

# DPANet: Depth Potentiality-Aware Gated Attention Network for RGB-D Salient Object Detection

Zuyao Chen, Runmin Cong, *Member, IEEE*, Qianqian Xu, *Senior Member, IEEE*,  
and Qingming Huang, *Fellow, IEEE*

**Abstract**—There are two main issues in RGB-D salient object detection: (1) how to effectively integrate the complementarity from the cross-modal RGB-D data; (2) how to prevent the contamination effect from the unreliable depth map. In fact, these two problems are linked and intertwined, but the previous methods tend to focus only on the first problem and ignore the consideration of depth map quality, which may yield the model fall into the sub-optimal state. In this paper, we address these two issues in a holistic model synergistically, and propose a novel network named *DPANet* to explicitly model the potentiality of the depth map and effectively integrate the cross-modal complementarity. By introducing the depth potentiality perception, the network can perceive the potentiality of depth information in a learning-based manner, and guide the fusion process of two modal data to prevent the contamination occurred. The gated multi-modality attention module in the fusion process exploits the attention mechanism with a gate controller to capture long-range dependencies from a cross-modal perspective. Experimental results compared with 16 state-of-the-art methods on 8 datasets demonstrate the validity of the proposed approach both quantitatively and qualitatively. <https://github.com/JosephChenHub/DPANet>

**Index Terms**—Salient object detection, RGB-D images, Depth potentiality perception, Gated multi-modality attention.

Manuscript received Mar. 2020. This work was supported in part by the National Key R&D Program of China under Grant 2018AAA0102003, in part by the Beijing Nova Program under Grant Z201100006820016, in part by the National Natural Science Foundation of China under Grant 62002014, Grant 61620106009, Grant U1636214, Grant 61931008, Grant 61836002, Grant 61672514, Grant 61976202, in part by Key Research Program of Frontier Sciences under Grant CAS: QYZDJ-SSW-SYS013, in part by the Strategic Priority Research Program of Chinese Academy of Sciences under Grant XDB28000000, in part by Beijing Natural Science Foundation under Grant 4182079, and in part by Youth Innovation Promotion Association CAS, in part by Hong Kong Scholars Program, in part by China Postdoctoral Science Foundation under Grant 2020T130050, Grant 2019M660438, and in part by the Fundamental Research Funds for the Central Universities under Grant 2019RC039. (Zuyao Chen and Runmin Cong contributed equally to this work.) (Corresponding authors: Qianqian Xu and Qingming Huang.)

Z. Chen is with the School of Computer Science and Technology, University of Chinese Academy of Sciences, Beijing 100190, China (e-mail: chenzuyao17@mails.ucas.ac.cn).

R. Cong is with the Institute of Information Science, Beijing Jiaotong University, Beijing 100044, China, also with the Beijing Key Laboratory of Advanced Information Science and Network Technology, Beijing 100044, China, and also with the Department of Computer Science, City University of Hong Kong, Hong Kong SAR, China (e-mail: rmcong@bjtu.edu.cn).

Q. Xu is with the Key Laboratory of Intelligent Information Processing, Institute of Computing Technology, Chinese Academy of Sciences, Beijing 100190, China (e-mail: xuqianqian@ict.ac.cn).

Q. Huang is with the School of Computer Science and Technology, University of Chinese Academy of Sciences, Beijing 101408, China, also with the Key Laboratory of Big Data Mining and Knowledge Management (BDKM), University of Chinese Academy of Sciences, Beijing 101408, China, also with the Key Laboratory of Intelligent Information Processing, Institute of Computing Technology, Chinese Academy of Sciences, Beijing 100190, China, and also with Peng Cheng Laboratory, Shenzhen 518055, China (e-mail: qmhuang@ucas.ac.cn).

## I. INTRODUCTION

**S**ALIENT object detection (SOD) aims to locate interesting regions that attract human attention most in an image [1], [2]. As a pre-processing technique, SOD benefits a variety of applications including object segmentation [3], person re-identification [4], image understanding [5], thumbnail creation [6], image quality assessment [7], and image enhancement [8]. In the past years, CNN-based methods have achieved promising performances in the SOD task owing to its powerful representation ability of CNN [9]. Most of them [10]–[23] focused on detecting the salient objects from the RGB image, while it is hard to achieve better performance in some challenging and complex scenarios using only one single modal data, such as similar appearance between the foreground and background (see the first row in Fig. 1), the cluttered background interferences (see the second row in Fig. 1). Recently, depth information has become increasingly popular thanks to the affordable and portable devices, *e.g.*, Microsoft Kinect and iPhone XR, which can provide many useful and complementary cues in addition to the color appearance information, such as shape structure and boundary information and has been successfully applied in many vision tasks [24]–[32]. Introducing depth information into SOD does address these challenging scenarios to some degree. However, as shown in Fig. 1, there exists a conflict that depth maps are sometimes inaccurate and would contaminate the results of SOD. Previous works generally integrate the RGB and depth information in an indiscriminate manner, which may induce negative results when encountering the inaccurate or blurred depth maps. Moreover, it is often insufficient to fuse and capture complementary information of different modal from the RGB image and depth map via simple strategies such as cascading, multiplication, which may degrade the saliency result. Hence, there are two main issues in RGB-D SOD to be addressed, *i.e.*, 1) how to prevent the contamination from unreliable depth information; 2) how to effectively integrate the multi-modal information from the RGB image and corresponding depth map.

As a remedy for the above-mentioned issues, we propose a Depth Potentiality-Aware Gated Attention Network (*DPANet*) that can simultaneously model the potentiality of the depth map, and optimize the fusion process of RGB and depth information in a gated attention mechanism. Instead of indiscriminately integrating multi-modal information from the RGB image and depth map, we focus on adaptively fusing two modal data by considering the depth potentiality perception

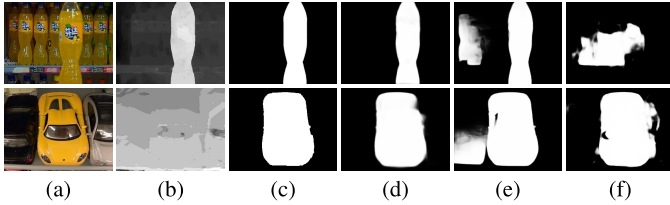


Fig. 1. Sample results of our method compared with others. RGB-D methods are marked in **boldface**. (a) RGB image; (b) Depth map; (c) Ground truth; (d) **Ours**; (e) BASNet [14]; (f) CPFPP [33].

in a learning-based manner. The depth potentiality perception works as a controller to guide the aggregation of cross-modal information and prevent the contamination from the unreliable depth map. For the depth potentiality perception, we take the saliency detection as the task orientation to measure the relationship between the binary depth map and the corresponding saliency mask. If the binary depth map obtained by a simple thresholding method (*e.g.*, Ostu [34]) is close to the ground truth, the reliability of the depth map is high, and therefore a higher depth confidence response should be assigned for this depth input. The depth confidence response will be adopted in the gated multi-modality fusion to better integrate the cross-modal information and prevent the contamination.

Considering the complementarity and inconsistency of RGB and depth information, we propose a gated multi-modality attention (GMA) module to capture long-range dependencies from a cross-modal perspective. Concatenating or summing the cross-modal features from RGB-D images not only has information redundancy, but also makes the truly effective complementary information submerged in a large number of data features. Therefore, the GMA module we designed utilizes the spatial attention mechanism to extract the most discriminative features, and can adaptively use the gate function to control the fusion rate of the cross-modal information and reduce the negative impact caused by the unreliable depth map. Moreover, we design a multi-level feature fusion mechanism to better integrate different levels of features including single-modality and multi-modality features. As is shown in Fig. 1, the proposed network can handle some challenging scenarios, such as the background disturbances (similar appearance), complex scenes, and unreliable depth maps.

In summary, our main contributions are listed as follows:

- For the first time, we address the unreliable depth map in the RGB-D SOD network in an end-to-end formulation, and propose the *DPANet* by incorporating the depth potentiality perception into the cross-modality integration pipeline.
- Without increasing the training label (*i.e.*, depth quality label), we model a task-orientated depth potentiality perception module that can adaptively perceive the potentiality of the input depth map, and further weaken the contamination from unreliable depth information.
- We propose a GMA module to effectively aggregate the cross-modal complementarity of the RGB and depth images, where the spatial attention mechanism aims at reducing the information redundancy, and the gate func-

tion controller focuses on regulating the fusion rate of the cross-modal information.

- Without any pre-processing (*e.g.*, HHA [35]) or post-processing (*e.g.*, CRF [36]) techniques, the proposed network outperforms 16 state-of-the-art methods on 8 RGB-D SOD datasets in quantitative and qualitative evaluations.

## II. RELATED WORK

In this section, we will review the salient object detection models for RGB and RGB-D images, especially the deep learning based methods, which have achieved impressive progress in recent years.

### A. RGB Salient Object Detection

The last decades have witnessed the prosperous development and improvement of salient object detection for RGB image. Specifically, from the bottom-up models [37]–[40] to top-down models [10]–[20], the performance is constantly being refreshed and surpassed. The bottom-up models used some priors (*e.g.*, background prior [38], compactness prior [39]) and properties (*e.g.*, sparse [37], low rank [40]) to define the saliency. Recently, deep learning based SOD methods have gradually become mainstream and achieved significant performance improvements. Hou *et al.* [10] introduced short connections into the skip-layer structures within the holistically-nested edge detector network architecture to achieve accurate saliency detection. The edge or boundary cue is embedded in the deep model to highlight the boundary region of salient object and improve the performance of SOD, such as EGNNet [15], BASNet [14], and AFNet [16]. Attention mechanism has been successfully applied in SOD to learn more discriminative feature representation, such as PiCANet [12], PFANet [17], and GCPANet [18]. In addition, some studies focus on real-time SOD or weakly-/un-supervised learning. Liu *et al.* [13] explored the role of pooling in neural networks for SOD, which can yield detail enriched saliency maps. Zhang *et al.* [19] used the noisy label to train the SOD model and achieved comparable performance with the state-of-the-art supervised deep models. Zeng *et al.* [20] proposed a unified framework to train SOD models with diverse weak supervision sources including category labels, captions, and unlabelled data.

