

# CNN-based RGB-D Salient Object Detection: Learn, Select and Fuse

Hao Chen and Youfu Li, *Senior Member, IEEE*

**Abstract**—The goal of this work is to present a systematic solution for RGB-D salient object detection, which addresses the following three aspects with a unified framework: modal-specific representation learning, complementary cue selection and cross-modal complement fusion. To learn discriminative modal-specific features, we propose a hierarchical cross-modal distillation scheme, in which the well-learned source modality provides supervisory signals to facilitate the learning process for the new modality. To better extract the complementary cues, we formulate a residual function to incorporate complements from the paired modality adaptively. Furthermore, a top-down fusion structure is constructed for sufficient cross-modal interactions and cross-level transmissions. The experimental results demonstrate the effectiveness of the proposed cross-modal distillation scheme in zero-shot saliency detection and pre-training on a new modality, as well as the advantages in selecting and fusing cross-modal/cross-level complements.

**Index Terms**—RGB-D, salient object detection, convolutional neural network, cross-modal distillation



## 1 INTRODUCTION

THE availability of depth sensors (e.g., in Microsoft Kinect and Intel RealSense) allows the RGB-based computer vision systems with more accurate and robust performance, hence nurturing a wide range of applications [1]. Complementary to the RGB data, the synchronized depth information carries additional geometry cues, which are immune to appearance changes, illumination varying and subtle background movements. The joint inference with RGB and depth information could benefit various computer vision tasks [2], [3]. A good example is the salient object detection [4] of identifying the most visually attractive object/objects in a scene, which has been widely applied in image retrieval [5] and object tracking [6]. The RGB-based methods are very likely to fail when the salient object and background present similar appearance [7], [8], [9], [10]. From another perspective, the corresponding depth map, which supplies auxiliary saliency cues, opens up a new opportunity to solve this challenge.

By fusing the RGB and depth data, a rich amount of algorithms [12], [13], [14], [15], [16], [17], [18], [19], [20], [21], [22], [23], [24], [25], [26] have been proposed for RGB-D salient object detection. Some previous works [17], [21], [23], [24] focus on crafting RGB-D features with prior knowledge, as the salient object tends to pop-out its surroundings. However, these nontrivial assumptions cannot be well generalized to all contexts. Another line of works [4], [18], [19], [20] infer saliency from each modality separately and then solve the multi-modal fusion problem by straightforward combination schemes. However, the cross-modal complements are not well integrated for better representations. Recently, the success of deep learning techniques [27] in various computer vision tasks motivates more researchers to design RGB-D systems based on deep learning tools.

A popular architecture is the “two-stream” Convolutional Neural Network (CNN) [25], [26], [28], [29], [30], [31], in which the paired RGB and depth images work independently and then aggregate in an early or late stage. In these networks, the depth stream is typically trained from scratch or initialized with the well-trained RGB CNN. Nonetheless, these training schemes typically end with insufficient depth-specific learning due to the scarcity of the labeled depth data.

Without carefully selecting the real complementary cues, the direct combination strategy in previous two-stream networks is also confronted with ambiguous and uninformative fusion. Moreover, with a single fusion layer as done in [25], [26], [28], [31], it is unlikely to explore both the contextual and spatial cross-modal complementarities existed in multiple levels. Thus, the systematic solution for understanding RGB-D data still remains as an open issue. As illustrated in Fig. 1, we argue that an ideal RGB-D fusion system should successfully achieve the following three goals:

(1) **Learn:** In some scenarios, the specialists in one modality (e.g., geometry cues in the depth map) are missing in its counterpart (i.e., the RGB image). Accordingly, an informative RGB-D combination firstly calls for carefully extracting discriminative modal-specific features from each modality. Otherwise, knowledge from one modality may not assist and even mislead the inference for its counterpart. However, we are often confronted with an imbalanced amount of labeled data prepared for each modality. Thus, the challenge lies on how to learn rich modal-specific representations from the new modality with limited labeled data.

(2) **Select:** An informative multi-modal fusion process should be attentive to the real complementary components. This awareness mechanism enables the cross-modal fusion to select complementary representations and ignore the redundant ones.

