

Bilateral Attention Network for RGB-D Salient Object Detection

Zhao Zhang, Zheng Lin, Jun Xu, Wenda Jin, Shao-Ping Lu, and Deng-Ping Fan

Abstract—Most existing RGB-D salient object detection (SOD) methods focus on the foreground region when utilizing the depth images. However, the background also provides important information in traditional SOD methods for promising performance. To better explore salient information in both foreground and background regions, this paper proposes a Bilateral Attention Network (BiANet) for the RGB-D SOD task. Specifically, we introduce a Bilateral Attention Module (BAM) with a complementary attention mechanism: foreground-first (FF) attention and background-first (BF) attention. The FF attention focuses on the foreground region with a gradual refinement style, while the BF one recovers potentially useful salient information in the background region. Benefitted from the proposed BAM module, our BiANet can capture more meaningful foreground and background cues, and shift more attention to refining the uncertain details between foreground and background regions. Additionally, we extend our BAM by leveraging the multi-scale techniques for better SOD performance. Extensive experiments on six benchmark datasets demonstrate that our BiANet outperforms other state-of-the-art RGB-D SOD methods in terms of objective metrics and subjective visual comparison. Our BiANet can run up to 80fps on 224×224 RGB-D images, with an NVIDIA GeForce RTX 2080Ti GPU. Comprehensive ablation studies also validate our contributions.

Index Terms—Bilateral attention, salient object detection, RGB-D image.

I. INTRODUCTION

SALIENT object detection (SOD) aims to segment the most attractive objects in an image. As an fundamental computer vision task, SOD has been widely applied into many vision applications, such as visual tracking [28], [32], image segmentation [20], [23], [43], and video analysis [53], [47], etc. Most of existing SOD methods [21], [31], [51] mainly deal with RGB images. However, they usually produce inaccurate SOD results on the scenarios of similar texture, complex background, or homogeneous objects [46], [52]. With the popularity of depth sensors in smartphones, the depth information, e.g., 3D layout and spatial cues, is crucial for reducing the ambiguity in the RGB images, and serves as important supplements to improve the SOD performance [25].

Recently, RGB-D SOD has received increasing research attention [4], [37]. Early RGB-D SOD works [35], [40], [42] introduced the depth contrast as an important prior for the SOD task. The recent work of CFPF [55] utilized the depth contrast prior to design an effectiveness loss. These methods essentially explore depth information to shift more priority on the

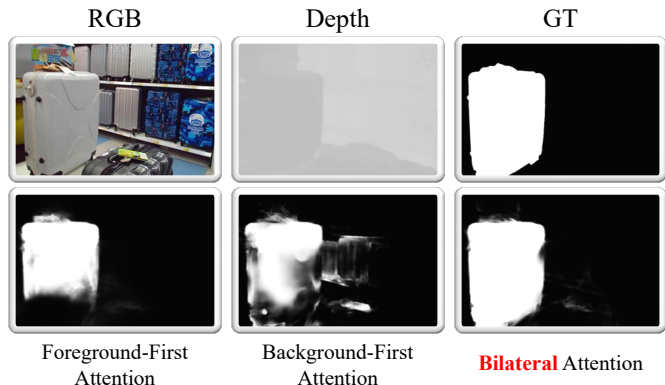


Fig. 1. Comparison of RGB-D SOD results by Foreground-First, Background-First, and our Bilateral attention mechanisms. Depth information provides rich foreground and background relationships. Paying more attention to foreground helps to predict high-confidence foreground objects, but may produce incomplete results. Focusing more on background finds more complete objects, but may introduce unexpected noise. Our BiANet jointly explores foreground and background cues, and achieves complete foreground prediction with little background noise.

foreground region [3], [2]. However, as demonstrated in [29], [48], [49], understanding what background is can also promote the SOD performance. Several traditional methods [26], [50] predict salient objects jointly from the complementary foreground and background information, which is largely ignored by current RGB-D SOD networks.

In this paper, we propose a Bilateral Attention Network (BiANet) to collaboratively learn complementary foreground and background features from both RGB and depth streams for better RGB-D SOD performance. As shown in Figure 2, our BiANet employs a two-stream architecture, and the side outputs from the RGB and depth streams are concatenated in multiple stages. Firstly, we use the high-level semantic features F_6 to locate the foreground and background regions S_6 . However, the initial saliency map S_6 is coarse and in low-resolution. To enhance the coarse saliency map, we design a Bilateral Attention Module (BAM), which is composed of the complementary foreground-first (FF) attention and background-first (BF) attention mechanisms. The FF shifts attention on the foreground region to gradually refine its saliency prediction, while the BF focuses on the background region to recover the potential salient regions around the boundaries. By bilaterally exploring the foreground and background cues, the model helps predict more accurately as shown in Figure 1. Secondly, we propose a multi-scale extension of BAM (MBAM) to effectively learn multi-scale contextual information, and capture both local and global saliency information to further improve the SOD performance. Extensive experiments on six

Zhao Zhang (e-mail: z Zhang@mail.nankai.edu.cn), Zheng Lin, Jun Xu, Shao-Ping Lu, and Deng-Ping Fan are with the TKLNDST, College of Computer Science, Nankai University. Wenda Jin is with Tianjin University. Shao-Ping Lu is the corresponding author (e-mail: sl Lu@nankai.edu.cn).

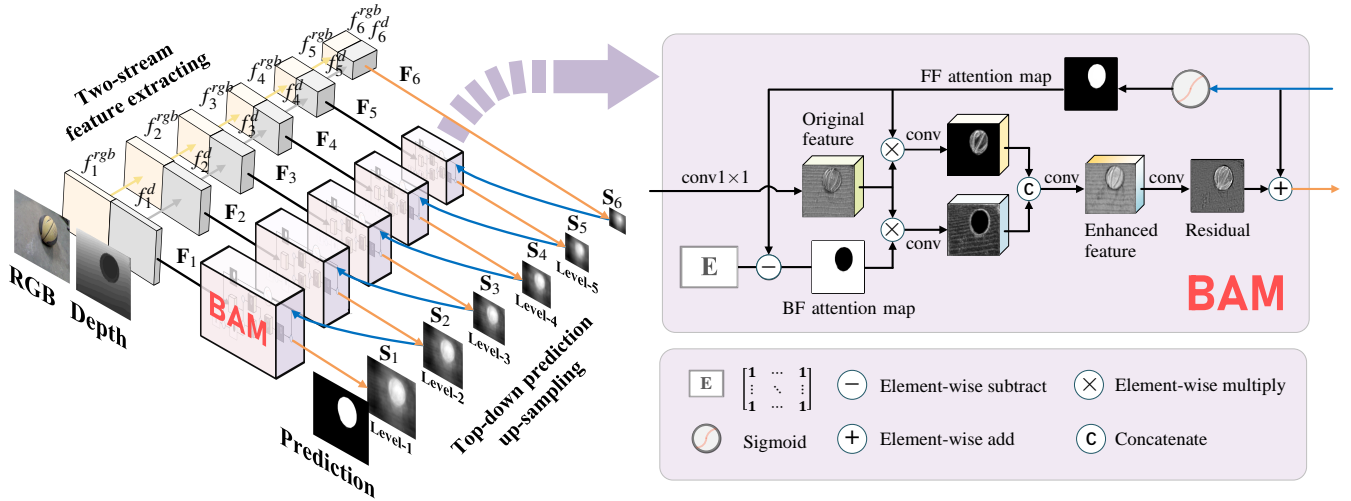


Fig. 2. **The overall architecture of our BiANet.** BAM denotes the proposed Bilateral Attention Module, and it also can be selectively replaced by its multi-scale extension (MBAM). BiANet contains three main steps: two-stream feature extracting, top-down prediction up-sampling, and bilateral attention residual compensation (by BAM). Specifically, it first extracts the multi-level features $\{f_i^{rgb}, f_i^d\}_{i=1}^6$ from the RGB and depth streams, and concatenates them to $\{F_i\}_{i=1}^6$. We take the top feature F_6 to predicate a coarse salient map S_6 . To obtain the accurate and high-resolution result, we up-sample the initial salient map and compensate the details by BAMs in a top-down manner. BAMs receive the higher-level prediction S_{i+1} and current level feature F_i as inputs. In a BAM, the foreground-first attention map A_i^F and the background-first attention map A_i^B can be calculated according to S_{i+1} . We apply the dual complementary attention maps to explore the foreground and background cues bilaterally, and jointly infer the residual for refining the up-sampled saliency map.

benchmark datasets demonstrate that our BiANet achieves better performance than previous state-of-the-arts on RGB-D SOD, and is very fast owing to our simple architecture.

In summary, our main contributions are three-fold:

- **We propose a simple yet effective Bilateral Attention Module (BAM)** to explore the foreground and background cues collaboratively with the rich foreground and background information from the depth images.
- **Our BiANet achieves better performance on six popular RGB-D SOD datasets** under nine standard metrics, and presents better visual effects (*e.g.*, contains more details and sharp edges) than the state-of-the-art methods.
- **Our BiANet runs at 34fps~80fps** on an NVIDIA GeForce RTX2080Ti GPU under different settings, and is a feasible solution for real-world applications.

The remainder of this paper is organized as follows. In §II, we briefly survey the related work. In §III, we present the proposed Bilateral Attention Network (BiANet) for RGB-D Salient Object Detection. Extensive experiments are conducted in §IV to evaluate its performance when compared with state-of-the-art RGB-D SOD methods on six benchmark datasets. The conclusion is given in §V.

II. RELATED WORK

A. RGB-D Salient Object Detection

RGB-D salient object detection (SOD) aims to segment the most attractive object(s) in a pair of cross-modal RGB and depth images. Early methods mainly focus on extracting low-level saliency cues from RGB and depth images, exploring object distance [25], difference of Gaussian [22], graph knowledge [9], multi-level discriminative saliency fusion [42], multi-contextual contrast [8], [35], and background enclosure [13], *etc.* However, these methods often produce inaccurate saliency predictions, due to the lack of high-level feature representation.