### B. RGB-D Salient Object Detection

For the unsupervised SOD method for RGB-D images, some handcrafted features are designed, such as depth contrast, depth measure. Peng *et al.* [41] proposed to fuse the RGB and depth at first and feed it to a multi-stage saliency model. Song *et al.* [42] proposed a multi-scale discriminative saliency fusion framework. Feng *et al.* [43] proposed a Local Background Enclosure (LBE) to capture the spread of angular directions. Ju *et al.* [44] proposed a depth-induced method based on using anisotropic center-surround difference. Fan *et al.* [45] combined the region-level depth, color and spatial information to achieve saliency detection in stereoscopic images. Cheng *et*

*al.* [46] measured the salient value using color contrast, depth contrast, and spatial bias.

Lately, deep learning based approaches have gradually become a mainstream trend in RGB-D saliency detection. Qu *et al.* [47] proposed to fuse different low-level saliency cues into hierarchical features, including local contrast, global contrast, background prior and spatial prior. Zhu *et al.* [48] designed a master network to process RGB values, and a sub-network for depth cues and incorporate the depth-based features into the master network. Chen *et al.* [49] designed a progressive network attempting to integrate cross-modal complementarity. Zhao *et al.* [33] integrated the RGB features and enhanced depth cues using depth prior for SOD. Piao *et al.* [50] proposed a depth-induced multi-scale recurrent attention network. However, these efforts attempting to integrate the RGB and depth information indiscriminately ignore contaminations from inaccurate or blurred depth maps. Fan *et al.* [51] attempted to address this issue by designing a depth deparator unit to abandon the low-quality depth maps.

### III. METHODOLOGY

#### A. Overview of the Proposed Network

As shown in Fig. 2, the proposed network is a symmetrical two-stream encoder-decoder architecture. To be concise, we denote the output features of RGB branch in the encoder component as  $rb_i$  ( $i = 1, 2, 3, 4, 5$ ), and the features of depth branch in the encoder component as  $db_i$  ( $i = 1, 2, 3, 4, 5$ ). The feature  $rb_i$  ( $i = 2, 3, 4, 5$ ) and feature  $db_i$  ( $i = 2, 3, 4, 5$ ) are fed into a GMA module to obtain the corresponding enhanced feature  $rf_i$ ,  $df_i$ , respectively. In GAM module, the weight of the gate is learned by the network in a supervised way. Specifically, the top layers' feature  $rb_5$  and  $db_5$  are passed through a global average pooling (GAP) layer and two fully connected layers to learn the predicted score of depth potentiality via the regression loss with the help of the pseudo labels. Then, the decoder of two branches integrates multi-scale features progressively. Finally, we aggregate the two decoders' output and generate the saliency map by using the Multi-scale and Multi-modality Feature Fusion Modules. To facilitate the optimization, we add auxiliary loss branches at each sub-stage, *i.e.*,  $rd_i$  and  $dd_i$  ( $i = 5, 4, 3, 2$ ).

#### B. Depth Potentiality Perception

Most previous works [33], [48]–[50], [52] generally integrate the multi-modal features from RGB and corresponding depth information indiscriminately. However, as mentioned before, there exist some contaminations when depth maps are unreliable. To address this issue, Fan *et al.* [51] proposed a depth deparator unit to switch the RGB path and RGB-D path in a mechanical and unsupervised way. Different from the work [51], our proposed network can explicitly model the confidence response of the depth map and control the fusion process in a soft manner rather than directly discard the low-quality depth map.

Since we do not hold any labels for depth map quality assessment, we model the depth potentiality perception as a saliency-oriented prediction task, that is, we train a model

to automatically learn the relationship between the binary depth map and the corresponding saliency mask. The above modeling approach is based on the observation that if the binary depth map segmented by a threshold is close to the ground truth, the depth map is highly reliable, so a higher confidence response should be assigned to this depth input. Specifically, we first apply Otsu [34] to binarize the depth map  $I$  into a binary depth map  $\tilde{I}$ , which describes the potentiality of the depth map from the saliency perspective. Then, we design a measurement to quantitatively evaluate the degree of correlation between the binary depth map and the ground truth. IoU (intersection over union) is adopted to measure the accuracy between the binary map  $\tilde{I}$  and the ground truth  $G$ , which can be formulated as:

$$D_{\text{iou}} = \frac{|\tilde{I} \cap G|}{|\tilde{I} \cup G|}, \quad (1)$$

where  $|\cdot|$  denotes the area. However, in some cases, the coarse binary depth map will contain the background, which causes the  $D_{\text{iou}}$  tending to small even if the final saliency map is very closed to the ground truth. Hence, we define another metric to relax the strong constraint of IoU, which is defined as:

$$D_{\text{cov}} = \frac{|\tilde{I} \cap G|}{|G|}. \quad (2)$$

This metric  $D_{\text{cov}}$  reflects the ratio of intersection area to the ground truth, which indicates that the binary depth map is expected to cover the salient object more complete. Finally, inspired by the F-measure [53], we combine these two metrics to measure the potentiality of depth map for SOD task, *i.e.*,

$$D(\tilde{I}, G) = \frac{(1 + \gamma) \cdot D_{\text{iou}} \cdot D_{\text{cov}}}{D_{\text{iou}} + \gamma \cdot D_{\text{cov}}}, \quad (3)$$

where  $\gamma$  is a weighting coefficient. Considering the noise and inaccuracy that may be caused by threshold segmentation, we put more emphasis on the completeness of covering regions when combining the IoU metric and COV metric, thus we set the  $\gamma$  to 0.3 to emphasis the  $D_{\text{cov}}$  over  $D_{\text{iou}}$  following the setting in [53].

To learn the potentiality of the depth map, we provide  $D(\tilde{I}, G)$  as the pseudo label  $g$  to guide the training of the regression process. Specifically, the top layers features of two branches' backbone are concatenated after passing through GAP, and then two fully connected layers are applied to obtain the estimation  $\hat{g}$ . The  $D(\tilde{I}, G)$  is only available in the training phase. As  $\hat{g}$  reflects the potentiality confidence of depth map, we introduce it in the GMA module to prevent the contamination from unreliable depth map in the fusion process, which will be explained in the GMA module.

#### C. Gated Multi-modality Attention Module

Taken into account that there exist complementarity and inconsistency of the cross-modal RGB-D data, directly integrating the cross-modal information may induce negative results, such as contaminations from unreliable depth maps. Besides, the features of the single modality usually are affluent in spatial or channel aspect, but also include information

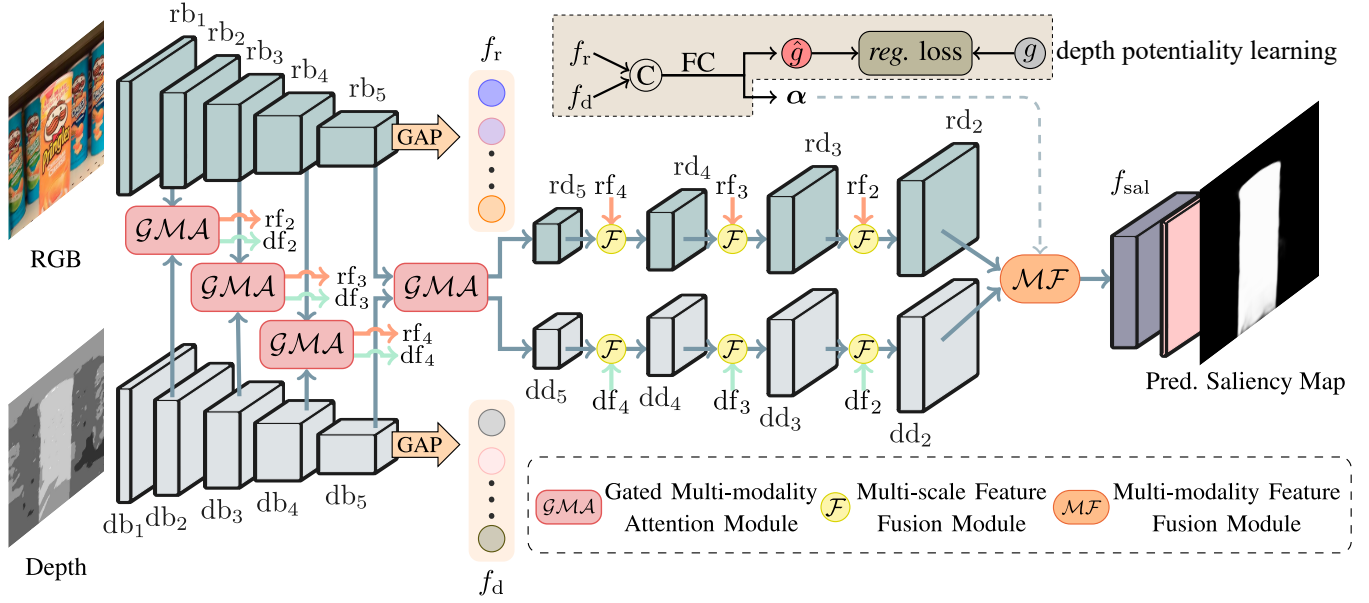


Fig. 2. Architecture of *DPANet*. For better visualization, we only display the modules and features of each stage.  $rb_i, db_i$  ( $i = 1, 2, \dots, 5$ ) denote the features generated by the backbone of the two branches respectively, and  $rd_i, dd_i$  ( $i = 5, 4, 3, 2$ ) represent the features of decoder stage.  $rf_i, df_i$  ( $i = 2, 3, 4, 5$ ,  $rf_5 = rd_5$ ,  $df_5 = dd_5$ ) refer to the output of the GMA module.  $f_{sal}$  is the generated final saliency map. “C” and “FC” refer to the concatenation and fully-connected layers respectively, and “GAP” represents the global average pooling.