(3) **Fuse:** The last step is to fuse the selected cross-modal information sufficiently. The complementarities between

---

• Hao Chen and Youfu Li are with the Department of Mechanical Engineering, City University of Hong Kong.  
E-mail: meyfli@cityu.edu.hk (Youfu Li is the corresponding author)

Manuscript received April 19, 2005; revised August 26, 2015.

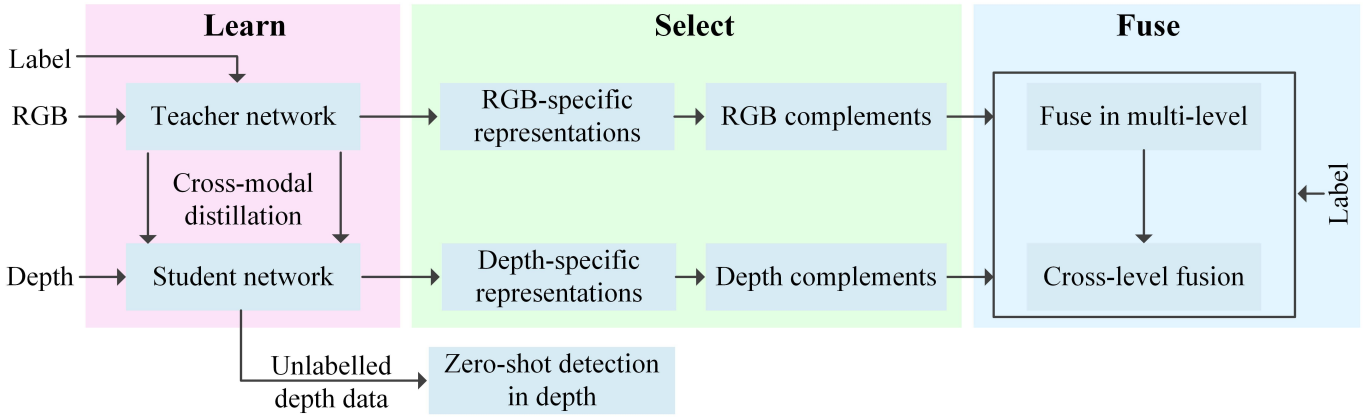


Fig. 1: Our pipeline for RGB-D salient object detection.