Recently, deep neural networks (DNNs) [18] have been employed to investigate high-level representations of cross-modal fusion of RGB and depth images, with much better SOD performance. Most of these DNNs [5], [17], [44] first extract the RGB and depth features separately and then fuse them in the shallow, middle, or deep layers of the network. The methods of [3], [4], [27], [37] further improved the SOD performance by fusing cross-modal features in multi-level stages instead of as a one-off integration. Zhao *et al.* [55] also took the enhanced depth image as attention maps to boost RGB features in multiple stages with better SOD performance.

B. Foreground and Background Cues

There are great differences in the distribution of foreground and background, so it is necessary to explore their respective cues. In traditional methods, some works focus on reasoning salient areas in foreground and background jointly. Yang *et al.* [50] proposed a two-stage method for SOD. It first takes the four boundaries of the inputs as background seeds to infer foreground queries via a graph-based manifold ranking. Then, it ranks the graph depending on the foreground seeds in the same manner for final detection. This method is enlightening, but it has obvious limitations: 1) It is inappropriate to use the four boundaries directly as background, because the foreground is likely to be connected to the boundaries. 2) Aggregation at the super-pixel level also results in rough outputs. For the limitation 1), Liang *et al.* [29] introduce the depth map to distinguish foreground and background regions instead of only assuming the boundaries as background. The depth map shows clear disparity in most senses; thus, it can support more precise locating. For the limitation 2), Li *et al.* [26] further used the regularized random walks ranking to formulate pixel-wised saliency maps, which improves the scaling effect caused by super-pixel aggregation. Nevertheless, only depending on these

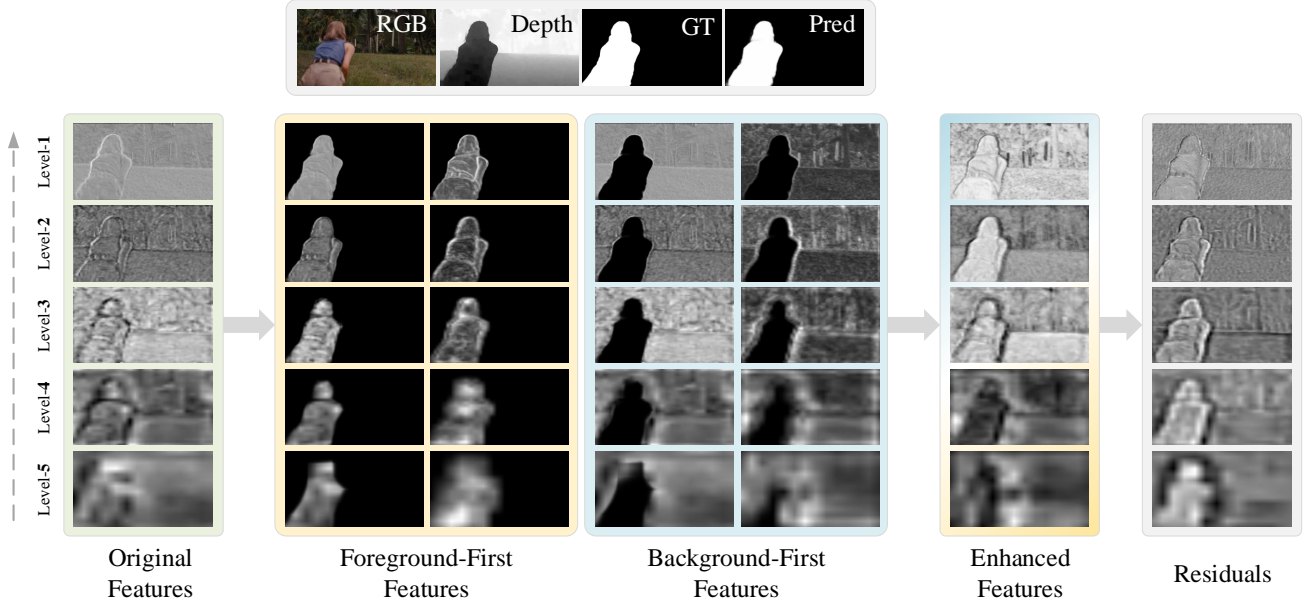


Fig. 3. **Visualizing the working mechanism of bilateral attention.** The original features are the averaged side-output features in each levels. We show the original features directly multiplied by foreground- and background-first attention maps in left columns of yellow and blue boxes. The right columns of the two boxes are the further convoluted features in two branches. As can be seen, the foreground-first features focus on foreground region to explore the saliency cues; while the background-first features shift more attention to the background regions to mine the potentially significant objects. No matter in the features of foreground- or background-first features, more priority is shifted to the uncertain areas caused by the up sampling. When fusing the two branches and jointly inferring, we can see the bilaterally enhanced features have a more accurate understanding where the foreground or background is. Due to obtaining more attention, the uncertain areas are reassigned to the right attribution by the residual with strong contrast. 'Pred' is the prediction of the model.

low-level priors, traditional methods cannot accurately locate the initial region of foreground and background.

Recently, Chen *et al.* [6] proposed to gradually explore saliency regions from the background using reverse attention, but they ignored the contribution of foreground cues to the final detection. As far as we know, how to jointly refine the salient objects from the foreground and background regions is still an open problem in deep RGB-D SOD methods.

III. PROPOSED BIANET FOR RGB-D SOD

In this section, we first introduce the overall architecture of our BiANet, and then present the bilateral attention module (BAM) as well as its multi-scaled extension (MBAM).

A. Architecture Overview

As shown in Figure 2, our Bilateral Attention Network (BiANet) contains three main steps: feature extracting, prediction up-sampling, and bilateral attention residual compensation. We extract the multi-level features from the RGB and depth streams. With increasing network depth, the high-level features (e.g., \mathbf{F}_4) will be more potent for capturing global context, while it loses the object details. When we up-sample the high-level predictions, the saliency maps (e.g., \mathbf{S}_5) will be blurred, e.g. the edges will become uncertain. Thus, we use the proposed Bilateral Attention Module (BAM) to distinguish foreground and background regions.

1) *Feature extracting*: We encode RGB and depth information with two streams. Specifically, both the RGB and depth streams employ five convolutional blocks from VGG-16 [41]

as the standard backbone and attach an additional convolution group with three convolutional layers to predict the saliency maps, respectively. Unlike previous works [17], [57], [5], we explore the cross-modal fusion of RGB and depth features at multiple stages, rather than fusing them once in low or high stage. The i -th side output f_i^{rgb} from the RGB stream and f_i^d from the depth stream are concatenated as a feature tensor \mathbf{F}_i . Note that, \mathbf{F}_6 is concatenated by $M(f_5^{rgb})$ and $M(f_5^d)$, where $M(\cdot)$ denotes the max-pooling operation. The coarse saliency map \mathbf{S}_6 is derived from \mathbf{F}_6 , and $\{\mathbf{F}_1, \mathbf{F}_2, \dots, \mathbf{F}_5\}$ are prepared for the BAMs in our BiANet to further refine the up-sampled saliency maps, by distinguishing the uncertain regions as foreground or background in a top-down manner.

2) *Prediction up-sampling*: The initial saliency map predicted from the high-level features is coarse in low-resolution, but useful to predict the initial position of the foreground and background, since it contains rich semantic information. To refine the basic saliency map \mathbf{S}_6 , a lower-level feature \mathbf{F}_5 with more details is used to predict the residual component between the higher-level prediction and the ground-true (GT) with the help of BAM. We add the predicted residual component \mathbf{R}_5 to the up-sampled higher-level prediction \mathbf{S}_6 , and obtain a refined prediction \mathbf{S}_5 , etc., that is,

$$\mathbf{S}_i = \mathbf{R}_i + U(\mathbf{S}_{i+1}), i \in \{1, \dots, 5\}, \quad (1)$$

where $U(\cdot)$ means up-sampling. Finally, our BiANet obtains a saliency map by $\mathbf{S} = \sigma(\mathbf{S}_1)$, where $\sigma(\cdot)$ is a sigmoid function.

3) *Bilateral attention residual compensation*: To get better residuals and distinguish up-sampled foreground and background regions, we design a bilateral attention module (BAM)

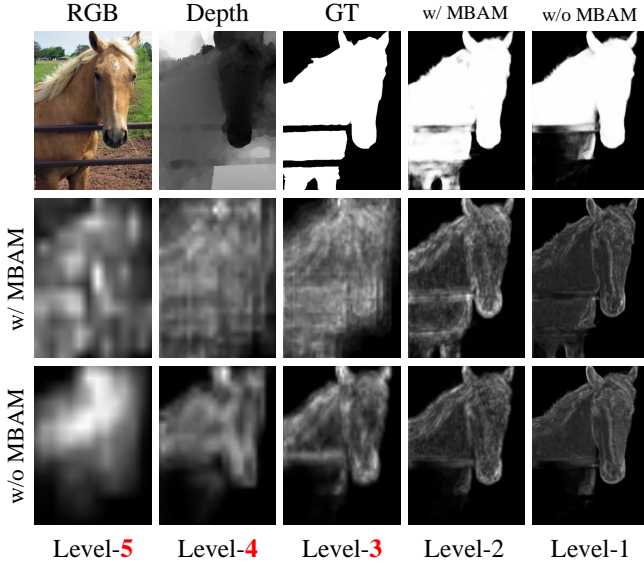


Fig. 4. **Comparison of the high-level features captured by MBAM and BAM.** The second row is the averaged foreground-first features from the model where the MBAMs are applied in the top three levels (marked with red numbers). The third row is the averaged foreground-first features obtained from the model in which all levels are armed with MBAMs. We can see that, compared with applying the MBAMs, MBAMs in higher levels capture more complete information, which is conducive to the object locating as shown in the first row.

to enable our BiANet to discriminate the foreground and background. In our BAM, the higher-level prediction serves as a foreground-first attention (FF) map, and the reversed prediction serves as background-first (BF) attention map to combine the bilateral attention on foreground and background. In Figure 3, one can see that the residual generated by BAM possesses high contrast at the object boundaries. More details are described in Sections III-B and III-C.

