L.: Contrast prior and fluid pyramid integration for RGBD salient object detection. In: *IEEE CVPR* (2019)
43. Zhou, Z., Wang, Z., Lu, H., Wang, S., Sun, M.: Global and local sensitivity guided key salient object re-augmentation for video saliency detection. *arXiv preprint [arXiv:1811.07480](https://arxiv.org/abs/1811.07480)* (2018)
44. Zhu, C., Cai, X., Huang, K., Li, T.H., Li, G.: PDNet: prior-model guided depth-enhanced network for salient object detection. In: *IEEE ICME* (2019)
45. Zhu, C., Li, G., Wang, W., Wang, R.: An innovative salient object detection using center-dark channel prior. In: *IEEE ICCVW* (2017)