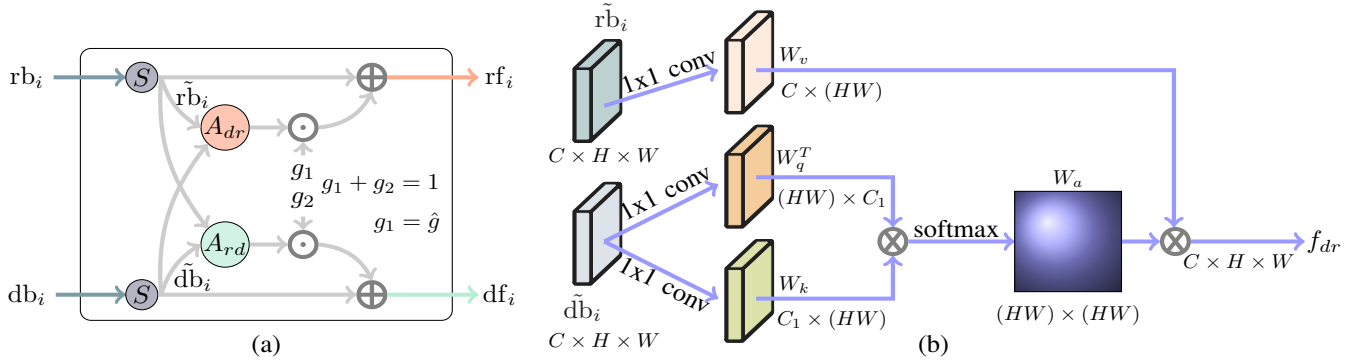


Fig. 3. Illustration of GMA module. (a) shows the construction of GMA module, and (b) represents the operation  $A_{dr}$ . The operation  $A_{rd}$  is symmetrical to the  $A_{dr}$  (exchange the position of  $\tilde{r}_i$  and  $\tilde{d}_i$ ). For conciseness, we just show the  $A_{dr}$ . And the  $\hat{g}$  refers to the prediction score of the depth potentiality.

redundancy. To cope with these issues, we design a GMA module that exploits the attention mechanism to automatically select and strengthen important features for saliency detection, and incorporate the gate controller into the GMA module to prevent the contamination from the unreliable depth map.

To reduce the redundancy of single-modal features and highlight the feature response on the salient regions, we apply spatial attention (see ‘S’ in Fig. 3) to the input feature  $rb_i$  and  $db_i$ , respectively. The process can be described as:

$$f = conv_1(f_{in}), \quad (4)$$

$$(W; B) = conv_2(f), \quad (5)$$

$$f_{out} = \delta(W \odot f + B), \quad (6)$$

where  $f_{in}$  represents the input feature of the RGB branch or depth branch (i.e.,  $rb_i$  or  $db_i$ ),  $conv_i$  ( $i = 1, 2$ ) refers to the convolution layer,  $\odot$  denotes element-wise multiplication,  $\delta$  is the ReLU activation function, and  $f_{out}$  represents the modified RGB/depth feature (i.e.,  $\tilde{r}_i$  or  $\tilde{d}_i$ ). The channels of modified

feature  $\tilde{r}_i, \tilde{d}_i$  are unified into 256 dimensions at each stage. Note that, the weights are not shared for the RGB and depth branches in our model.

Further, inspired by the success of self-attention [54], [55], we design two symmetrical attention sub-modules to capture long-range dependencies from a cross-modal perspective. Taking  $A_{dr}$  in Fig. 3 as an example,  $A_{dr}$  exploits the depth information to generate a spatial weight for RGB feature  $\tilde{r}_i$ , as depth cues usually can provide helpful information (e.g., the coarse location of salient objects) for RGB branch. Technically, we first apply  $1 \times 1$  convolution operation to project the  $\tilde{d}_i$  into  $W_q \in \mathbb{R}^{C_1 \times (HW)}$ ,  $W_k \in \mathbb{R}^{C_1 \times (HW)}$ , and project the  $\tilde{r}_i$  into  $W_v \in \mathbb{R}^{C \times (HW)}$ , where  $C, H, W$  refer to the channel, height, width of the feature  $W_v$ , respectively, and  $C_1$  is set to  $1/8$  of  $C$  for computation efficiency. We compute



the enhanced feature as follows:

$$W_a = \text{softmax}(W_q^T \otimes W_k), \quad (7)$$

$$f_{dr} = W_v \otimes W_a, \quad (8)$$

where  $\text{softmax}$  is applied in the column of  $W_a$ ,  $\otimes$  represents matrix multiplication. The enhanced feature  $f_{dr}$  is then reshaped into  $C \times H \times W$ . The another sub-module  $A_{rd}$  is symmetric to  $A_{dr}$ . These two attention modules aim to capture the long-range dependencies from a cross-modal perspective, where  $A_{dr}$  exploits the depth information to generate a spatial weight for the RGB feature, and  $A_{rd}$  refines the depth feature by using the spatial weight generated from the RGB feature.

Finally, we introduce the gates  $g_1$  and  $g_2$  with the constraint of  $g_1 + g_2 = 1$  to control the interaction of the enhanced features and modified features, which can be formulated as

$$\text{rf}_i = \tilde{\text{rb}}_i + g_1 \cdot f_{dr}, \quad (9)$$

$$\text{df}_i = \tilde{\text{db}}_i + g_2 \cdot f_{rd}. \quad (10)$$

where  $\text{rf}_i$  and  $\text{df}_i$  are the refined feature of the RGB and depth branches, respectively, which is used in the decoder stage. In equation (9), we use the gate to control the interaction of the enhanced features and modified features. For the enhanced feature, the weight  $g_1 = \hat{g}$ , where  $\hat{g}$  is learned under the supervision of the pseudo label  $g$ . When  $\hat{g}$  is closed to 1, it means that the depth map is highly reliable, and more depth information will be introduced to RGB branch to reduce the background disturbances; When  $\hat{g}$  is closed to 0, RGB branch will be the dominant branch and less depth information will be adopted, and the RGB information will play a more important role to prevent the contamination. Some feature visualization maps are shown in Fig. 7. Taking the third image as an example, the quality of the depth map is poor intuitively, so the enhanced feature  $f_{dr}$  is wrongly focused on the lower left areas. But with the constraint of the weight  $g_1$ , the encoder feature we get can effectively highlight the salient region and suppress the effect of enhanced feature  $f_{dr}$ . More details will be discussed in Section IV-E.

#### D. Multi-level Feature Fusion

Feature fusion plays a more critical role in RGB-D saliency detection due to cross-modal information, often directly affecting the performance. In order to obtain more comprehensive and discriminative fusion features, we consider two aspects for multi-level feature integration. First, the features at different scales contain different information, which can complement each other. Therefore, we use a multi-scale progressive fusion strategy to integrate the single-modal feature from coarse to fine. Second, for the multi-modality features, we exploit the designed GMA module to enhance the features separately instead of early fusing the RGB and depth features, which can reduce the interference of different modality. Finally, we aggregate the two modal features by using multi-modality feature fusion to obtain the saliency map.

**Multi-scale Feature Fusion.** Low-level features can provide more detail information, such as boundary, texture, and spatial structure, but may be sensitive to the background

noises. Contrarily, high-level features contain more semantic information, which is helpful to locate the salient object and suppress the noises. Different from previous works [13], [14] generally fuse the low-level features and high-level features by concatenation or summation operation, we adopt a more aggressive yet effective operation, *i.e.*, multiplication. The multiplication operation can strengthen the response of salient objects, meanwhile suppress the background noises. Specifically, taking the fusion of higher level feature  $\text{rd}_5$  and lower level feature  $\text{rf}_4$  as an instance, the multi-scale feature fusion can be described as

$$f_1 = \delta(\text{upsample}(\text{conv}_3(\text{rd}_5)) \odot \text{rf}_4), \quad (11)$$

$$f_2 = \delta(\text{conv}_4(\text{rf}_4) \odot \text{upsample}(\text{rd}_5)), \quad (12)$$

$$f_F = \delta(\text{conv}_5([f_1, f_2])), \quad (13)$$

where  $\text{upsample}$  is the up-sampling operation via bilinear interpolation, and  $[\cdot, \cdot]$  represents the concatenation operation. The fusion result  $f_F$  is exactly the higher level feature of the next fusion stage.

**Multi-modality Feature Fusion.** During the multi-modality feature fusion, we consider two issues: (1) How to select the most useful and complementary information from the RGB and depth features. Thus, we learn the weight  $\alpha$  to balance the complementarity when combining the complementary information from two modal data. The weight  $\alpha$  is learned from the top layers features of two branches backbone, which demonstrates the channel importance of the multi-modality features. (2) How to prevent the contamination caused by the unreliable depth map during fusing. Thus, the weight  $\hat{g}$  is used to control the introduction ratio of the depth information. The weight  $\hat{g}$  is learned under the supervision of the pseudo label  $g$  and reflects the potentiality confidence of depth map.

Specifically, to fuse the cross-modal features  $\text{rd}_2$  and  $\text{dd}_2$ , we design a weighted channel attention mechanism to automatically select useful channels, which can be formulated as

$$f_3 = \alpha \odot \text{rd}_2 + \hat{g} \cdot (1 - \alpha) \odot \text{dd}_2, \quad (14)$$

$$f_4 = \text{rd}_2 \odot \text{dd}_2, \quad (15)$$

$$f_{\text{sal}} = \delta(\text{conv}([f_3, f_4])), \quad (16)$$

where  $\alpha \in \mathbb{R}^{256}$  is the weight vector learned from RGB and depth information (see Fig. 2),  $\hat{g}$  is the learned weight of the gate as mentioned before. The equation (15) reflects the common response for salient objects, while equation (14) combines the two modal features via channel selection ( $\alpha$ ) and gate mechanism ( $\hat{g}$ ) for considering the complementarity and inconsistency.

#### E. Loss Function

For training the network, we consider the classification loss and regression loss to define the loss function, where the classification loss is used to constrain the saliency prediction, and the regression loss aims to model the depth potentiality response.