4) *Loss function:* Deep supervision is widely used in the SOD task [14], [19]. It clarifies the optimization goals for each step of the network, and accelerates the convergence of training. For quick convergence, we also apply deep supervision in the depth stream output \mathbf{S}_d , RGB stream output \mathbf{S}_{rgb} , and each top-down side output $\{\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_6\}$. The total loss function of our BiANet is

$$\mathcal{L} = \sum_{i=1}^6 w_i \mathcal{L}_{ce}(\sigma(\mathbf{S}_i), \mathbf{GT}) + w_d \mathcal{L}_{ce}(\sigma(\mathbf{S}_d), \mathbf{GT}) + w_{rgb} \mathcal{L}_{ce}(\sigma(\mathbf{S}_{rgb}), \mathbf{GT}), \quad (2)$$

in which w_i, w_d , and w_{rgb} are the weight coefficients and simply set to 1 in our experiments. $\mathcal{L}_{ce}(\cdot)$ is the binary cross entropy loss, which is formulated as

$$\mathcal{L}_{ce}(\mathbf{X}, \mathbf{Y}) = -\frac{1}{N} \sum_{i=1}^N (y_i \log(x_i) + (1 - y_i) \log(1 - x_i)). \quad (3)$$

In the above equation, $x_i \in \mathbf{X}$ and $y_i \in \mathbf{Y}$, and N denotes the total pixel number.

B. Bilateral Attention Module (BAM)

Given the initial foreground and background, how to refine the prediction using higher-resolution cross-modal features is the focus of this paper. Considering that the distribution of

foreground and background are quite different, we design a bilateral attention module using a pair of reversed attention components to learn features from the foreground and background respectively, and then jointly refine the prediction. As can be seen in Figure 2, to focus more on the foreground, we use the up-sampled prediction from the higher-level as foreground-first attention (FF) maps $\{\mathbf{A}^F\}_{i=1}^5$ after they are activated by sigmoid, and the background-first attention (BF) maps $\{\mathbf{A}^B\}_{i=1}^5$ are generated by subtracting FF maps from matrix \mathbf{E} , in which all the elements are 1.

$$\begin{cases} \mathbf{A}_i^F = \sigma(U(\mathbf{S}_{i+1})), \\ \mathbf{A}_i^B = \mathbf{E} - \sigma(U(\mathbf{S}_{i+1})), \end{cases} \quad i \in \{1, 2, 3, 4, 5\}. \quad (4)$$

Then, as shown in Figure 2, we apply FF and BF to weight the side-output features in two branches, respectively, and further predict the residual component jointly.

$$\mathbf{R}_i = \mathcal{P}_R([\mathcal{P}(\hat{\mathbf{F}}_i \odot \mathbf{A}_i^F), \mathcal{P}(\hat{\mathbf{F}}_i \odot \mathbf{A}_i^B)]), \quad (5)$$

where $\hat{\mathbf{F}}_i$ is the channel-reduced feature of \mathbf{F}_i using 32 1×1 convolutions to reduce the computational cost. \mathcal{P} represents the feature extraction operation consisting of 32 convolution kernels with a size of 3×3 and a ReLU layer. The two branches do not share parameters. $[\cdot, \cdot]$ means concatenation. \mathcal{P}_R is the prediction layer to output a single channel residual map via a 3×3 kernel after the same feature extraction operation with \mathcal{P} . Once the \mathbf{R}_i is obtained, the refined prediction \mathbf{S}_i is obtained via Equation 1.

To better understand the working mechanism of BAM, in Figure 3, we visualize the channel-wise averaged features from BAMs in different levels. In BAM, the original features will be first fed into two branches by multiply the FF and BF attention maps, respectively. The result of the direct multiplication is illustrated in the left half of the yellow (FF features) and blue (BF features) boxes. We can see that FF branch shifts attention to the foreground area predicted from its higher level to explore foreground saliency cues. After a convolution layer, more priority is given to the uncertain area. Complementarily, BF branch focuses on the background area to explore the background cues, looking for possible salient objects within it. In our BiANet, the top-down prediction up-sampling is a process in which the resolution of salient objects is gradually increased. It will result in uncertain coarse edges. We can see that both of FF and BF features focus on the uncertain area (such as object boundaries). The low-level and high-resolution FF branch will eliminate the overflow of the uncertain area, while the BF branch will eliminate the uncertain area which does not belong to the background. That is an important reason why BiANet performs better on detail and is prone to predicting sharp edges. After the joint inferring, we can see the bilaterally enhanced features contain more discriminative spatial information of foreground and background. The generated residual components are with sharp contrast on the edges, and then suppress the background area and strengthen the foreground regions.

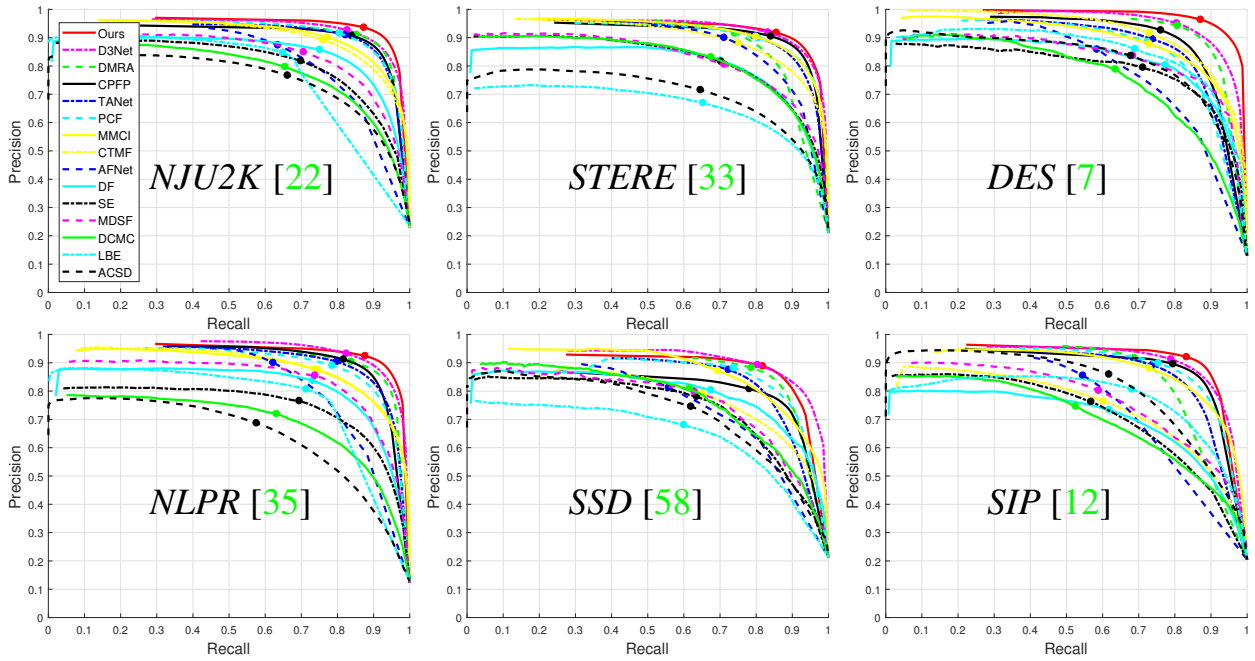


Fig. 5. PR curves of our BiANet and other 14 state-of-the-art methods across 6 datasets. The node on each curve denotes the precision and recall value used for calculating max F-measure.

C. Multi-Scale Extension of BAM (MBAM)

Salient objects in a scene are various in location, size, and shape. Thus, exploring the multi-scaled context in high-level layers benefits for understanding the scene [45], [56]. To this end, we extend our BAM with a multi-scale version, in which groups of dilated convolutions are used to extract pyramid representations from the undetermined foreground and background areas. Specifically, the module can be described as

$$\mathbf{R}_i = \mathcal{P}_R \left(\left[\sqcup_{i=1}^4 \mathcal{D}_i (\mathbf{F}_i \odot \mathbf{A}_i^F), \sqcup_{i=1}^4 \mathcal{D}_i (\mathbf{F}_i \odot \mathbf{A}_i^B) \right] \right), \quad (6)$$

where \sqcup means a concatenate operation. \mathcal{D}_1 consists of 1×1 kernels with 32 channels and a ReLU layer. $\{\mathcal{D}_i\}_{i=2}^4$ is a group of dilated convolutions, with rates of 3, 5, and 7. They all consist of 3×3 kernels with 32 channels and a ReLU layer.

We recommend applying the MBAM in high-level cross-modal features, such as $\{\mathbf{F}_3, \mathbf{F}_4, \mathbf{F}_5\}$, which need different sizes of receptive fields to explore multi-scale context. MBAM effectively improves the detection performance but introduces a certain computational cost. Thus, the number of MBAM should be a trade-off in practical applications. In Section IV-C3, we discuss in detail how the number of MBAM changes the detection effect and calculation cost.

In order to intuitively observe the gain effect brought by MBAM, we visualize the averaged foreground-first feature maps from MBAMs and BAMs in Figure 4. In the second row, the feature maps are obtained from the model with three MBAMs in its top three levels, while in the last row, all the feature maps are collected from BAMs. We can see the target object (horse) account for a large proportion of the scene. Without the ability to perceive multi-scale information effectively, the BAM fails to capture the accurate global salient regions in high levels and leads to incomplete prediction finally. When introducing the multi-scale extension, we can see

higher-level features achieve stronger spatial representation, which supports to locate more complete salient object.