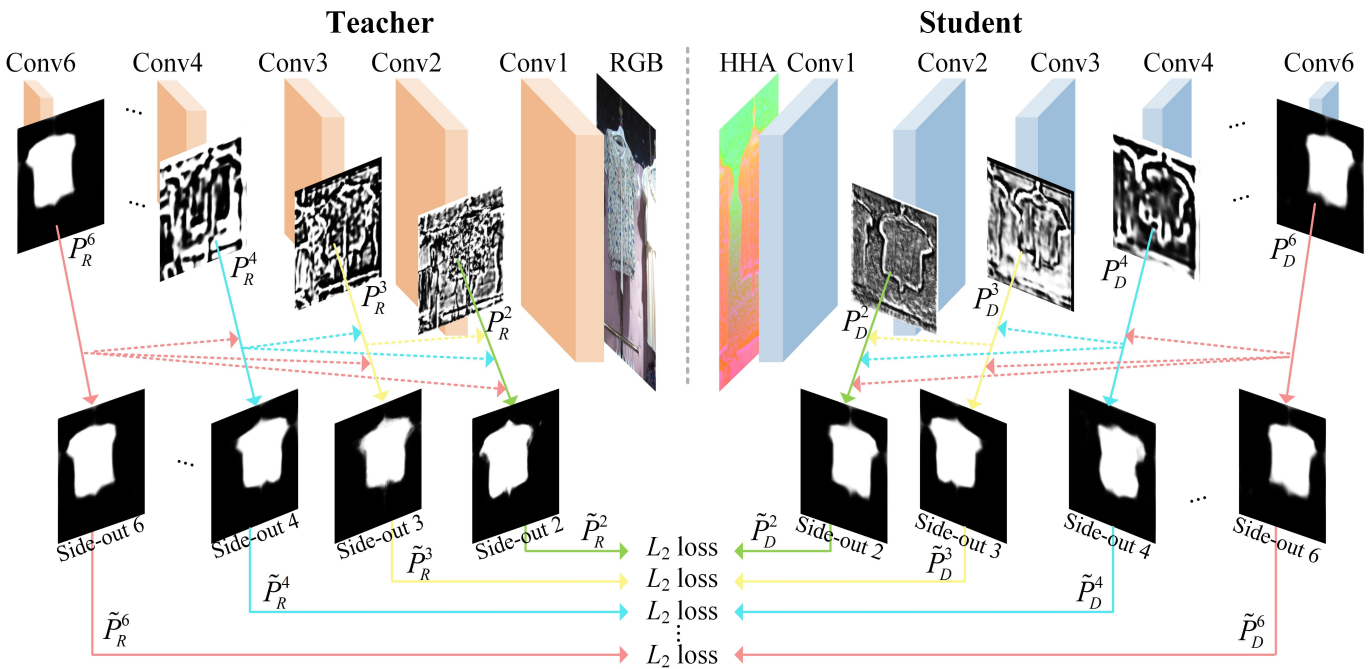


Fig. 2: The architecture of the hierarchical cross-modal distillation network. Adaptation layers in each level are omitted for simplification. When training the teacher network, the L2 loss in each level is replaced with the cross-entropy loss between the side output and the ground-truth mask. For the depth data, we follow the previous approach [11] to encode it as 3-channel HHA (horizontal, aboveground height and surface normal angle) representations.

RGB and depth data exist in both high-level contexts and low-level spatial details. Consequently, a sufficient RGB-D fusion process is in demand to associate both the low-level and high-level cross-modal features for joint decision.

Considering the unavailability of large-scale labeled data in the depth modality, we leverage structured knowledge provided by the source modality (i.e., RGB) to aid the learning of the new modality. Specifically, we use the side outputs of the source modality (i.e., RGB) as supervision to learn the target modality (i.e., depth). We term our scheme *hierarchical cross-modal distillation*, which eschews the reliance on saliency ground-truths in the new modality.

To render an effective fusion process, we explicitly en-

code the cross-modal complements with the residual function and the goal of selecting cross-modal features is formulated as asymptotically approximating the residual. Different from the direct concatenation of multi-modal features, such a cross-modal residual connection is more likely to expose the desired complementarity.

Concerning sufficient multi-modal fusion, we adopt a top-down fusion manner, in which cross-modal features are combined in each level and the integrated RGB-D representations, in turn, guide the inference of shallower layers. The resulting network demonstrates rich multi-level RGB-D representations for joint inference and consequently, the saliency map quality is improved progressively from coarse

to fine.

Our preliminary studies [16] discussed the above-mentioned “Select” and “Fuse” problems. However, the problem “Learn” remains under-studied. In this work, we extend [16] by investigating the problem “Learn” and propose the hierarchical cross-modal distillation method.

In summary, this work has the following five contributions:

(1) We systematically analyze the key issues in interpreting RGB-D data, which guides the system design process.

(2) We propose a cross-modal distillation scheme, which allows zero-shot detection or favors better learning of new modalities with limited labeled data.

(3) The residual function is designed to explicitly capture the cross-modal complementarity.

(4) We propose a progressively top-down cross-modal cross-level fusion topology. Thus, the inference path comes to be aware of modal-specific and level-specific contributions.

(5) This work achieves state-of-the-art performance on three benchmarks consistently.

## 2 RELATED WORK

### 2.1 RGB-D Saliency Detection and Other RGB-D Systems

A large body of earlier works focus on designing RGB-D features or combining unimodal predictions, which are termed as “feature fusion” and “result fusion” solutions respectively. A common wisdom in crafting depth-induced saliency cues is that human fixations prefer closer depth ranges. Based on this prior, Lang et al. [18] use Gaussian Mixture Models to model the distribution of depth-induced saliency. This prior is useful but is easily confused by nearer backgrounds. On the other hand, two regions, sharing the same depth may be in different contexts and should be differentiated. Considering the scene structures, Ju et al. [21] use relative depth instead of the absolute one for evaluation and propose the anisotropic center-surround difference for measurement. Desingh et al. [19] adopt the global-contrast method [9] used in the RGB-induced saliency detection with depth values as inputs. A similar framework is also used in [32]. Different from these global-contrast paradigms, Feng et al. [17] propose to measure the distinction of one region in a local context. They design a local background enclosure feature, which estimates the proportion of the object popping out the background. Peng et al. [4] then propose a hybrid framework that incorporates global-contrast and local-contrast strategies. To further enrich the representative ability of RGB-D data, Song et al. [23] segment the RGB-D pair into superpixels with different scales. The features are then combined as multi-scale representations.