**Classification Loss.** In saliency detection, binary cross-entropy loss is commonly adopted to measure the relation

between predicted saliency map and the ground truth, which can be defined as

$$\ell = -\frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W [G_{ij} \log(S_{ij}) + (1 - G_{ij}) \log(1 - S_{ij})], \quad (17)$$

where  $H$ ,  $W$  refer to the height and width of the image respectively,  $G$  denotes the ground truth, and  $S$  represents the predicted saliency map. To facilitate the optimization of the proposed network, we add auxiliary loss at four decoder stages. Specifically, a  $3 \times 3$  convolution layer is applied for each stage ( $rd_i$ ,  $dd_i$ ,  $i = 5, 4, 3, 2$ ) to squeeze the channel of the output feature maps to 1. Then these maps are up-sampled to the same size as the ground truth via bilinear interpolation and sigmoid function is used to normalize the predicted values into  $[0, 1]$ . Thus, the whole classification loss consists of two parts, *i.e.*, the dominant loss corresponding to the output and the auxiliary loss of each sub-stage.

$$\ell_{cls} = \ell_{dom} + \sum_{i=1}^8 \lambda_i \ell_{aux}^i, \quad (18)$$

where  $\lambda_i$  denotes the weight of different loss, and  $\ell_{dom}$ ,  $\ell_{aux}^i$  denote the dominant and auxiliary loss, respectively. The auxiliary loss branches only exist during the training stage.

**Regression Loss.** To model the potentiality of depth map, the smooth L1 loss [56] is used as the supervision signal. The smooth L1 loss is defined as

$$\ell_{reg} = \begin{cases} 0.5(g - \hat{g})^2, & \text{if } |g - \hat{g}| < 1 \\ |g - \hat{g}| - 0.5, & \text{otherwise} \end{cases}, \quad (19)$$

where  $g$  is the pseudo label as mentioned in the depth potentiality perception, and  $\hat{g}$  denotes the estimation of the network as shown in Fig. 2.

**Final Loss.** The final loss is the linear combination of the classification loss and regression loss,

$$\ell_{final} = \ell_{cls} + \lambda \ell_{reg}, \quad (20)$$

where  $\lambda$  is weight of  $\ell_{reg}$ , which is set to 1 in our model. The whole training process is conducted in an end-to-end way.

## IV. EXPERIMENTS

### A. Datasets

We evaluate the proposed method on 8 public RGB-D SOD datasets with the corresponding pixel-wise ground-truth.

**NJUD** [44] consists of 2003 RGB images and corresponding depth images with diverse objects and complex scenarios. The depth images are estimated from the stereo images. **NLPR** [41] contains 1000 RGB-D images captured by Kinect. Moreover, there exist multiple salient objects in an image of this dataset. **STEREO797** [57] contains 797 stereoscopic images collected from the Internet, where the depth maps are estimated from the stereo images. **LFSD** [58] includes 100 RGB-D images, in which the depth map is captured by Lytro light field camera. **RGBD135** [46] contains 135 RGB-D images captured by Kinect simultaneously. **SSD** [59] contains 80 images picked up from three stereo movies, where the depth map is generated by depth estimation method. **DUT**

[50] consists of 1200 paired images containing more complex scenarios, such as multiple or transparent objects, *etc.*. This dataset is split into 800 training data and 400 testing data. **SIP** [51] contains 929 high-resolution person RGB-D images captured by Huawei Meta10.

### B. Evaluation Metrics

To quantitatively evaluate the effectiveness of the proposed method, precision-recall (PR) curves, F-measure ( $F_\beta$ ) score and curves, Mean Absolute Error (MAE), and S-measure ( $S_m$ ) are adopted. Thresholding the saliency map at a series of values, pairs of precision-recall value can be computed by comparing the binary saliency map with the ground truth. The F-measure<sup>1</sup> is a comprehensive metric that takes both precision and recall into account, which is defined as:

$$F_\beta = \frac{(1 + \beta^2) \cdot Precision \cdot Recall}{\beta^2 \cdot Precision + Recall} \quad (21)$$

where  $\beta^2$  is set to 0.3 to emphasize the precision over recall, as suggested by [53].

MAE is defined as the average pixel-wise absolute difference between the saliency map and the ground truth:

$$MAE = \frac{1}{H \times W} \sum_{y=1}^H \sum_{x=1}^W |S(x, y) - G(x, y)| \quad (22)$$

where  $S$  is the saliency map,  $G$  denotes the ground truth,  $H$  and  $W$  are the height and width of the saliency map, respectively.

S-measure evaluates the structural similarity between the saliency map and the ground truth [60], which is defined as:

$$S_m = \alpha * S_o + (1 - \alpha) * S_r \quad (23)$$

where  $\alpha$  is set to 0.5 to balance the object-aware structural similarity ( $S_o$ ) and region-aware structural similarity ( $S_r$ ).

### C. Implementation Details

Following [49], we take 1400 images from NJUD [44] and 650 images from NLPR [41] as the training, and 100 images from NJUD dataset and 50 images from NLPR dataset as the validation set. Moreover, we unified the depth map with the necessary flips before training. To reduce the overfitting, we use multi-scale resizing and random horizontal flipping augmentation. During the inference stage, images are simply resized to  $256 \times 256$ , and then fed into the network to obtain prediction without any other post-processing (*e.g.*, CRF [36]) or pre-processing techniques (*e.g.*, HHA [35]). We use Pytorch to implement our model, and the ResNet-50 [61] is used as our backbone. Mini-batch stochastic gradient descent (SGD) is used to optimize the network with the batch size of 32, the momentum of 0.9, and the weight decay of  $5e-4$ . The loss weights are set as  $\{\lambda_1, \lambda_2, \dots, \lambda_8\} = \{1.0, 0.8, 0.6, 0.4, 1.0, 0.8, 0.6, 0.4\}$ . We use the warm-up and linear decay strategies with the maximum learning rate  $5e-3$  for the backbone and 0.05 for other parts and stop training after 30 epochs.

<sup>1</sup>We report the maximum F-measure in our results

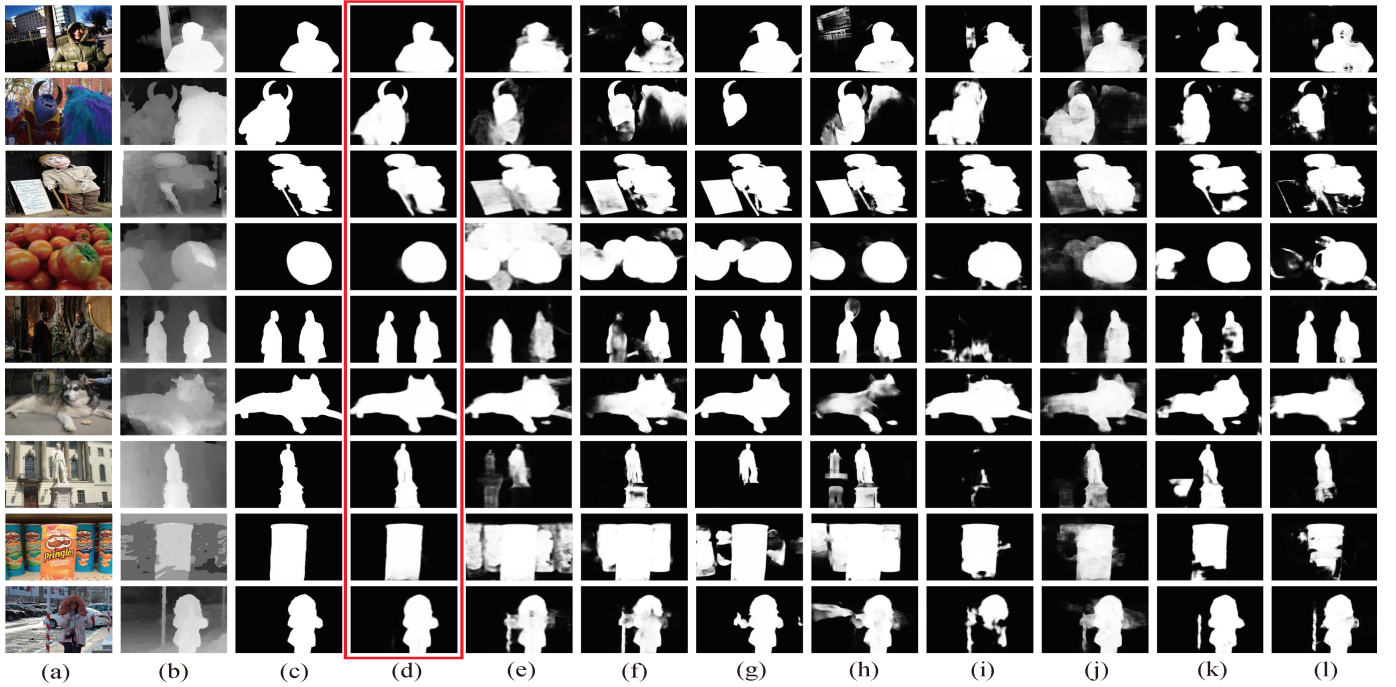


Fig. 4. Qualitative comparison of the proposed approach with some state-of-the-art RGB and RGB-D SOD methods, in which our results are highlighted by a red box. (a) RGB image. (b) Depth map. (c) GT. (d) DPANet. (e) PiCAR. (f) PoolNet. (g) BASNet. (h) EGNet. (i) CPPF. (j) PDNet. (k) DMRA. (l) AF-Net.