D. Implementation Details

1) *Settings*: We apply the MBAM in the high-level side outputs $\{\mathbf{F}_3, \mathbf{F}_4, \mathbf{F}_5\}$ during implementation, and use bilinear interpolation in all interpolation operations. The initial parameters of our backbone are loaded from a VGG-16 network pre-trained on ImageNet. Our BiANet is based on PyTorch [34].

2) *Training*: Following D3Net [12], we use the training set containing 1485 and 700 image pairs from the NJU2K [22] and NLPR [35] datasets, respectively. We employ the Adam optimizer [24] with an initial learning rate of 0.0001, $\beta_1 = 0.9$, and $\beta_2 = 0.99$. The batch size is set to 8, and we train our BiANet for 25 epochs in total. The training images are resized to 224×224 , also during the test. The output saliency maps are resized back to the original size for evaluation. Accelerated by an NVIDIA GeForce RTX 2080Ti, our BiANet takes about 2 hours for training, and runs at 34~80fps (with different numbers of MBAMs) for the inputs with 224×224 resolution.

IV. EXPERIMENTS

A. Evaluation Protocols

1) *Evaluation datasets*: We conduct experiments on six widely used RGB-D based SOD datasets. NJU2K [22] and NLPR [35] are two popular large-scale RGB-D SOD datasets containing 1985 and 1000 images, respectively. DES [7] contains 135 indoor images with fine structures collected with Microsoft Kinect [54]. STERE [33] contains 1000 internet images, and the corresponding depth maps are generated by stereo images using a sift flow algorithm [30]. SSD [58] is

TABLE I

QUANTITATIVE COMPARISONS OF OUR BIANET WITH NINE DEEP-LEARNING-BASED METHODS AND FIVE TRADITIONAL METHODS ON SIX POPULAR DATASETS IN TERM OF S-MEASURE (S_α), MAXIMUM F-MEASURE (MAX F_β), MEAN F-MEASURE (MEAN F_β), ADAPTIVE F-MEASURE (ADP F_β), MAXIMUM E-MEASURE (MAX E_ξ), MEAN E-MEASURE (MEAN E_ξ), ADAPTIVE E-MEASURE (ADP E_ξ), AND MEAN ABSOLUTE ERROR (MAE, \mathcal{M}). F_β AND E_ξ REPRESENT MAX F_β AND MAX E_ξ BY DEFAULT. \uparrow MEANS THAT THE LARGER THE NUMERICAL VALUE, THE BETTER THE MODEL, WHILE \downarrow MEANS THE OPPOSITE. FOR TRADITIONAL METHODS, THE STATISTICS ARE BASED ON OVERALL DATASETS RATHER ON THE TEST SET.

Metric	ACSD	LBE	DCMC	MDSF	SE	DF	AFNet	CTMF	MMCI	PCF	TANet	CPFP	DMRA	D3Net	BiANet	
	ICIP14 [22]	CVPR16 [13]	SPL16 [9]	TIP17 [42]	ICME16 [16]	TIP17 [39]	arXiv19 [44]	TOC18 [17]	PR19 [5]	CVPR18 [3]	TIP19 [4]	CVPR19 [55]	ICCV19 [37]	arXiv19 [12]	2020 Ours	
NJU2K [22]	$S_\alpha \uparrow$	0.699	0.695	0.686	0.748	0.664	0.763	0.772	0.849	0.858	0.877	0.878	0.879	0.886	0.893	0.915
	max $F_\beta \uparrow$	0.711	0.748	0.715	0.775	0.748	0.804	0.775	0.845	0.852	0.872	0.874	0.877	0.886	0.887	0.920
	mean $F_\beta \uparrow$	0.512	0.606	0.556	0.628	0.583	0.650	0.764	0.779	0.793	0.840	0.841	0.850	0.873	0.859	0.903
	adp $F_\beta \uparrow$	0.696	0.740	0.717	0.757	0.734	0.784	0.768	0.788	0.812	0.844	0.844	0.837	0.872	0.840	0.892
	max $E_\xi \uparrow$	0.803	0.803	0.799	0.838	0.813	0.864	0.853	0.913	0.915	0.924	0.925	0.926	0.927	0.930	0.948
	mean $E_\xi \uparrow$	0.593	0.655	0.619	0.677	0.624	0.696	0.826	0.846	0.851	0.895	0.895	0.910	0.920	0.910	0.934
	adp $E_\xi \uparrow$	0.786	0.791	0.791	0.812	0.772	0.835	0.846	0.864	0.878	0.896	0.893	0.895	0.908	0.894	0.926
$\mathcal{M} \downarrow$	0.202	0.153	0.172	0.157	0.169	0.141	0.100	0.085	0.079	0.059	0.060	0.053	0.051	0.051	0.039	
STERE [33]	$S_\alpha \uparrow$	0.692	0.660	0.731	0.728	0.708	0.757	0.825	0.848	0.873	0.875	0.871	0.879	0.835	0.889	0.904
	max $F_\beta \uparrow$	0.669	0.633	0.740	0.719	0.755	0.757	0.823	0.831	0.863	0.860	0.861	0.874	0.847	0.878	0.898
	mean $F_\beta \uparrow$	0.478	0.501	0.590	0.527	0.610	0.617	0.806	0.758	0.813	0.818	0.828	0.841	0.837	0.841	0.879
	adp $F_\beta \uparrow$	0.661	0.595	0.742	0.744	0.748	0.742	0.807	0.771	0.829	0.826	0.835	0.830	0.844	0.829	0.873
	max $E_\xi \uparrow$	0.806	0.787	0.819	0.809	0.846	0.847	0.887	0.912	0.927	0.925	0.923	0.925	0.911	0.929	0.942
	mean $E_\xi \uparrow$	0.592	0.601	0.655	0.614	0.665	0.691	0.872	0.841	0.873	0.887	0.893	0.912	0.879	0.906	0.926
	adp $E_\xi \uparrow$	0.793	0.749	0.831	0.830	0.825	0.838	0.886	0.864	0.901	0.897	0.906	0.903	0.900	0.902	0.926
$\mathcal{M} \downarrow$	0.200	0.250	0.148	0.176	0.143	0.141	0.075	0.086	0.068	0.064	0.060	0.051	0.066	0.054	0.043	
DES [7]	$S_\alpha \uparrow$	0.728	0.703	0.707	0.741	0.741	0.752	0.770	0.863	0.848	0.842	0.858	0.872	0.900	0.898	0.931
	max $F_\beta \uparrow$	0.756	0.788	0.666	0.746	0.741	0.766	0.728	0.844	0.822	0.804	0.827	0.846	0.888	0.880	0.926
	mean $F_\beta \uparrow$	0.513	0.576	0.542	0.523	0.617	0.604	0.713	0.756	0.735	0.765	0.790	0.824	0.873	0.851	0.910
	adp $F_\beta \uparrow$	0.717	0.796	0.702	0.744	0.726	0.753	0.730	0.778	0.762	0.782	0.795	0.829	0.866	0.863	0.915
	max $E_\xi \uparrow$	0.850	0.890	0.773	0.851	0.856	0.870	0.881	0.932	0.928	0.893	0.910	0.923	0.943	0.935	0.971
	mean $E_\xi \uparrow$	0.612	0.649	0.632	0.621	0.707	0.684	0.810	0.826	0.825	0.838	0.863	0.889	0.933	0.902	0.948
	adp $E_\xi \uparrow$	0.855	0.911	0.849	0.869	0.852	0.877	0.874	0.911	0.904	0.912	0.919	0.927	0.944	0.946	0.975
$\mathcal{M} \downarrow$	0.169	0.208	0.111	0.122	0.090	0.093	0.068	0.055	0.065	0.049	0.046	0.038	0.030	0.033	0.021	
NLPR [35]	$S_\alpha \uparrow$	0.673	0.762	0.724	0.805	0.756	0.802	0.799	0.860	0.856	0.874	0.886	0.888	0.899	0.905	0.925
	max $F_\beta \uparrow$	0.607	0.745	0.648	0.793	0.713	0.778	0.771	0.825	0.815	0.841	0.863	0.867	0.879	0.885	0.914
	mean $F_\beta \uparrow$	0.429	0.626	0.543	0.649	0.624	0.664	0.755	0.740	0.737	0.802	0.819	0.840	0.864	0.852	0.894
	adp $F_\beta \uparrow$	0.535	0.736	0.614	0.665	0.692	0.744	0.747	0.724	0.730	0.795	0.796	0.823	0.854	0.832	0.881
	max $E_\xi \uparrow$	0.780	0.855	0.793	0.885	0.847	0.880	0.879	0.929	0.913	0.925	0.941	0.932	0.947	0.945	0.961
	mean $E_\xi \uparrow$	0.578	0.719	0.684	0.745	0.742	0.755	0.851	0.840	0.841	0.887	0.902	0.918	0.940	0.923	0.948
	adp $E_\xi \uparrow$	0.742	0.855	0.786	0.812	0.839	0.868	0.884	0.869	0.872	0.916	0.916	0.924	0.941	0.931	0.956
$\mathcal{M} \downarrow$	0.179	0.081	0.117	0.095	0.091	0.085	0.058	0.056	0.059	0.044	0.041	0.036	0.031	0.033	0.024	
SSD [58]	$S_\alpha \uparrow$	0.675	0.621	0.704	0.673	0.675	0.747	0.714	0.776	0.813	0.841	0.839	0.807	0.857	0.865	0.867
	max $F_\beta \uparrow$	0.682	0.619	0.711	0.703	0.710	0.735	0.687	0.729	0.781	0.807	0.810	0.766	0.844	0.846	0.849
	mean $F_\beta \uparrow$	0.469	0.489	0.572	0.470	0.564	0.624	0.672	0.689	0.721	0.777	0.773	0.747	0.828	0.815	0.832
	adp $F_\beta \uparrow$	0.656	0.613	0.679	0.674	0.693	0.724	0.694	0.710	0.748	0.791	0.767	0.726	0.821	0.790	0.821
	max $E_\xi \uparrow$	0.785	0.736	0.786	0.779	0.800	0.828	0.807	0.865	0.882	0.894	0.897	0.852	0.906	0.907	0.916
	mean $E_\xi \uparrow$	0.566	0.574	0.646	0.576	0.631	0.690	0.762	0.796	0.796	0.856	0.861	0.839	0.897	0.886	0.896
	adp $E_\xi \uparrow$	0.765	0.729	0.786	0.772	0.778	0.812	0.803	0.838	0.860	0.886	0.879	0.832	0.892	0.885	0.902
$\mathcal{M} \downarrow$	0.203	0.278	0.169	0.192	0.165	0.142	0.118	0.099	0.082	0.062	0.063	0.082	0.058	0.059	0.050	
SIP [12]	$S_\alpha \uparrow$	0.732	0.727	0.683	0.717	0.628	0.653	0.720	0.716	0.833	0.842	0.835	0.850	0.806	0.864	0.883
	max $F_\beta \uparrow$	0.763	0.751	0.618	0.698	0.661	0.657	0.712	0.694	0.818	0.838	0.830	0.851	0.821	0.861	0.890
	mean $F_\beta \uparrow$	0.542	0.571	0.499	0.568	0.515	0.464	0.702	0.608	0.771	0.814	0.803	0.821	0.811	0.830	0.873
	adp $F_\beta \uparrow$	0.727	0.733	0.645	0.694	0.662	0.673	0.705	0.684	0.795	0.825	0.809	0.819	0.819	0.829	0.875
	max $E_\xi \uparrow$	0.838	0.853	0.743	0.798	0.771	0.759	0.819	0.829	0.897	0.901	0.895	0.903	0.875	0.910	0.925
	mean $E_\xi \uparrow$	0.614	0.651	0.598	0.645	0.592	0.565	0.793	0.705	0.845	0.878	0.870	0.893	0.844	0.893	0.913
	adp $E_\xi \uparrow$	0.827	0.841	0.786	0.805	0.756	0.794	0.815	0.824	0.886	0.899	0.893	0.899	0.863	0.901	0.920
$\mathcal{M} \downarrow$	0.172	0.200	0.186	0.167	0.164	0.185	0.118	0.139	0.086	0.071	0.075	0.064	0.085	0.063	0.052	