Despite the effectiveness of these handcrafted features, they lack high-level reasoning and suffer from limited generalization ability. To address this limitation, recent works resort to deep learning techniques. Qu et al. [25] combine the low-level features from RGB and depth modalities as the joint input to train a CNN from scratch. Compared to the previous works based on handcrafted features, this method achieves encouraging performance gains. However, it may be difficult to fully leverage the power of CNNs by feeding

the crafted features rather than the raw image pair as inputs. In contrast, Han et al. [26] design a “two-stream” late fusion architecture, in which the RGB and depth images are learned separately and their deep representations are combined by a joint fully connected layer for collaborative decision. Compared to [25], [26] achieves large improvement due to the better combination of high-level contexts. Despite this, the low-level cross-modal complements are underexplored in [26] and the resulting saliency maps are severely blurred. In summary, both [25] and [26] fail to combine the low-level and high-level cross-modal complementarity simultaneously. Recently, Chen et al. [33] propose a multi-branch fusion network with fully connected layers for global reasoning and dilated convolutional layers for capturing local details. The results from two branches are combined by direct summation. However, the network is not a fully convolutional one and fails to utilize the information from all layers for joint inference.

Deep learning techniques especially CNNs are also popular solutions for other RGB-D systems. Among which, the “two-stream” late fusion architecture is likewise the most widely-used one. In [28] and [31], the multi-modal fusion layer combines the decisions from RGB and depth by modeling their consistency and independency. More recent works [29], [34], [35], [36] also investigate the cross-modal complementarity in multiple levels. Although the community of CNN-based RGB-D systems has achieved encouraging improvements, a comprehensive analysis on the RGB-D fusion problem is lacking, which, in our view, will benefit future works on multi-modal systems or new unlabeled modalities a lot.

### 2.2 Cross-modal transfer

The transfer learning community mainly solve the domain adaptation problem in the same modality [37], [38], [39], [40]. In [37], Hinton et al. use the final soft outputs of the large well-trained teacher network as targets of the small student network. Subsequent works [38], [39], [40] extend this idea by encouraging the student to mimic the deep representations from the teacher. Our topic lies in the cross-modal transfer problem, which is more difficult due to the severe cross-modal discrepancy. Notable works include [41], [42], [43], [44], [45], [46]. [42], [43] aims at learning joint representation by mapping the features from different modalities to a shared feature space. [41] learns a mapping from the source modality to the unlabelled new ones to hallucinate modalities, while [44], [45] design a hallucinate network to distill depth features. Gupta et al. [46] generalize the idea in [38], [39], [40] to the cross-modal domain. However, due to the cross-modal discrepancy, a considerable part of source feature maps (e.g., texture and color changes in RGB images) are inaccessible for the target modality (i.e., depth). Hence, it is too strict for the new modality to mimic the high-dimensionality features from the teacher. These RGB-specific features provide uninformative and even negative supervision for the student. As a result, the student network is hard to converge especially when it is deep. In this work, the goal of the student network is relaxed to mimic the side outputs from each level in the teacher network.