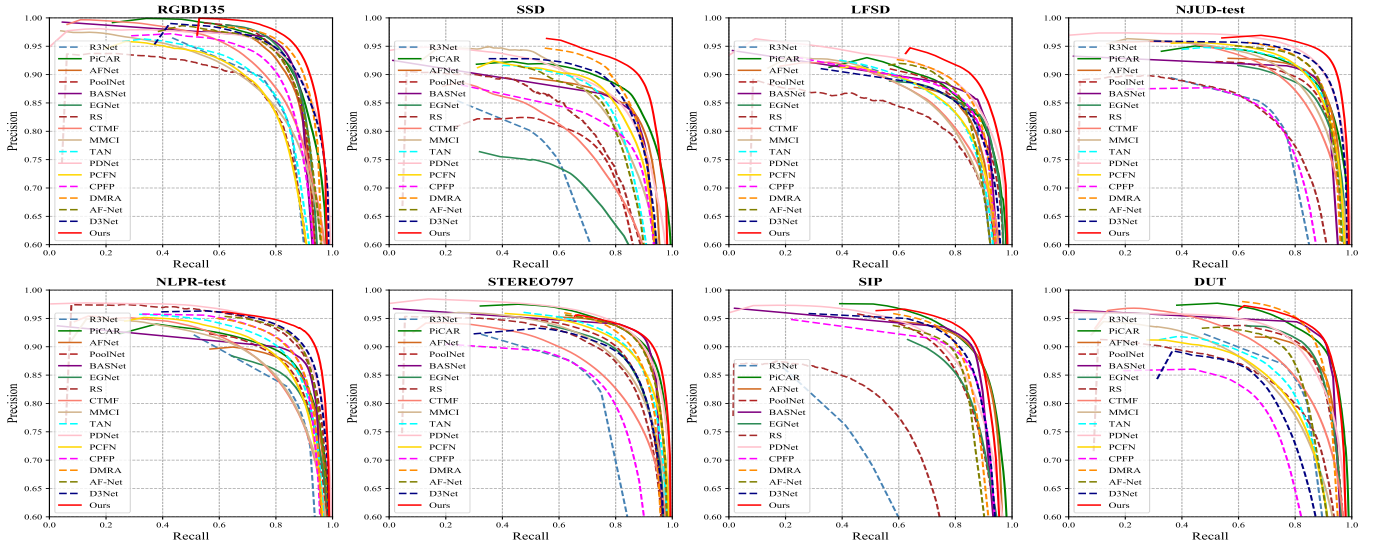


Fig. 5. Illustration of PR curves on different datasets. The closer the PR curve is to (1, 1), the better performance of the method.

#### D. Comparison with the State-of-the-arts

We compare the proposed model with 16 state-of-the-art methods, including 10 RGB-D saliency models (*i.e.*, D<sup>3</sup>Net [51], AF-Net [62], DMRA [50], CPPF [33], PCFN [49], PDNet [48], TAN [63], MMCI [64], CTMF [52], and RS [65]) and 6 latest RGB saliency models (*i.e.*, EGNet [15], BASNet [14], PoolNet [13], AFNet [16], PiCAR [12], and R<sup>3</sup>Net [11]). For fair comparisons, we use the released code and default parameters to reproduce the saliency maps or the saliency maps provided by the authors.

**1) Qualitative Evaluation:** To further illustrate the advantages of the proposed method, we provide some visual examples of different methods. As shown in Fig. 4, our proposed network obtains a superior result with precise saliency location, clean background, complete structure, and sharp boundaries, and also can address various challenging scenarios, such as low contrast, complex scene, background disturbance, and multiple objects. To be specific,

(a) Our model achieves more complete structure and sharp boundaries in the results. For example, in the second image, the horns and body of the salient object cannot be completely detected by the comparison methods, such as PoolNet [13],

TABLE I

PERFORMANCE COMPARISON ON 8 PUBLIC DATASETS. THE BEST RESULTS ON EACH DATASET ARE HIGHLIGHTED IN **BOLDFACE**. FROM TOP TO BOTTOM: OUR METHOD, CNN-BASED RGB-D SOD METHODS, AND THE LATEST RGB SOD METHODS.

Method	RGBD135 Dataset			SSD Dataset			LFSD Dataset			NJUD-test Dataset		
	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$
DPANet (ours)	<b>0.933</b>	<b>0.919</b>	<b>0.023</b>	<b>0.898</b>	<b>0.889</b>	<b>0.042</b>	<b>0.891</b>	<b>0.864</b>	<b>0.072</b>	<b>0.930</b>	<b>0.921</b>	<b>0.035</b>
D <sup>3</sup> Net (TNNLS20)	0.909	0.898	0.031	0.861	0.857	0.058	0.840	0.825	0.095	0.909	0.900	0.047
AF-Net (Arxiv19)	0.904	0.892	0.033	0.828	0.815	0.077	0.857	0.818	0.091	0.900	0.883	0.053
DMRA (ICCV19)	0.921	0.911	0.026	0.874	0.857	0.055	0.865	0.831	0.084	0.900	0.880	0.052
CPFP (CVPR19)	0.882	0.872	0.038	0.801	0.807	0.082	0.850	0.828	0.088	0.799	0.798	0.079
PCFN (CVPR18)	0.842	0.843	0.050	0.845	0.843	0.063	0.829	0.800	0.112	0.887	0.877	0.059
PDNet (ICME19)	0.906	0.896	0.041	0.844	0.841	0.089	0.865	0.846	0.107	0.912	0.897	0.060
TAN (TIP19)	0.853	0.858	0.046	0.835	0.839	0.063	0.827	0.801	0.111	0.888	0.878	0.060
MMCI (PR19)	0.839	0.848	0.065	0.823	0.813	0.082	0.813	0.787	0.132	0.868	0.859	0.079
CTMF (TCyb18)	0.865	0.863	0.055	0.755	0.776	0.100	0.815	0.796	0.120	0.857	0.849	0.085
RS (ICCV17)	0.841	0.824	0.053	0.783	0.750	0.107	0.795	0.759	0.130	0.796	0.741	0.120
EGNet (ICCV19)	0.913	0.892	0.033	0.704	0.707	0.135	0.845	0.838	0.087	0.867	0.856	0.070
BASNet (CVPR19)	0.916	0.894	0.030	0.842	0.851	0.061	0.862	0.834	0.084	0.890	0.878	0.054
PoolNet (CVPR19)	0.907	0.885	0.035	0.764	0.749	0.110	0.847	0.830	0.095	0.874	0.860	0.068
AFNet (CVPR19)	0.897	0.878	0.035	0.847	0.859	0.058	0.841	0.817	0.094	0.890	0.880	0.055
PiCAR (CVPR18)	0.907	0.890	0.036	0.864	0.871	0.055	0.849	0.834	0.104	0.887	0.882	0.060
R <sup>3</sup> Net (IJCAI18)	0.857	0.845	0.045	0.711	0.672	0.144	0.843	0.818	0.089	0.805	0.771	0.105

TABLE II  
CONTINUATION OF TABLE I.

Method	NLPR-test Dataset			STEREO797 Dataset			SIP Dataset			DUT Dataset		
	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$
DPANet (ours)	<b>0.922</b>	<b>0.928</b>	<b>0.024</b>	<b>0.915</b>	<b>0.911</b>	<b>0.041</b>	<b>0.904</b>	<b>0.883</b>	<b>0.051</b>	<b>0.916</b>	0.899	0.048
D <sup>3</sup> Net (TNNLS20)	0.907	0.912	0.030	0.869	0.868	0.058	0.881	0.860	0.063	0.797	0.775	0.097
AF-Net (Arxiv19)	0.904	0.903	0.032	0.905	0.893	0.047	0.870	0.844	0.071	0.862	0.831	0.077
DMRA (ICCV19)	0.887	0.889	0.034	0.895	0.874	0.052	0.883	0.850	0.063	0.913	0.880	0.052
CPFP (CVPR19)	0.888	0.888	0.036	0.815	0.803	0.082	0.870	0.850	0.064	0.771	0.760	0.102
PCFN (CVPR18)	0.864	0.874	0.044	0.884	0.880	0.061	–	–	–	0.809	0.801	0.100
PDNet (ICME19)	0.905	0.902	0.042	0.908	0.896	0.062	0.863	0.843	0.091	0.879	0.859	0.085
TAN (TIP19)	0.877	0.886	0.041	0.886	0.877	0.059	–	–	–	0.824	0.808	0.093
MMCI (PR19)	0.841	0.856	0.059	0.861	0.856	0.080	–	–	–	0.804	0.791	0.113
CTMF (TCyb18)	0.841	0.860	0.056	0.827	0.829	0.102	–	–	–	0.842	0.831	0.097
RS (ICCV17)	0.900	0.864	0.039	0.857	0.804	0.088	–	–	–	0.807	0.797	0.111
EGNet (ICCV19)	0.845	0.863	0.050	0.872	0.853	0.067	0.846	0.825	0.083	0.888	0.867	0.064
BASNet (CVPR19)	0.882	0.894	0.035	0.914	0.900	0.041	0.894	0.872	0.055	0.912	<b>0.902</b>	<b>0.041</b>
PoolNet (CVPR19)	0.863	0.873	0.045	0.876	0.854	0.065	0.856	0.836	0.079	0.883	0.864	0.067
AFNet (CVPR19)	0.865	0.881	0.042	0.905	0.895	0.045	0.891	0.876	0.055	0.880	0.868	0.065
PiCAR (CVPR18)	0.872	0.882	0.048	0.906	0.903	0.051	0.890	0.878	0.060	0.903	0.892	0.062
R <sup>3</sup> Net (IJCAI18)	0.832	0.846	0.049	0.811	0.754	0.107	0.641	0.624	0.158	0.841	0.812	0.079

CPFP [33], and the background regions (*e.g.*, cartoon character on the right) are wrongly retained. Similarly, in the sixth image, most methods fail to detect the left front leg of the salient dog (*e.g.*, AF-Net [62], DMRA [50], EGNet [15], and BASNet [14]), and the detected object boundaries are blurred and inaccurate (*e.g.*, PDNet [48] and PiCAR [12]). By contrast, our method yields a more complete structure and sharp boundaries.

(b) Our model can address the complex and low contrast scenes. For example, in the fifth image, the low light makes the person in black on the left very close to the background, so many methods cannot accurately detect him, such as PoolNet [13], CPFP [33]. By contrast, our model can detect those two

persons with more complete structure and clean background. In the last image, the background is very complex, including a nearby indicator pole and multiple cars in the distance. Thus, the nearby indicator pole is detected as salient object by most methods, while our method can successfully suppress this region and obtain a better result with complete structure and clear boundaries.