a small-scale but high-resolution dataset with 400 images in 960×1080 resolution. *SIP* [12] is a high-quality RGB-D SOD dataset with 929 person images.

2) *Evaluation metrics*: We employ 9 metrics to comprehensively evaluate these methods. **Precision-Recall (PR) curve** [38] shows the precision and recall performances of the predicted saliency map at different binary thresholds. **F-measure** [1] is computed by the weighted harmonic mean of the thresholded precision and recall. We employ maximum F-measure (max F_β), mean F-measure (mean F_β), and adaptive F-measure (adp F_β). **Mean Absolute Error (MAE, \mathcal{M})** [36]

directly estimates the average pixel-wise absolute difference between the prediction and the binary ground-truth map. **S-measure** (S_α) [10] is an advanced metric, which takes the region-aware and object-aware structural similarity into consideration. **E-measure** [11] is the recent proposed Enhanced alignment measure in the binary map evaluation field, which combines local pixel values with the image level mean value in one term, jointly capturing image-level statistics and local pixel matching information. Similar to F_β , we employ the maximum E-measure (max E_ξ), mean E-measure (mean E_ξ), and adaptive E-measure (adp E_ξ).

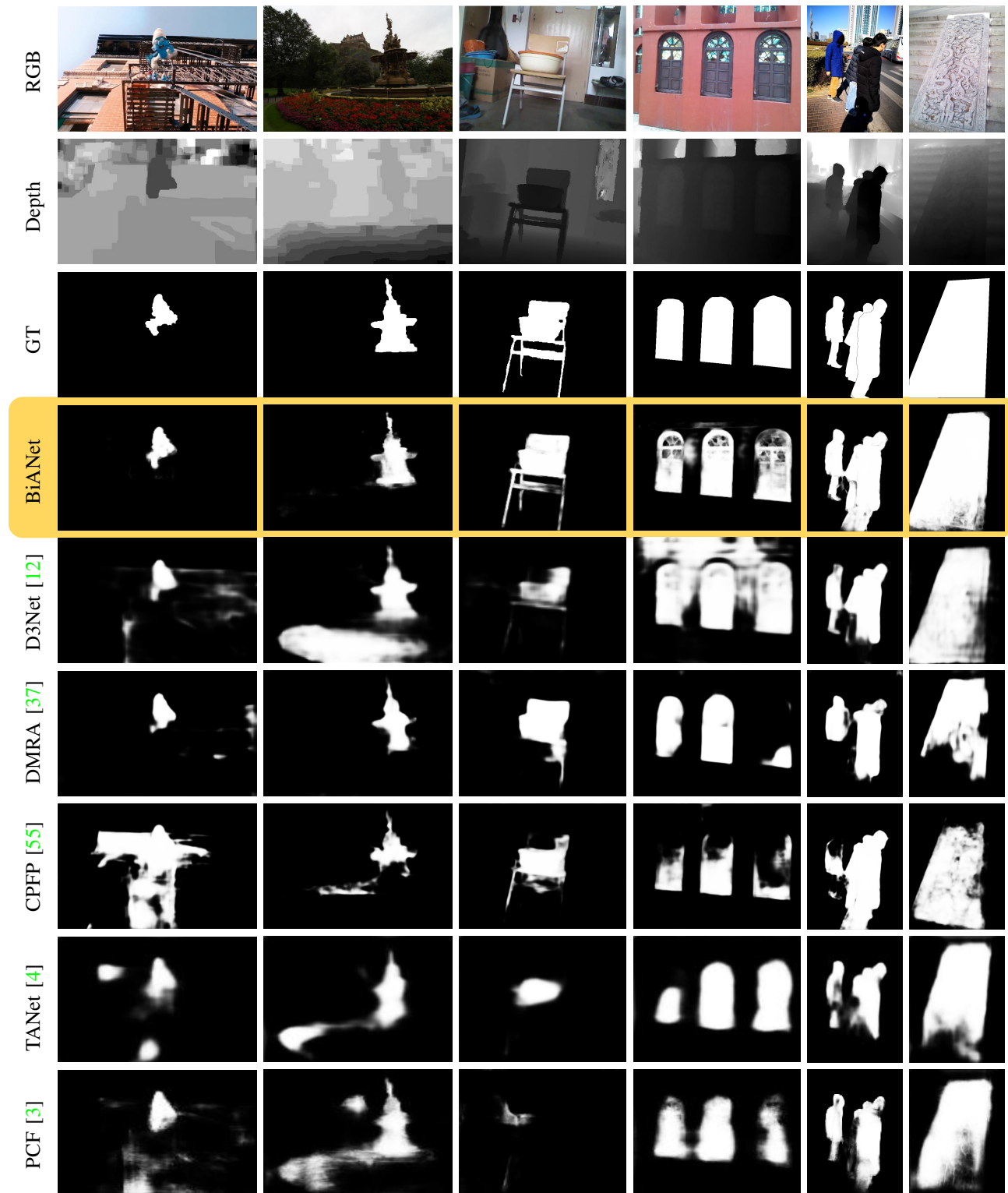


Fig. 6. **Visual comparison of BiANet with other top 5 methods.** The inputs include difficult scenes of tiny objects (column 1), complex background (column 1 and 2), complex texture (column 3), low contrast (column 2 and 6), low-quality or confusing depth (column 2, 4, and 6), and multiple objects (column 4 and 5).

B. Comparison with State-of-the-Arts

1) *Comparison methods:* We compared with 14 state-of-the-art RGB-D SOD methods, including 5 traditional methods: ACSD [22], LBE [13], DCMC [9], MDSF [42], and SE [16],

and 9 DNN-based methods: DF [39], AFNet [44], CTMF [17], MMCI [5], PCF [3], TANet [4], CPFP [55], DMRA [37], and D3Net [12]. The codes and saliency maps of these methods are provided by the authors.

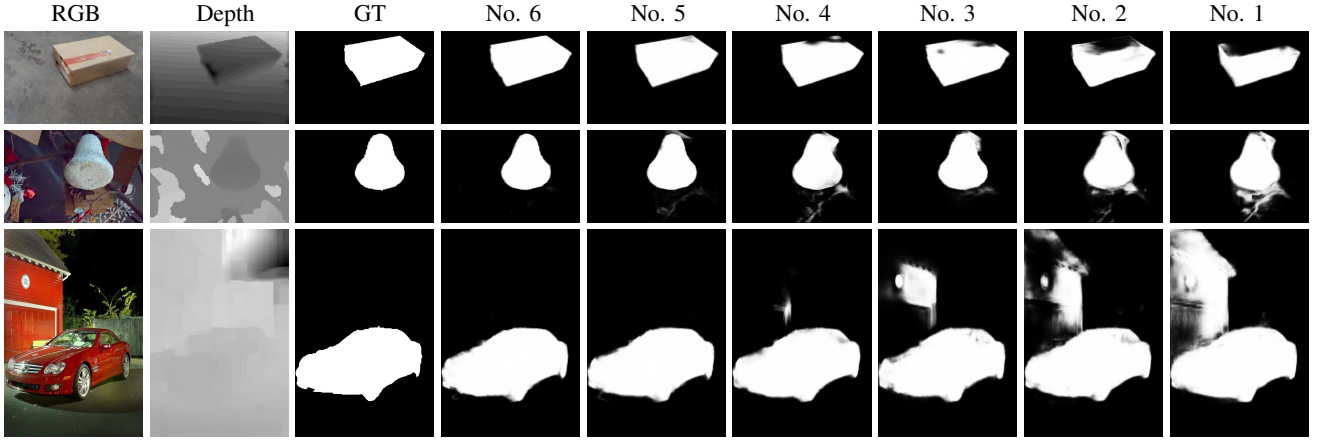


Fig. 7. **Visual comparison in the ablation studies.** The candidate mechanisms are deep information (Dep), foreground-first attention (FF), background-first attention (BF), and multi-scale extension (ME). No. 6: (Dep + FF + BF + ME). No. 5: (Dep + FF + BF). No. 4: (Dep + BF). No. 3: (Dep + FF). No. 2: Dep. No. 1: Baseline.