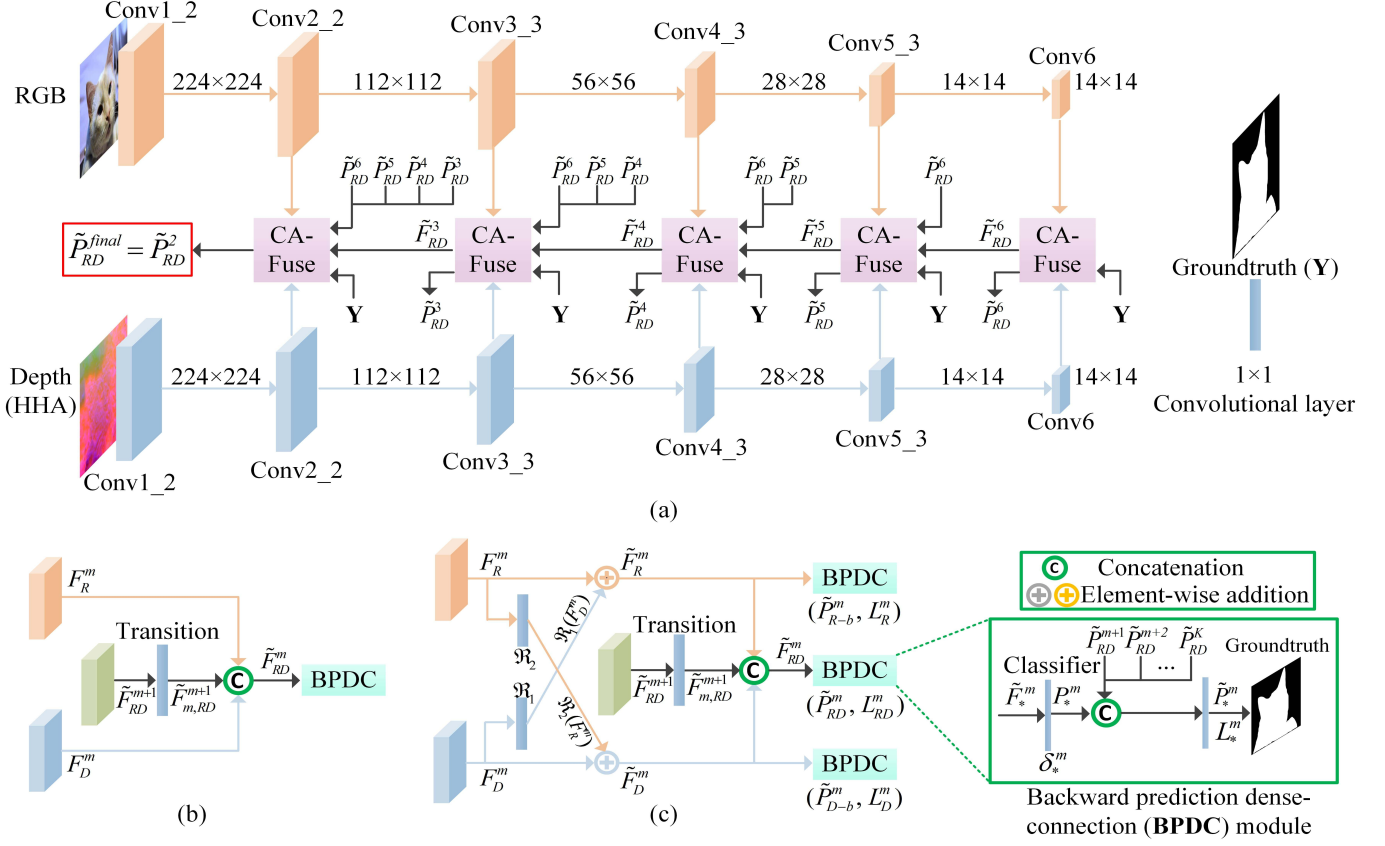


Fig. 3: (a) The architecture of the RGB-D salient object detection network. (b) The details of direct concatenation of cross-modal cross-level features without the cross-modal residual designs. (c) The details of the complementarity-aware fusion (CA-Fuse) block.

### 3 THE PROPOSED METHOD

The proposed model can fulfill zero-shot detection in new unlabeled modalities as well as multi-modal joint inference. We address the three learning objectives with following stages: the teacher network with RGB images and ground-truth masks; the hierarchical cross-modal distillation network with un-annotated RGB-D pairs; and the multi-modal fusion network for RGB-D saliency detection with RGB-D pairs and ground-truths. In the following sections, we will follow the training sequence to introduce each network and discuss how the proposed solution behaves to learn, select and fuse cross-modal complementarity.