(c) Our model can well handle the disturbance of similar appearances between the salient object and backgrounds. In the fourth image, there are many tomatoes in the image, but the slightly green tomato in the front is more prominent than others. Without the help of depth cues, all RGB SOD models fail to suppress these interferences, resulting in inaccurate



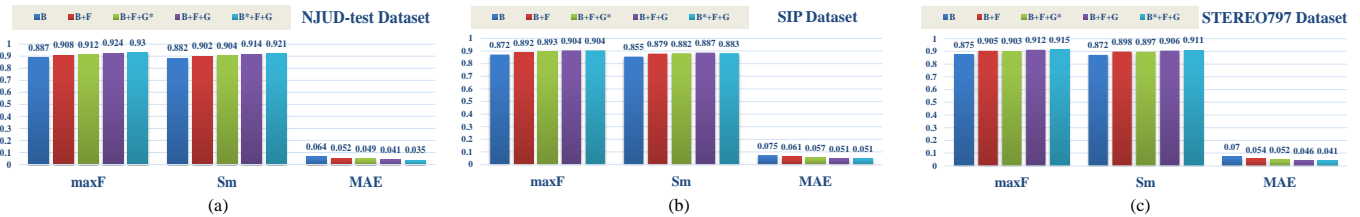


Fig. 6. Ablation study of module verification on NJUD-test, SIP, and STEREO797 datasets.

detection results. Moreover, the existing RGB-D SOD methods also cannot address this challenging case very well. By contrast, our method shows a competitive advantage in terms of the completeness, sharpness, and accuracy. Analogously, our model can integrally highlight the sculpture and its pedestal in the seventh image, even the salient object has the similar color appearances with the background.

(d) Our model can produce more robust result even when confronting with the inaccurate or blurred depth information (e.g., the third and sixth images), which illustrates the power of the GMA module. In these challenging scenarios, it can be seen that the network is capable to utilize cross-modal complementary information and prevent the contamination from the unreliable depth map.

2) **Quantitative Evaluation:** For a more intuitive comparison of algorithm performance, we report the PR curves in Fig. 5 and the quantitative metrics including the maximum F-measure, S-measure, and MAE score in Tables I and II. From the PR curves shown in Fig. 5, we can see that the proposed method achieves both higher precision and recall scores against other compared methods with an obvious margin on the eight datasets. As reported in Tables I and II, our method outperforms all the compared methods in terms of all the measurements, except that the MAE score on the DUT dataset achieves the second-best performance. It is worth mentioning that the performance gain is significant against the compared methods. For example, compared with the latest RGB-D SOD method *DMRA*, our algorithm still achieves competitive performance in the case of a small number of training samples (i.e., the training samples do not include 800 training data of the DUT dataset). On the SSD dataset, compared with the *DMRA* method, the percentage gain of our *DPANet* achieves 2.7% in terms of maximum F-measure, 3.7% in terms of S-measure, and 23.6% in terms of MAE score. Compared with the *second best* method in different datasets, our method still achieves obvious performance gain. For example, on the NJUD-test dataset, the percentage gain reaches 2.0% for F-measure, 2.3% for S-measure, and 25.5% for MAE score. On the NLPD-test dataset, the *minimum percentage gain* reaches 1.7% for F-measure, 1.8% for S-measure, and 20.0% for MAE score. Thus, all the quantitative measures demonstrate the effectiveness of the proposed model.

3) **Run-time Comparison:** We provide the run-time comparison results in the Table III. From it, we can see that our model is faster than most of these deep learning based RGB-D SOD methods, which demonstrates the efficiency of the

TABLE III  
COMPARISONS OF INFERENCE TIME OF DIFFERENT DEEP LEARNING BASED RGB-D SOD METHODS.

	CTMF	MMCI	TAN	PDNet	PCFN
Time (s)	0.63	0.05	0.07	0.07	0.06
	CPFP	AF-Net	DMRA	D <sup>3</sup> Net	Ours
Time (s)	0.17	0.03	0.06	0.05	0.03

proposed *DPANet*.

### E. Ablation Study

1) **Module Verification:** In this section, we conduct the ablation study to demonstrate the effectiveness of each key components designed in the proposed network on three different datasets: NJUD-test, SIP, and STEREO797, and the results are shown in Fig. 6. We choose the network that removes the GMA module and regression loss, replaces the multi-scale feature fusion module with concatenation, and replaces the multi-modality feature fusion with multiplication and summation as the baseline (denoted as ‘B’). From Fig. 6, compared the ‘B’ with the ‘B+F’, the multi-scale feature fusion (denoted as ‘F’) improves the baseline about 2 ~ 3 points in terms of the maximum F-measure. After adding the GMA module without regression loss constrain (denoted as ‘G\*’), the F-measure rises up to 0.912 on NJUD-test dataset, which is comparable with the state-of-the-art results. Furthermore, the performance is significantly enhanced after adding regression loss (B+F+G), which achieves the percentage gain of 16.3% on NJUD-test compared with the (B+F+G\*) case in terms of MAE score. Finally, adding the gate module in the multi-modality feature fusion (i.e., B\*+F+G), our method yields the best performance result with the percentage gain of 4.8% in terms of F-measure and 44.4% in terms of MAE compared with the original baseline on NJUD-test dataset. The experiments on another two datasets, i.e., SIP, and STEREO797 also demonstrate the effectiveness of the designed components.

To better understand the attention mechanism designed in the GMA module, we visualize some feature maps and the corresponding heat-maps in Fig. 7. Taking the fourth GMA module as an example, as expected, the GMA module should learn the cross-modal complementarity from a cross-modal perspective and prevent the contamination from the unreliable depth map. Recall that the output of the fourth GMA module is denoted as  $rf_5 = r\hat{b}_5 + g_1 \cdot f_{dr}$ ,  $df_5 = d\hat{b}_5 + g_2 \cdot f_{rd}$ , where  $g_1 = \hat{g}$ , and  $g_1 + g_2 = 1$ . As Fig. 7 shows,

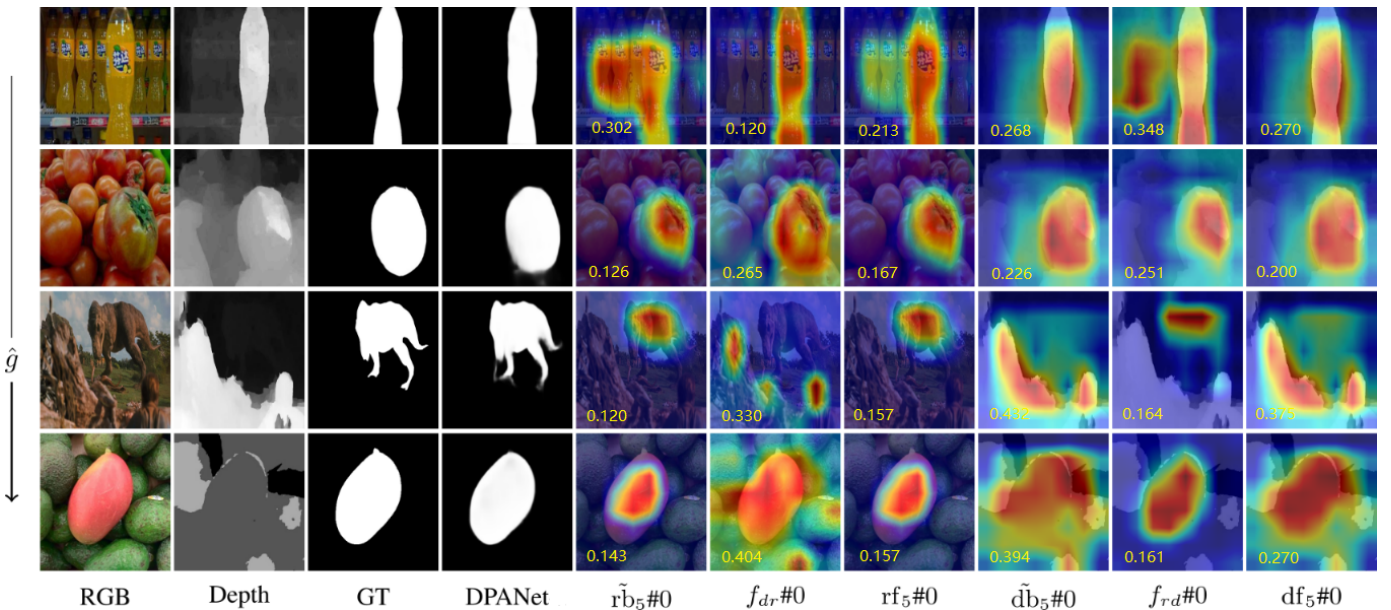


Fig. 7. Visualization of the GMA module. “#0” refers to the first channel of the features. The number at the bottom left of each image is the MAE value of the corresponding feature map.

TABLE IV  
ABLATION STUDIES ON NJUD-TEST, SIP, AND STEREO797 DATASETS.

	NJUD-test Dataset			SIP Dataset			STEREO797 Dataset		
	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$
DPANet	0.930	0.921	0.035	0.904	0.883	0.051	0.915	0.911	0.041
concatenation	0.919	0.914	0.039	0.904	0.876	0.056	0.912	0.905	0.044
summation	0.923	0.915	0.038	0.906	0.881	0.054	0.910	0.904	0.045
hard manner	0.908	0.902	0.047	0.893	0.868	0.064	0.905	0.899	0.050
w/o depth	0.908	0.903	0.043	0.864	0.837	0.074	0.913	0.908	0.042

- when the depth map is reliable, but the RGB image encounters interference from the similar background (*e.g.*, the first and second images), the features of the RGB branch ( $r_{b5}\#0$ ) usually fail to well focus on the area of salient object. By contrast, the features of depth map ( $d_{b5}\#0$ ) can provide complementary information to enhance the feature maps ( $f_{dr}\#0$ ) by suppressing the background noises. Further, the features of the RGB branch ( $r_{b5}\#0$ ) and the enhanced features ( $f_{dr}\#0$ ) are combined with the weight of  $g_1$  to obtain the features ( $rf_5\#0$ ) of the decoder stage.
- when the depth map tends to unreliable (*e.g.*, the third and fourth images), we can see that the imperfect depth information almost has no impact on the features of the RGB branch thanks to the gate controller  $g_2$  (compare  $r_{b5}\#0$  and  $rf_5\#0$ ), and introducing the RGB information to the depth branch leads the features of depth branch more convergent (from  $d_{b5}\#0$  to  $df_5\#0$ ).