TABLE II
ABLATION ANALYSIS FOR THE PROPOSED ARCHITECTURE ON THE *NJU2K* AND *STERE* DATASETS. THE CANDIDATE MECHANISMS ARE DEEP INFORMATION (DEP), FOREGROUND-FIRST ATTENTION (FF), BACKGROUND-FIRST ATTENTION (BF), AND MULTI-SCALE EXTENSION (ME). ME IS APPLIED ON THE TOP THREE LEVEL FEATURES.

#	Candidates				<i>NJU2K</i> [22]		<i>STERE</i> [33]	
	Dep	FF	BF	ME	$F_{\beta} \uparrow$	$S_{\alpha} \uparrow$	$F_{\beta} \uparrow$	$S_{\alpha} \uparrow$
No. 1					0.881	0.885	0.882	0.893
No. 2	✓				0.903	0.904	0.887	0.894
No. 3	✓	✓			0.908	0.908	0.895	0.901
No. 4	✓		✓		0.910	0.908	0.892	0.900
No. 5	✓	✓	✓		0.915	0.913	0.897	0.903
No. 6	✓	✓	✓	✓	0.920	0.915	0.898	0.904

TABLE III
IMPROVEMENTS OF ACCURACY BY OUR BAM IN EACH SIDE OUTPUTS COMPARED WITH NO. 2 (WITHOUT BAM & MBA).

BAM	Level-1	Level-2	Level-3	Level-4	Level-5	No. 2
$S_{\alpha} \uparrow$	0.908	0.909	0.908	0.906	0.904	0.904
$F_{\beta} \uparrow$	0.910	0.911	0.909	0.905	0.904	0.903
$E_{\xi} \uparrow$	0.944	0.945	0.943	0.943	0.941	0.942
$\mathcal{M} \downarrow$	0.043	0.043	0.044	0.044	0.045	0.046

TABLE IV
IMPROVEMENTS OF ACCURACY BY OUR MBAM IN EACH SIDE OUTPUTS COMPARED WITH NO. 2 (WITHOUT BAM & MBAM).

MBAM	Level-1	Level-2	Level-3	Level-4	Level-5	No. 2
$S_{\alpha} \uparrow$	0.908	0.909	0.910	0.910	0.910	0.904
$F_{\beta} \uparrow$	0.909	0.912	0.909	0.911	0.911	0.903
$E_{\xi} \uparrow$	0.944	0.945	0.945	0.946	0.947	0.942
$\mathcal{M} \downarrow$	0.044	0.043	0.042	0.042	0.042	0.046

2) *Quantitative evaluation*: The complete quantitative evaluation results are listed in Table I. The comparison methods are presented from right to left according to the comprehensive performance of these metrics, where the lower the value of MAE (\mathcal{M}), the better the effect of the model. The other metrics are the opposite. We also plot the PR curves of these methods in Figure 5. One can see that our BiANet achieves remarkable

advantages over the comparison methods. DMRA [37] and D3Net [12] are well-matched in these datasets. On the large-scaled *NJU2K* [22] and *NLPR* [35] datasets, our BiANet outperforms the second best with $\sim 3\%$ improvement on max F_{β} . On the *DES* [7] dataset, Compared to methods which are heavily dependent on depth information, our proposed BiANet also has a 3.8% improvement on max F_{β} . This indicates that our BiANet can make more efficient use of depth information. Although the *SSD* [58] dataset is high-resolution, the quality of the depth map is poor. Our BiANet still exceeds D3Net [12], which is specifically designed for robustness to low-quality depth maps. Our BiANet also performs the best on the *SIP* [12], which is a challenging dataset with complex scenes and multiple objects.

3) *Qualitative results*: To further demonstrate the effectiveness of our BiANet, we visualized the saliency maps of our BiANet and other top 5 methods in Figure 6. One can see that the target object in the 1st column is tiny, and its white shoes and hat are hard to distinguish from the background. Our BiANet effectively utilizes the depth information, while the others are disturbed by RGB background clutter. The inputs in the 2nd column are challenging because the depth map is mislabeled, and the RGB image was taken in a dark environment with low contrast. Our BiANet successfully detects the target sculpture and eliminates the interference of flowers and the base of the sculpture, while D3Net mistakenly detects a closer rosette, and DMRA loses the part of the object that is similar to the background. The 3rd column shows the ability of our BiANet to detect complex structures of salient objects. Among these methods, only our BiANet completely discover the chairs, including the fine legs. The 4th column is a multi-object scene. Because there are no depth differences between the three salient windows below and the wall, they are not reflected on the depth map, but the three windows above are clearly observed on the depth map. In this case, the depth map will mislead subsequent segmentation. Our BiANet detects multiple objects from RGB images with less noise. The 5th column is also a multi-object scene. The bottom half of depth map is confused with the interference from the ground. Thus, detecting the legs of these persons in the image is very

TABLE V

ACCURACY AND CALCULATION COST ANALYSIS FOR MBAM. $\times 0 \sim \times 5$ MEANS THE NUMBER OF MBAMS, WHICH ARE APPLIED FROM HIGH LEVELS TO LOW LEVELS. FPS DENOTES FRAMES PER SECOND. PARAMS MEANS THE SIZE OF PARAMETERS. FLOPS = FLOATING POINT OPERATIONS. THE ACCURACY METRICS F_β AND \mathcal{M} ARE EVALUATED ON THE *NJU2K* DATASET. THE CALCULATION COST METRICS FPS AND FLOPS ARE TESTED AT 224×224 RESOLUTION. NOTE THAT, $\times 3$ IS THE DEFAULT SETTING IN SECTION IV-B.

	$\times 0$	$\times 1$	$\times 2$	$\times 3$	$\times 4$	$\times 5$	D3Net [12]	DMRA [37]
$F_\beta \uparrow$	0.914	0.917	0.918	0.920	0.920	0.921	0.887	0.886
$\mathcal{M} \downarrow$	0.041	0.040	0.040	0.039	0.038	0.039	0.051	0.051
FPS \uparrow	~ 80	~ 65	~ 55	~ 50	~ 42	~ 34	~ 55	~ 40
Params \downarrow	45.0M	46.9M	48.7M	49.6M	50.1M	50.4M	145.9M	59.7M
FLOPs \downarrow	34.4G	35.0G	36.2G	39.1G	45.2G	58.4G	55.7G	121.0G

TABLE VI

PERFORMANCES OF OUR BiANet BASED ON DIFFERENT BACKBONES. VGG-11 AND VGG-16 IS THE VGG NETWORK PROPOSED IN [41]. RESNET-50 IS PROPOSED IN [18]. RES2NET-50 IS PROPOSED IN [15].

Backbone	VGG-11	VGG-16	ResNet-50	Res2Net-50	
FPS	60	50	25	23	
<i>NJU2K</i> [22]	$S_\alpha \uparrow$	0.912	0.915	0.917	0.923
	$F_\beta \uparrow$	0.913	0.920	0.920	0.925
	$E_\xi \uparrow$	0.947	0.948	0.949	0.952
	$\mathcal{M} \downarrow$	0.040	0.039	0.036	0.034
<i>STERE</i> [33]	$S_\alpha \uparrow$	0.899	0.904	0.905	0.908
	$F_\beta \uparrow$	0.892	0.898	0.899	0.904
	$E_\xi \uparrow$	0.941	0.942	0.943	0.942
	$\mathcal{M} \downarrow$	0.045	0.043	0.040	0.039
<i>DES</i> [7]	$S_\alpha \uparrow$	0.943	0.931	0.930	0.942
	$F_\beta \uparrow$	0.938	0.926	0.927	0.942
	$E_\xi \uparrow$	0.979	0.971	0.968	0.978
	$\mathcal{M} \downarrow$	0.019	0.021	0.021	0.017
<i>NLPR</i> [35]	$S_\alpha \uparrow$	0.927	0.925	0.926	0.929
	$F_\beta \uparrow$	0.914	0.914	0.917	0.919
	$E_\xi \uparrow$	0.951	0.961	0.962	0.963
	$\mathcal{M} \downarrow$	0.024	0.024	0.023	0.023
<i>SSD</i> [58]	$S_\alpha \uparrow$	0.861	0.867	0.863	0.863
	$F_\beta \uparrow$	0.839	0.849	0.843	0.843
	$E_\xi \uparrow$	0.899	0.916	0.911	0.901
	$\mathcal{M} \downarrow$	0.054	0.050	0.048	0.050
<i>SIP</i> [12]	$S_\alpha \uparrow$	0.877	0.883	0.887	0.889
	$F_\beta \uparrow$	0.882	0.890	0.890	0.893
	$E_\xi \uparrow$	0.924	0.925	0.926	0.928
	$\mathcal{M} \downarrow$	0.054	0.052	0.047	0.047

difficult. However, our BiANet successfully detected all the legs. The last row is a large-scale object whose color and depth map are not distinguished. Large scale, low color contrast and lack of discriminative depth information make the scene very challenging. Fortunately, our BiANet is robust on this scene.