#### 3.1 Hierarchical Cross-modal Distillation

For the cross-modal transfer learning, suitable supervisory signals should be customized for the student network. If a excessively strict constraint is set, it turns out to be difficult to ensure training convergence. In contrast, an over-relaxed constraint, such as appropriating the final class distributions, appears too weak to learn the shallow layers effectively. A well-balanced knowledge distillation method should provide sufficient supervisory signals while allows the exploration of specialists in a new modality. Intuitively, the features from two modalities, though discrepant, can make consistent inference for the same task. So our primitive

choice is to use the inference from each level in the teacher network as supervision. However, as the observations in [47], the lower layers are more modal-specific and task-agnostic, while the deeper layers hold opposite characteristics. As a result, the shallow layers will hardly produce coherent inference by different modalities. Specifically, the shallower layers trained with RGB images are activated by texture/color changes, which are immune for the depth modality. We further consider that with the global guidance from the deep layers, the discrepancy between the combined side-outs in shallow layers across different modalities can be effectively reduced with respect to the individual inference in each level. Also, it is hard to optimize multiple levels from scratch jointly in a deep network. However, the progressive enhancement inferred from the teacher reveals level-specific contributions and cross-level collaborations, which are pretty informative supervisory signals for the student. To this end, we let the student progressively approach the side-outs from the teacher. We call it "hierarchical cross-modal distillation scheme". Such a design presents several distinguished advantages:

(a) Compared to the feature-based objective function, the relaxed inference-based supervisory signals allow more flexibility for the student to explore specialists.

(b) These supervision signals decouple the joint learning of multiple layers and define level-specific optimization

objectives for each level independently, which simplifies the training process for a deep student network.

(c) This method promotes a better transfer learning across modalities and even promises the feasibility of using deep CNNs to zero-shot detection on new visual modalities. Given a new modality  $\mathcal{M}_D^X$  with unlabeled training samples  $\mathcal{X}^D$ , our goal lies on learning model-specific features from  $\mathcal{X}^D$  by transferring knowledge from a different modality  $\mathcal{M}_R^X$  with large-scale labeled images.

Denote the  $K$  layered representations  $\Psi = \{\varphi_R^i, i = 1, \dots, K\}$ , where  $\varphi_R^i$  denotes the representation in the  $i^{\text{th}}$  layer for the modality  $\mathcal{M}_R^X$ . Based on  $\varphi_R^i$ , a reliable classifier  $\delta_r^i$  is learned for level-specific inference  $P_R^i = \delta_r^i(\varphi_R^i)$ .

Now, suppose we have a dataset  $\mathcal{D}_{r,d}$  which contains sufficient un-annotated paired images from  $\mathcal{M}_R^X$  and  $\mathcal{M}_D^X$ . We implement this idea by densely skip-connecting the inferences from the deeper layers to all lower layers to generate collaborative side outputs. As shown in Fig.2, the cross-modal distillation network contains two parts: (a) Unimodal cross-level connections (Note that the teacher and student share same cross-level connections), which are described with dotted lines; (b) Cross-modal connections, which are illustrated by solid lines. The inference of a deep layer will be combined with all shallower sideouts (e.g., in the teacher net,  $P_R^6$  will be fed to  $P_{R'}^5, P_{R'}^4, P_{R'}^3, P_{R'}^2$  for combination). Formally,

$$\tilde{P}_R^i = \begin{cases} \mathbf{w}_R^i P_R^i + \sum_{k=i+1}^K \tilde{\mathbf{w}}_{R,i}^k \tilde{P}_R^k, & i = 1, \dots, K-1 \\ P_R^i, & i = K \end{cases} \quad (1)$$

where  $\mathbf{w}_R^i$  and  $\tilde{\mathbf{w}}_{R,i}^k$  are the weights for  $P_R^i$  and the side-out  $\tilde{P}_R^k$  from the  $k^{\text{th}}$  level.

Similarly for the counterpart modality  $\mathcal{M}_D^X$ :

$$\tilde{P}_D^i = \begin{cases} \mathbf{w}_D^i P_D^i + \sum_{k=i+1}^K \tilde{\mathbf{w}}_{D,i}^k \tilde{P}_D^k, & i = 1, \dots, K-1 \\ P_D^i, & i = K \end{cases} \quad (2)$$

where  $P_D^i = \delta_d^i(\phi_D^i)$ ,  $\Phi = \{\phi_D^i, i = 1, \dots, K\}$  denotes the  $K$  layered representations,  $\delta_d^i$  is the learned classifier,  $\mathbf{w}_D^i$  and  $\tilde{\mathbf{w}}_{D,i}^k$  are the weights for  $P_D^i$  and the side-outs  $\tilde{P}_D^k$  from the deeper levels, respectively.