In a summary, the designed GMA module can learn the complementarity from a cross-modal perspective and prevent the contamination from the unreliable depth map. The output features of the GMA module will be aggregated progressively in the decoder stage via multi-scale feature fusion separately.

Finally, to obtain the saliency map, the features of the RGB and depth branches are aggregated via the multi-modality feature fusion. All these ablation studies demonstrate the effectiveness of main components in our network, including the GMA module, regression loss, and multi-level feature fusion.

2) **Fusion Method Verification:** For fusing the low-level features and high-level features, we adopt a more aggressive yet effective operation to replace the concatenation or summation operation, *i.e.*, multiplication. The main reason is that the multiplication operation can strengthen the response of salient objects, meanwhile suppress the background noises. Some comparisons of the different fusion methods are shown in the Table IV, replacing the multiplication in equations (11) and (12) with the concatenation or summation operation. From it, we can see that using multiplication is more effective than concatenation or summation. These experiments verified the hypothesis that multiplication can strengthen the response of the salient object and suppress the background noise.

3) **Depth-related Information Verification:** In this subsection, we conduct the ablation studies of the depth-related information.

First, we investigate the difference between the soft manner

TABLE V  
DEPTH POTENTIALITY (DP) SCORES ON DIFFERENT DATASETS

	RGBD135	SSD	LFSD	NJUD-test
DP score	0.495	0.675	0.755	0.657
minPG-F	1.3%	2.7%	1.9%	2.0%
	NLPR-test	STEREO797	SIP	DUT
DP score	0.481	0.635	0.683	0.650
minPG-F	1.7%	0.1%	1.1%	0.3%

and hard manner in the learning stage of the depth potentiality. The soft manner is the default setting that we use the regression as the supervision. The hard manner changes the regression as a classification problem, *i.e.*, the pseudo label  $g$  is binarized by a fixed threshold 0.5, and binary cross entropy loss is adopted as the supervision signal, which forces the predict depth potentiality  $\hat{g}$  tending to 0 or 1 rather than a smooth contiguous value. From Table IV, we can see that the soft manner (*DPANet*) is better than the hard manner (row ‘hard manner’) in a holistic view. In addition, the soft manner is not affected by the threshold setting and is more robust.

Second, in order to verify the importance of depth information in RGB-D SOD, we report the experimental result of using only RGB image, which is denoted as ‘w/o depth’ in Table IV. As shown in Table IV, with the help of the depth map, the final saliency performance is obviously improved. For example, on the NJUD-test dataset, the F-measure is improved from 0.908 to 0.930 with the percentage gain of 2.4%, and the MAE score is improved from 0.043 to 0.035 with the percentage gain of 18.6%.

Third, to further evaluate the depth potentiality, we calculate the mean value of the depth potentiality (DP) score using equations (1)-(3) on 8 datasets, which is shown in the Table V. In Table V, we also report the minimum percentage gain of the F-measure, which is denoted as ‘minPG-F’. For the dataset with the worst depth map quality (*i.e.*, NLPR-test dataset with the DP score of 0.481), the minimum gain of our *DPANet* reaches 1.7%. In general, our *DPANet* can well handle the poor depth perception potential as well as the good.

#### F. Failure Cases

For future researchers to develop better algorithms, we provide some failure cases as shown in Fig. 8. It can be seen that it is difficult to locate salient objects perfectly in the following three aspects, whether it is based on the RGB SOD method (BASNet [14]) or the RGB-D SOD method (Ours and DMRA [50]): (1) Long-distance small and multiple salient objects. In the first image, when there are obvious differences in the size of salient targets in a scene, especially when the salient objects are multiple and small, and the depth map cannot provide the effective depth information of long-distance objects, it is difficult to completely detect all salient objects. (2) The conflict between the depth information and salient objects. In the second image, the highlighted object in the depth map (*i.e.*, the close-range person on the right) is not the final salient object. For this case, it is difficult for the algorithm to suppress the interference of close-range

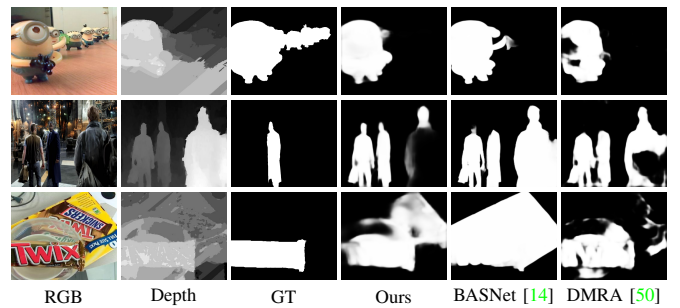


Fig. 8. Failure examples.

objects, leading to the false alarms. (3) The complex and cluttered background regions. When confronting with the complex scenarios (*e.g.*, the third image), especially when the depth map fails to provide accurate information, our algorithm is difficult to effectively extract the salient object from the complex background.

#### V. CONCLUSION

In this paper, we propose a novel framework *DPANet* to achieve RGB-D SOD. Considering the contamination from the unreliable depth map, we model a saliency-orientated depth potentiality perception module to evaluate the potentiality of the depth map and weaken the contamination. To effectively aggregate the cross-modal complementarity, we propose a GMA module to highlight the saliency response and regulate the fusion rate of the cross-modal information. Finally, the multi-stage and multi-modality feature fusion are used to generate the discriminative RGB-D features and produce the saliency map. Experiments on eight RGB-D datasets demonstrate that the proposed network outperforms other 16 state-of-the-art methods under different evaluation metrics.

#### REFERENCES

- [1] W. Wang, Q. Lai, H. Fu, J. Shen, and H. Ling, “Salient object detection in the deep learning era: An in-depth survey,” in *arXiv preprint arXiv:1904.09146*, 2019. 1
- [2] R. Cong, J. Lei, H. Fu, M.-M. Cheng, W. Lin, and Q. Huang, “Review of visual saliency detection with comprehensive information,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 28, no. 10, pp. 2941–2959, 2019. 1
- [3] W. Wang, J. Shen, R. Yang, and F. Porikli, “Saliency-aware video object segmentation,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 1, pp. 20–33, 2018. 1
- [4] R. Zhao, W. Ouyang, and X. Wang, “Unsupervised salience learning for person re-identification,” in *CVPR*, 2013, pp. 3586–3593. 1
- [5] F. Zhang, B. Du, and L. Zhang, “Saliency-guided unsupervised feature learning for scene classification,” *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 4, pp. 2175–2184, 2014. 1
- [6] W. Wang, J. Shen, Y. Yu, and K.-L. Ma, “Stereoscopic thumbnail creation via efficient stereo saliency detection,” *IEEE Trans. Vis. Comput. Graph.*, vol. 23, no. 8, pp. 2014–2027, 2017. 1
- [7] Q. Jiang, F. Shao, W. Lin, K. Gu, G. Jiang, and H. Sun, “Optimizing multi-stage discriminative dictionaries for blind image quality assessment,” *IEEE Trans. Multimedia*, vol. 20, no. 8, pp. 2035–2048, 2018. 1
- [8] C. Li, J. Guo, R. Cong, Y. Pang, and B. Wang, “Underwater image enhancement by dehazing with minimum information loss and histogram distribution prior,” *IEEE Trans. Image Process.*, vol. 25, no. 12, pp. 5664–5677, 2016. 1
- [9] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *NeurIPS*, 2012, pp. 1097–1105. 1