C. Ablation Study

In this section, we mainly investigate: 1) the benefits of bilateral attention mechanism to our BiANet; 2) the effectiveness of BAM in different levels to our BiANet for RGB-D SOD; 3) the further improvements of MBAM in different levels to

our BiANet; 4) the benefits of combining BAM and MBAM for RGB-D SOD; and 5) the impact of different backbones to our BiANet for RGB-D SOD.

1) *Effectiveness of bilateral attention*: We conduct ablation studies on the large-scaled *NJU2K* and *STERE* datasets to investigate the contributions of different mechanisms in the proposed method. The baseline model used here contains a VGG-16 backbones and a residual refine structure. It takes RGB images as input without depth information. The performance of our basic network without any additional mechanisms is illustrated in Table II No. 1. Based on the network, we gradually add different mechanisms and test various combinations. These candidates are depth information (Dep), foreground-first attention (FF), background-first attention (BF), and multi-scale extension (ME). In Table II No. 3, by applying FF, the performance is improved to some extent, It benefits from the foreground cues being learned effectively by shifting the attention to the foreground objects. This is also reflected in Figure 7. The foreground objects are detected more accurately; however, without good understanding on background cues, it tend to mistake some background objects, such as the red house in the third row, or cannot find complete foreground objects as lack of exploration on background regions. We get a similar accuracy when using the BF only, as shown in No. 4. It excels at distinguishing between salient areas and non-salient areas in the background, and can help to find more complete regions of the salient object in the uncertain background; however, too much attention focusing on the background and without a good understanding of the foreground cues, it leads that sometimes background noise is introduced. When we combine FF together with BF to form our BAM and apply it in all side outputs, the performance boosts. We can see that BAM increases S-measure by 0.9% and max F-measure by 1.2% compared with No. 2. When we replace the top three levels BAMs with MBAMs, the performance further improved. In Figure 7, compared to the performance of No. 2 without BAM, the detected salient objects of No. 6 possess higher confidence, sharper edges, and less background noise.

2) *Effectiveness of BAM with different levels*: In order to verify that our BAM module is effective at each feature level, we apply BAM to each side output of the No. 2 model’s feature extractor, respectively. That is, in each experiment, BAM is applied to one side output, while the others undergo general convolutions. From Table III, we can see that the BAMs in

every layer facilitate a universal improvement on detection performance. In addition, we find that BAM applied in the lower levels contributes more to the results.

3) *Effectiveness of MBAM in different levels*: In Table II, compared with No. 5, No. 6 carry out multi-scaled extension on its higher three levels $\{F_3, F_4, F_5\}$. This extension effectively improves the performance of the model. In order to better show the gain of MBAM in each level features, similar to Table III, we apply MBAM to each side output of the No. 2 model, respectively. The experimental results are recorded in Table IV, where different levels of MBAM bring different degrees of improvement to the results. Comparing Table III and Table IV, we can see a more interesting phenomenon that the BAM applied in the lower level brings more improvement while the MBAM applied in the higher level is more effective.

4) *Cooperation between BAM and MBAM*: The observation above guides us that when using BAM and MBAM in cooperation, we should give priority to multi-scale expansion of higher-level BAM. Therefore, we expand BAM from top to bottom until all BAMs are converted into MBAMs. We record the final detection performance and calculation cost during the gradual expansion in Table V. We start from the highest level, and gradually increase the number of MBAMs to three. We can see that the effect on the model is a steady improvement, but the computing cost is also increased. At the lower levels, adding MBAM has no obvious effect. This phenomenon is in line with our expectation. Besides, due to the high resolution, the extension of lower-level BAM will increase the calculation cost and reduce the robustness. The selection of the number of MBAM needs to balance the accuracy and speed requirements of the application scenario. In scenarios with higher speed requirements, we recommend not to use MBAM. Our most lightweight model can achieve ~ 80 fps while ensuring significant performance advantages. The parameter size and FLOPs are superior to the SOTA methods D3Net [12] and DMRA [37]. In scenarios where high accuracy is required, we suggest applying less than three MBAMs on higher-level features.

5) *Performances under different backbones*: We implement the BiANet based on some other widely-used backbones to demonstrate the effectiveness of the proposed bilateral attention mechanism on different feature extractors. Specifically, in addition to VGG-16 [41], we provide the results of BiANet on VGG-11 [41], ResNet-50 [18], and Res2Net-50 [15]. Compared with VGG-16, VGG-11 is a lighter backbone. As shown in Table VI, although the accuracy is slightly lower than VGG-16, it still reaches SOTA with a faster speed. BiANet with stronger backbones will bring more remarkable improvements. For example, when we employ ResNet-50 like D3Net [12] as backbone, our BiANet brings 1.5% improvement on *NJU2K* [22] in terms of the MAE compared with the D3Net [12]. When armed with Res2Net-50 [15], BiANet achieves 3.8% improvement on *NJU2K* [22] in terms of the max F-measure compared with the SOTA methods.

D. Failure Case Analysis

In Figure 8, we illustrate some failure cases when our BiANet works in some extreme environments. BiANet ex-

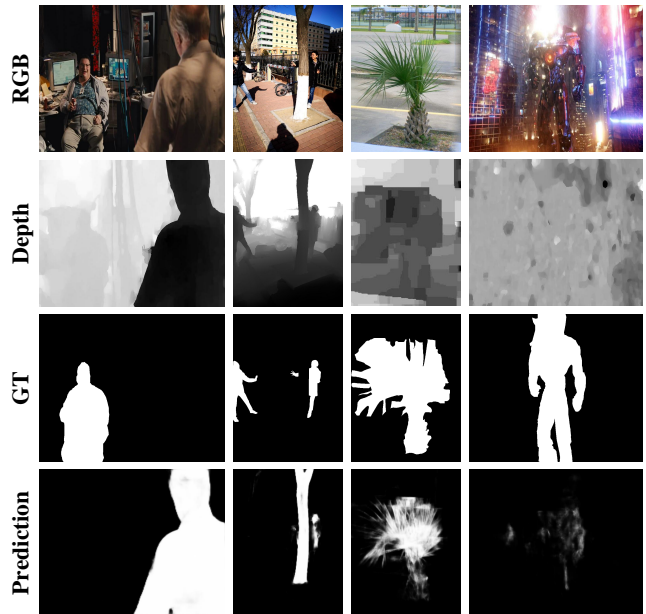


Fig. 8. **Failure cases of BiANet in extreme environments.** In the first two columns, as the objects closer to the observer are not the targets, the depth maps provide misleading information. In the last two columns, the BiANet fails lead by the confusing RGB information and coarse depth maps.

plores the saliency cues bilaterally in the foreground and background regions with the relationship provided by depth information. However, when the foreground regions indicated by depth information do not belong to the salient object, it is likely to mislead the prediction. The first two columns in Figure 8 are typical examples, where our BiANet mistakenly takes the object close to the observer as the target, and gives the wrong prediction. The other situation that may cause failure is when BiANet encounters coarse depth maps in complex scenarios (see the last two columns). In the third column, the depth map provides inaccurate spatial information, which affects the detection of details. In the last column, the inaccurate depth map and the confusing RGB information make BiANet fail to locate the target object.

V. CONCLUSION

In this paper, we propose a fast yet effective bilateral attention network (BiANet) for RGB-D saliency object detection (SOD) task. To better utilize the foreground and background information, we propose a bilateral attention module (BAM) to comprise the dual complementary of foreground-first attention and background-first attention mechanisms. To fully exploit the multi-scale techniques, we extend our BAM module to its multi-scale version (MBAM), capturing better global information. Extensive experiments on six benchmark datasets demonstrated that our BiANet, benefited by our BAM and MBAM modules, outperforms previous state-of-the-art methods on RGB-D SOD, in terms of quantitative and qualitative performance. The proposed BiANet runs at real-time speed on a single GPU, making it a potential solution for various real-world applications.

REFERENCES

- [1] Radhakrishna Achanta, Sheila Hemami, Francisco Estrada, and Sabine Susstrunk. Frequency-tuned salient region detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1597–1604, 2009.
- [2] Chenglizhao Chen, Jipeng Wei, Chong Peng, Weizhong Zhang, and Hong Qin. Improved saliency detection in rgb-d images using two-phase depth estimation and selective deep fusion. *IEEE Transactions on Image Processing*, 29:4296–4307, 2020.
- [3] Hao Chen and Youfu Li. Progressively Complementarity-Aware Fusion Network for RGB-D Salient Object Detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3051–3060, 2018.
- [4] Hao Chen and Youfu Li. Three-stream attention-aware network for RGB-D salient object detection. *IEEE Transactions on Image Processing*, 28(6):2825–2835, 2019.
- [5] Hao Chen, Youfu Li, and Dan Su. Multi-modal fusion network with multi-scale multi-path and cross-modal interactions for RGB-D salient object detection. *PR*, 86:376–385, 2019.
- [6] Shuhan Chen, Xiuli Tan, Ben Wang, and Xuelong Hu. Reverse attention for salient object detection. In *ECCV*, pages 234–250, 2018.
- [7] Yupeng Cheng, Huazhu Fu, Xingxing Wei, Jiangjian Xiao, and Xiaochun Cao. Depth enhanced saliency detection method. In *International Conference on Internet Multimedia Computing and Service*, page 23, 2014.
- [8] Runmin Cong, Jianjun Lei, Huazhu Fu, Junhui Hou, Qingming Huang, and Sam Kwong. Going from rgb to rgb-d saliency: A depth-guided transformation model. *IEEE Transactions on Cybernetics*, pages 1–13, 2019.
- [9] Runmin Cong, Jianjun Lei, Changqing Zhang, Qingming Huang, Xiaochun Cao, and Chunping Hou. Saliency detection for stereoscopic images based on depth confidence analysis and multiple cues fusion. *IEEE Signal Processing Letters*, 23(6):819–823, 2016.
- [10] Deng-Ping Fan, Ming-Ming Cheng, Yun Liu, Tao Li, and Ali Borji. Structure-measure: A new way to evaluate foreground maps. In *IEEE International Conference on Computer Vision (ICCV)*, pages 4548–4557, 2017.
- [11] Deng-Ping Fan, Cheng Gong, Yang Cao, Bo Ren, Ming-Ming Cheng, and Ali Borji. Enhanced-alignment measure for binary foreground map evaluation. In *International Joint Conference on Artificial Intelligence*, pages 698–704, 2018.
- [12] Deng-Ping Fan, Zheng Lin, Jia-Xing Zhao, Yun Liu, Zhao Zhang, Qibin Hou, Menglong Zhu, and Ming-Ming Cheng. Rethinking rgb-d salient object detection: Models, datasets, and large-scale benchmarks. *arXiv preprint arXiv:1907.06781*, 2019.
- [13] David Feng, Nick Barnes, Shaodi You, and Chris McCarthy. Local background enclosure for RGB-D salient object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2343–2350, 2016.
- [14] Mengyang Feng, Huchuan Lu, and Errui Ding. Attentive feedback network for boundary-aware salient object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1623–1632, 2019.
- [15] Shang-Hua Gao, Ming-Ming Cheng, Kai Zhao, Xin-Yu Zhang, Ming-Hsuan Yang, and Philip Torr. Res2net: A new multi-scale backbone architecture. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [16] Jingfan Guo, Tongwei Ren, and Jia Bei. Salient object detection for rgb-d image via saliency evolution. In *IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6, 2016.
- [17] Junwei Han, Hao Chen, Nian Liu, Chenggang Yan, and Xuelong Li. CNNs-based RGB-D saliency detection via cross-view transfer and multiview fusion. *IEEE Transactions on Cybernetics*, 28(6):2825–2835, 2018.
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [19] Qibin Hou, Ming-Ming Cheng, Xiaowei Hu, Ali Borji, Zhuowen Tu, and Philip Torr. Deeply supervised salient object detection with short connections. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(4):815–828, 2019.
- [20] Qibin Hou, Peng-Tao Jiang, Yunchao Wei, and Ming-Ming Cheng. Self-erasing network for integral object attention. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [21] Peng Jiang, Zhiyi Pan, Changhe Tu, Nuno Vasconcelos, Baoquan Chen, and Jingliang Peng. Super diffusion for salient object detection. *IEEE Transactions on Image Processing*, 2019.
- [22] Ran Ju, Ling Ge, Wenjing Geng, Tongwei Ren, and Gangshan Wu. Depth saliency based on anisotropic center-surround difference. In *IEEE Conference on Image Processing (ICIP)*, pages 1115–1119, 2014.
- [23] Chanhong Jung and Changick Kim. A unified spectral-domain approach for saliency detection and its application to automatic object segmentation. *IEEE Transactions on Image Processing*, 21(3):1272–1283, 2011.
- [24] Diederik P. Kingma and Jimmy Lei Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [25] Congyan Lang, Tam V. Nguyen, Harish Katti, Karthik Yadati, Mohan S. Kankanhalli, and Shuicheng Yan. Depth matters: influence of depth cues on visual saliency. In *European Conference on Computer Vision (ECCV)*, pages 101–115, 2012.
- [26] Changyang Li, Yuchen Yuan, Weidong Cai, Yong Xia, and David Dagan Feng. Robust saliency detection via regularized random walks ranking. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2710–2717, 2015.
- [27] Gongyang Li, Zhi Liu, and Haibin Ling. Icnnet: Information conversion network for rgb-d based salient object detection. *IEEE Transactions on Image Processing*, 29:4873–4884, 2020.
- [28] Peixia Li, Boyu Chen, Wanli Ouyang, Dong Wang, Xiaoyun Yang, and Huchuan Lu. GradNet: Gradient-guided network for visual object tracking. In *IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [29] Fangfang Liang, Lijuan Duan, Wei Ma, Yuanhua Qiao, Zhi Cai, and Laiyun Qing. Stereoscopic saliency model using contrast and depth-guided-background prior. *Neurocomputing*, 275:2227–2238, 2018.
- [30] Ce Liu, Jenny Yuen, and Antonio Torralba. Sift flow: Dense correspondence across scenes and its applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(5):978–994, 2011.
- [31] Yi Liu, Jungong Han, Qiang Zhang, and Caifeng Shan. Deep salient object detection with contextual information guidance. *IEEE Transactions on Image Processing*, 29:360–374, 2019.
- [32] V. Mahadevan and N. Vasconcelos. Saliency-based discriminant tracking. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1007–1013, 2009.
- [33] Yuzhen Niu, Yujie Geng, Xueqing Li, and Feng Liu. Leveraging stereopsis for saliency analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 454–461, 2012.
- [34] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 8024–8035, 2019.
- [35] Houwen Peng, Bing Li, Weihua Xiong, Weiming Hu, and Rongrong Ji. RGBD salient object detection: a benchmark and algorithms. In *European Conference on Computer Vision (ECCV)*, pages 92–109, 2014.
- [36] Federico Perazzi, Philipp Krähenbühl, Yael Pritch, and Alexander Hornung. Saliency filters: Contrast based filtering for salient region detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 733–740, 2012.
- [37] Yongri Piao, Wei Ji, Jingjing Li, Miao Zhang, and Huchuan Lu. Depth-induced multi-scale recurrent attention network for saliency detection. In *IEEE International Conference on Computer Vision (ICCV)*, pages 7254–7263, 2019.
- [38] David M W Powers. Evaluation: from Precision, Recall and F-measure to ROC, informedness, markedness and correlation. *Journal of Machine Learning Technologies*, 2(1):37–63, 2011.
- [39] Liangqiong Qu, Shengfeng He, Jiawei Zhang, Jiandong Tian, Yandong Tang, and Qingxiong Yang. RGBD salient object detection via deep fusion. *IEEE Transactions on Image Processing*, 26(5):2274–2285, 2017.
- [40] Jianqiang Ren, Xiaojin Gong, Lu Yu, Wenhui Zhou, and Michael Ying Yang. Exploiting global priors for rgb-d saliency detection. In *IEEE Conference on Computer Vision and Pattern Recognition Workshop*, pages 25–32, 2015.
- [41] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [42] Hangke Song, Zhi Liu, Huan Du, Guangling Sun, Olivier Le Meur, and Tongwei Ren. Depth-aware salient object detection and segmentation via multiscale discriminative saliency fusion and bootstrap learning. *IEEE Transactions on Image Processing*, 26(9):4204–4216, 2017.
- [43] Chung-Chi Tsai, Weizhi Li, Kuang-Jui Hsu, Xiaoning Qian, and Yen-Yu Lin. Image co-saliency detection and co-segmentation via progressive joint optimization. *IEEE Transactions on Image Processing*, 28(1):56–71, 2018.
- [44] Ningning Wang and Xiaojin Gong. Adaptive fusion for RGB-D salient object detection. *CoRR*, abs/1901.01369, 2019.

- [45] Wenguan Wang, Shuyang Zhao, Jianbing Shen, Steven C. H. Hoi, and Ali Borji. Salient object detection with pyramid attention and salient edges. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1448–1457, 2019.
- [46] Yupei Wang, Xin Zhao, Xuecai Hu, Yin Li, and Kaiqi Huang. Focal boundary guided salient object detection. *IEEE Transactions on Image Processing*, 28(6):2813–2824, 2019.
- [47] Ziqin Wang, Jun Xu, Li Liu, Fan Zhu, and Ling Shao. Ranet: Ranking attention network for fast video object segmentation. In *IEEE International Conference on Computer Vision (ICCV)*, Oct 2019.
- [48] Changqun Xia, Jia Li, Xiaowu Chen, Anlin Zheng, and Yu Zhang. What is and what is not a salient object? learning salient object detector by ensembling linear exemplar regressors. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [49] Xiaolin Xiao, Yicong Zhou, and Yue-Jiao Gong. Rgb-dsaliency detection with pseudo depth. *IEEE Transactions on Image Processing*, 28(5):2126–2139, 2018.
- [50] Chuan Yang, Lihe Zhang, Huchuan Lu, Xiang Ruan, and Ming-Hsuan Yang. Saliency detection via graph-based manifold ranking. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3166–3173, 2013.
- [51] Lihe Zhang, Jie Wu, Tiantian Wang, Ali Borji, Guohua Wei, and Huchuan Lu. A multistage refinement network for salient object detection. *IEEE Transactions on Image Processing*, 29:3534–3545, 2020.
- [52] Pingping Zhang, Wei Liu, Huchuan Lu, and Chunhua Shen. Salient object detection with lossless feature reflection and weighted structural loss. *IEEE Transactions on Image Processing*, 28(6):3048–3060, 2019.
- [53] Wei Zhang and Hantao Liu. Study of saliency in objective video quality assessment. *IEEE Transactions on Image Processing*, 26(3):1275–1288, 2017.
- [54] Zhengyou Zhang. Microsoft kinect sensor and its effect. *IEEE Transactions on Multimedia*, 19(2):4–10, 2012.
- [55] Jia-Xing Zhao, Yang Cao, Deng-Ping Fan, Ming-Ming Cheng, Xuan-Yi Li, and Le Zhang. Contrast prior and fluid pyramid integration for RGBD salient object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3927–3936, 2019.
- [56] Ting Zhao and Xiangqian Wu. Pyramid feature attention network for saliency detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3085–3094, 2019.
- [57] Chunbiao Zhu, Xing Cai, Kan Huang, Thomas H Li, and Ge Li. PDNet: Prior-model guided depth-enhanced network for salient object detection. In *ICME*, 2019.
- [58] Chunbiao Zhu and Ge Li. A Three-pathway Psychobiological Framework of Salient Object Detection Using Stereoscopic Technology. In *IEEE International Conference on Computer Vision Workshop*, pages 3008–3014, 2017.