- [10] Q. Hou, M.-M. Cheng, X. Hu, A. Borji, Z. Tu, and P. H. Torr, "Deeply supervised salient object detection with short connections," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 4, pp. 815–828, 2019. **1, 2**
- [11] Z. Deng, X. Hu, L. Zhu, X. Xu, J. Qin, G. Han, and P.-A. Heng, "R<sup>3</sup>Net: Recurrent residual refinement network for saliency detection," in *IJCAI*, 2018, pp. 684–690. **1, 2, 7**
- [12] N. Liu, J. Han, and M.-H. Yang, "PiCANet: Learning pixel-wise contextual attention for saliency detection," in *CVPR*, 2018, pp. 3089–3098. **1, 2, 7, 8**
- [13] J.-J. Liu, Q. Hou, M.-M. Cheng, J. Feng, and J. Jiang, "A simple pooling-based design for real-time salient object detection," in *CVPR*, 2019, pp. 3917–3926. **1, 2, 5, 7, 8**
- [14] X. Qin, Z. Zhang, C. Huang, C. Gao, M. Dehghan, and M. Jagersand, "BASNet: Boundary-aware salient object detection," in *CVPR*, 2019, pp. 7479–7489. **1, 2, 5, 7, 8, 11**
- [15] J.-X. Zhao, J.-J. Liu, D.-P. Fan, Y. Cao, J. Yang, and M.-M. Cheng, "EGNet: Edge guidance network for salient object detection," in *ICCV*, 2019, pp. 8778–8787. **1, 2, 7, 8**
- [16] M. Feng, H. Lu, and E. Ding, "Attentive feedback network for boundary-aware salient object detection," in *CVPR*, 2019, pp. 1623–1632. **1, 2, 7**
- [17] T. Zhao and X. Wu, "Pyramid feature attention network for saliency detection," in *CVPR*, 2019, pp. 3085–3094. **1, 2**
- [18] Z. Chen, Q. Xu, R. Cong, and Q. Huang, "Global context-aware progressive aggregation network for salient object detection," in *AAAI*, 2020, pp. 10599–10606. **1, 2**
- [19] J. Zhang, T. Zhang, Y. Dai, M. Harandi, and R. Hartley, "Deep unsupervised saliency detection: A multiple noisy labeling perspective," in *CVPR*, 2018, pp. 9029–9038. **1, 2**
- [20] Y. Zeng, Y. Zhuge, H. Lu, L. Zhang, M. Qian, and Y. Yu, "Multi-source weak supervision for saliency detection," in *CVPR*, 2019, pp. 6074–6083. **1, 2**
- [21] Q. Wang, Y. Yuan, P. Yan, and X. Li, "Saliency detection by multiple-instance learning," *IEEE Trans. Cybern.*, vol. 43, no. 2, pp. 660–672, 2013. **1**
- [22] Q. Wang, M. Chen, F. Nie, and X. Li, "Detecting coherent groups in crowd scenes by multiview clustering," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 1, pp. 46–58, 2018. **1**
- [23] Q. Zhang, R. Cong, J. Hou, C. Li, and Y. Zhao, "CoADNet: Collaborative aggregation-and-distribution networks for co-salient object detection," in *Proc. NeurIPS*, 2020. **1**
- [24] R. Cong, J. Lei, C. Zhang, Q. Huang, X. Cao, and C. Hou, "Saliency detection for stereoscopic images based on depth confidence analysis and multiple cues fusion," *IEEE Signal Process. Lett.*, vol. 23, no. 6, pp. 819–823, 2016. **1**
- [25] C. Choi and H. I. Christensen, "RGB-D object tracking: A particle filter approach on GPU," in *IROS*, 2013, pp. 1084–1091. **1**
- [26] R. Cong, J. Lei, H. Fu, Q. Huang, X. Cao, and C. Hou, "Co-saliency detection for RGBD images based on multi-constraint feature matching and cross label propagation," *IEEE Trans. Image Process.*, vol. 27, no. 2, pp. 568–579, 2018. **1**
- [27] R. Cong, J. Lei, H. Fu, W. Lin, Q. Huang, X. Cao, and C. Hou, "An iterative co-saliency framework for RGBD images," *IEEE Trans. Cybern.*, vol. 49, no. 1, pp. 233–246, 2019. **1**
- [28] F. Li, R. Cong, H. Bai, and Y. He, "RGB-D salient object detection with cross-modality modulation and selection," in *Proc. IJCAI*, 2020, pp. 534–543. **1**
- [29] R. Cong, J. Lei, H. Fu, Q. Huang, X. Cao, and N. Ling, "HSCS: Hierarchical sparsity based co-saliency detection for RGBD images," *IEEE Trans. Multimedia*, vol. 21, no. 7, pp. 1660–1671, 2019. **1**
- [30] R. Cong, J. Lei, H. Fu, J. Hou, Q. Huang, and S. Kwong, "Going from RGB to RGBD saliency: A depth-guided transformation model," *IEEE Trans. Cybern.*, vol. 50, no. 8, pp. 3627–3639, 2020. **1**
- [31] C. Li, R. Cong, S. Kwong, J. Hou, H. Fu, G. Zhu, D. Zhang, and Q. Huang, "ASIF-Net: Attention steered interweave fusion network for RGB-D salient object detection," *IEEE Trans. Cybern.*, vol. PP, no. 99, pp. 1–13, 2020. **1**
- [32] C. Li, R. Cong, Y. Piao, Q. Xu, and C. C. Loy, "RGB-D salient object detection with cross-modality modulation and selection," in *Proc. ECCV*, 2020, pp. 1–17. **1**
- [33] J.-X. Zhao, Y. Cao, D.-P. Fan, M.-M. Cheng, X.-Y. Li, and L. Zhang, "Contrast prior and fluid pyramid integration for RGBD salient object detection," in *CVPR*, 2019, pp. 3927–3936. **2, 3, 7, 8**
- [34] N. Otsu, "A threshold selection method from gray-level histograms," *IEEE Trans. Syst. Man Cybern.*, vol. 9, no. 1, pp. 62–66, 1979. **2, 3**
- [35] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik, "Learning rich features from RGB-D images for object detection and segmentation," in *ECCV*, 2014, pp. 345–360. **2, 6**
- [36] P. Krähenbühl and V. Koltun, "Efficient inference in fully connected crfs with gaussian edge potentials," in *NeurIPS*, 2011, pp. 109–117. **2, 6**
- [37] X. Li, H. Lu, L. Zhang, X. Ruan, and M.-H. Yang, "Saliency detection via dense and sparse reconstruction," in *ICCV*, 2013, pp. 2976–2983. **2**
- [38] W. Zhu, S. Liang, Y. Wei, and J. Sun, "Saliency optimization from robust background detection," in *CVPR*, 2014, pp. 2814–2821. **2**
- [39] L. Zhou, Z. Yang, Q. Yuan, Z. Zhou, and D. Hu, "Salient region detection via integrating diffusion-based compactness and local contrast," *IEEE Trans. Image Process.*, vol. 24, no. 11, pp. 3308–3320, 2015. **2**
- [40] H. Peng, B. Li, H. Ling, W. Hua, W. Xiong, and S. Maybank, "Salient object detection via structured matrix decomposition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 4, pp. 818–832, 2017. **2**
- [41] H. Peng, B. Li, W. Xiong, W. Hu, and R. Ji, "RGBD salient object detection: A benchmark and algorithms," in *ECCV*, 2014, pp. 92–109. **2, 6**
- [42] H. Song, Z. Liu, H. Du, G. Sun, O. Le Meur, and T. Ren, "Depth-aware salient object detection and segmentation via multiscale discriminative saliency fusion and bootstrap learning," *IEEE Trans. Image Process.*, vol. 26, no. 9, pp. 4204–4216, 2017. **2**
- [43] D. Feng, N. Barnes, S. You, and C. McCarthy, "Local background enclosure for RGB-D salient object detection," in *CVPR*, 2016, pp. 2343–2350. **2**
- [44] R. Ju, L. Ge, W. Geng, T. Ren, and G. Wu, "Depth saliency based on anisotropic center-surround difference," in *ICIP*, 2014, pp. 1115–1119. **2, 6**
- [45] X. Fan, Z. Liu, and G. Sun, "Salient region detection for stereoscopic images," in *DSP*, 2014, pp. 454–458. **2**
- [46] Y. Cheng, H. Fu, X. Wei, J. Xiao, and X. Cao, "Depth enhanced saliency detection method," in *ICIMCS*, 2014, p. 23. **3, 6**
- [47] L. Qu, S. He, J. Zhang, J. Tian, Y. Tang, and Q. Yang, "RGBD salient object detection via deep fusion," *IEEE Trans. Image Process.*, vol. 26, no. 5, pp. 2274–2285, 2017. **3**
- [48] C. Zhu, X. Cai, K. Huang, T. H. Li, and G. Li, "PDNet: Prior-model guided depth-enhanced network for salient object detection," in *ICME*, 2019, pp. 199–204. **3, 7, 8**
- [49] H. Chen and Y. Li, "Progressively complementarity-aware fusion network for RGB-D salient object detection," in *CVPR*, 2018, pp. 3051–3060. **3, 6, 7**
- [50] Y. Piao, W. Ji, J. Li, M. Zhang, and H. Lu, "Depth-induced multi-scale recurrent attention network for saliency detection," in *ICCV*, 2019, pp. 7253–7262. **3, 6, 7, 8, 11**
- [51] D.-P. Fan, Z. Lin, Z. Zhang, M. Zhu, and M.-M. Cheng, "Rethinking RGB-D salient object detection: Models, data sets, and large-scale benchmarks," *IEEE Trans. Neural Netw. Learn. Syst.*, 2020. **3, 6, 7**
- [52] J. Han, H. Chen, N. Liu, C. Yan, and X. Li, "CNNs-based RGB-D saliency detection via cross-view transfer and multiview fusion," *IEEE Trans. Cybern.*, vol. 48, no. 11, pp. 3171–3183, 2017. **3, 7**
- [53] R. Achanta, S. Hemami, F. Estrada, and S. Süsstrunk, "Frequency-tuned salient region detection," in *CVPR*, 2009, pp. 1597–1604. **3, 6**
- [54] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *NeurIPS*, 2017, pp. 5998–6008. **4**
- [55] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *CVPR*, 2018, pp. 7794–7803. **4**
- [56] R. Girshick, "Fast R-CNN," in *ICCV*, 2015, pp. 1440–1448. **6**
- [57] Y. Niu, Y. Geng, X. Li, and F. Liu, "Leveraging stereopsis for saliency analysis," in *CVPR*, 2012, pp. 454–461. **6**
- [58] N. Li, J. Ye, Y. Ji, H. Ling, and J. Yu, "Saliency detection on light field," in *CVPR*, 2014, pp. 2806–2813. **6**
- [59] C. Zhu and G. Li, "A three-pathway psychobiological framework of salient object detection using stereoscopic technology," in *ICCV*, 2017, pp. 3008–3014. **6**
- [60] D.-P. Fan, M.-M. Cheng, Y. Liu, T. Li, and A. Borji, "Structure-measure: A new way to evaluate foreground maps," in *ICCV*, 2017, pp. 4548–4557. **6**
- [61] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778. **6**
- [62] N. Wang and X. Gong, "Adaptive fusion for RGB-D salient object detection," *IEEE Access*, vol. 7, pp. 55277–55284, 2019. **7, 8**
- [63] H. Chen and Y. Li, "Three-stream attention-aware network for RGB-D salient object detection," *IEEE Trans. Image Process.*, vol. 28, no. 6, pp. 2825–2835, 2019. **7**



- [64] H. Chen, Y. Li, and D. Su, "Multi-modal fusion network with multi-scale multi-path and cross-modal interactions for RGB-D salient object detection," *Pattern Recognition*, vol. 86, pp. 376–385, 2019. 7
- [65] R. Shigematsu, D. Feng, S. You, and N. Barnes, "Learning RGB-D salient object detection using background enclosure, depth contrast, and top-down features," in *ICCV*, 2017, pp. 2749–2757. 7