Our scheme for learning sufficient modal-specific representations and inference from images in the modality  $\mathcal{M}_D^X$  is to train the representations  $\Phi$  and inference  $P_D^i(I_d)$  such that the combined side-out  $\tilde{P}_D^i(I_d)$  matches the one  $\tilde{P}_R^i(I_r)$  inferred from its paired image in the modality  $\mathcal{M}_R^X$ . Therefore, we measure the discrepancy between the side-outs from two modalities with a suitable loss  $g$ :

$$L_{HCD} = \sum_{\{I_r, I_d\} \in \mathcal{D}_{r,d}} \sum_{i=1}^K g(\tilde{P}_R^i(I_r), \tilde{P}_D^i(I_d)) \quad (3)$$

In our experiments, we adopt the L2 loss  $g(x, y) = \|x - y\|_2^2$  for measuring the distance. By minimizing  $L_{HCD}$ , the student network is encouraged to learn rich feature hierarchies for inference.

### 3.2 Complementarity-aware Cross-modal Fusion

Having learned modal-specific representations from each modality, the following step is to select the complementary ones for informative multi-modal fusion. To this end,

we propose the complementarity-aware cross-modal fusion (CA-Fuse) block to explicitly select cross-modal complements. Fig. 3(a) shows the architecture of the multi-modal fusion network and Fig. 3(c) exemplifies the CA-Fuse block in the  $m^{\text{th}}$  level. Formally, the adapted deep features from the RGB and depth streams are denoted as  $F_R^m$  and  $F_D^m$ , respectively. A  $1 \times 1$  convolutional layer, acting as a selector, is appended after  $F_D^m$  to select complementary information to enhance the RGB features via a cross-modal skip connection  $\tilde{F}_R^m = F_R^m + \mathfrak{R}_1(F_D^m)$ . It suggests that the target of using  $\mathfrak{R}_1(\cdot)$  to select complementary features from  $F_D^m$  can be posed as approximating the residual part, i.e.,  $\tilde{F}_R^m - F_R^m$  equivalently. Such a reformulation exposes the cross-modal complements explicitly and eases the incorporation. If  $F_R^m$  is competent for inference, the solver can simply adjust the residual towards zero. Otherwise,  $\mathfrak{R}_1(\cdot)$  will be pushed to distill complements from  $F_D^m$  to aid  $F_R^m$  for better prediction. To further encourage the determination of the residual part, the enhanced features  $\tilde{F}_R^m$  will infer saliency  $\tilde{P}_{R-b}^m$  and then compared to the ground-truth  $\mathbf{Y}$ . In minimizing the distance  $L_R^m$  between  $\tilde{P}_{R-b}^m$  and  $\mathbf{Y}$ ,  $\tilde{F}_R^m$  as well as  $\mathfrak{R}_1(F_D^m)$  will be optimized, thereby capturing the most complementary cues from the paired modality. A symmetric residual connection is also introduced from  $F_R^m$  to  $F_D^m$  to capture the complements from the RGB stream to enhance the depth features. Then these features across modalities are concatenated for joint prediction.

### 3.3 Progressively Top-down Cross-modal Cross-level Fusion Pattern

The last question regarding how to fuse the cross-modal complements sufficiently is solved by a progressively top-down fusion pattern, in which the cross-modal features are selected and combined by the CA-Fuse block in each level and the incorporated multi-modal features are selectively transmitted to the adjacent shallower layer for the cross-level combination. Concretely, the RGB-D representations  $\tilde{F}_{m, RD}^{m+1}$ , selected from the  $m+1$  layer  $\tilde{F}_{RD}^{m+1}$  by a transition layer (detailed parameters are illustrated in Table 1), will be upsampled by a fixed de-convolutional layer and then concatenated with  $\tilde{F}_R^m$  and  $\tilde{F}_D^m$  as a cross-level cross-modal representation community  $\tilde{F}_{RD}^m$ , which will be responsible for the prediction of the  $m^{\text{th}}$  CA-Fuse block by:

$$P_{RD}^m = \delta_{rd}^m(\tilde{F}_{RD}^m), \quad (4)$$

where  $\delta_{rd}^m$  are the parameters for fusing cross-modal cross-level features and performing joint inference. Another cross-level fusion strategy of skip-connecting the side-outs densely is also adopted in the multi-modal fusion network and implemented by the backward prediction dense-connection (BPDC) module. The combined side-out is denoted as

$$\tilde{P}_{RD}^m = \begin{cases} \mathbf{w}_{RD}^m P_{RD}^m + \sum_{k=m+1}^K \tilde{\mathbf{w}}_{RD,m}^k \tilde{P}_{RD}^k, & m = 2, \dots, K-1 \\ P_{RD}^m, & m = K \end{cases} \quad (5)$$

where  $\mathbf{w}_{RD}^m$  and  $\tilde{\mathbf{w}}_{RD,m}^k$  denotes the weights for the predictions from the current layer and all deeper layers, respectively. The joint loss function for the multi-modal fusion network consists of the side loss from each CA-Fuse block.

We also add a collaborative loss to encourage a informative combination of all the side-outs:

$$L_{final} = \sum_{m=1}^K \left( d(\tilde{P}_{R-b}^m, \mathbf{Y}) + d(\tilde{P}_{D-b}^m, \mathbf{Y}) + d(\tilde{P}_{RD}^m, \mathbf{Y}) \right) + d\left( \sum_{m=1}^K \tilde{\mathbf{w}}^m \tilde{P}_{RD}^m, \mathbf{Y} \right), \quad (6)$$

where  $d$  is an appropriate loss function,  $\mathbf{Y}$  is the ground-truth mask,  $\tilde{\mathbf{w}}^m$  is the weight for  $\tilde{P}_{RD}^m$ ,  $\tilde{P}_{R-b}^m$  and  $\tilde{P}_{D-b}^m$  denote the predictions by  $\tilde{F}_R^m$  and  $\tilde{F}_D^m$  in the CA-Fuse block, respectively. This joint loss enables the cross-modal and cross-level combinations to be complementary for better inference.

## 4 EXPERIMENTS

In this section, we will introduce the implementation details, experimental comparisons and ablation studies to verify the advantages of the proposed method to learn, select and fuse cross-modal complements and the promise in zero-shot saliency detection from depth images.

### 4.1 Dataset and Evaluation Metrics

We evaluate our model on three RGB-D benchmark datasets: **NLPR** [4] includes 1000 indoor/outdoor RGB-D pairs collected using Kinect; **NJUD** [21] and **STEREO** [24] datasets contains 2003 and 797 stereoscopic images respectively, which are generated from the Internet and 3D movies and an optical method is adopted to compute depth images. We follow the previous works [16], [26], [33] to randomly pick 650 and 1400 RGB-D pairs from the NLPR and NJUD datasets respectively and combine them as the training set. We adopt the Precision-Recall (PR) curve, the F-measure and the mean absolute error (MAE) scores as evaluation metrics. Concretely, each saliency map  $S$  will be binarized by a threshold. The converted binary mask will be compared to the ground-truth  $G$  to compute the precision and recall. By varying the threshold from 0 to 255, we can obtain a series of precision-recall pairs, thereby forming the PR curve. The formulation of the F-measure is

$$F_{\beta} = \frac{(1 + \beta^2) \cdot Precision \cdot Recall}{\beta^2 \cdot Precision + Recall}, \quad (7)$$

where  $\beta^2 = 0.3$  as suggested by [9], [10]. The saliency map and binary ground-truth are normalized to  $[0, 1]$  and the MAE is to measure the pixel-wise discrepancy between the saliency map  $\tilde{S}$  and the ground-truth mask  $\tilde{G}$  averagely:

$$MAE = \frac{1}{W \times H} \sum_{i=1}^W \sum_{j=1}^H |\tilde{S}(i, j) - \tilde{G}(i, j)|, \quad (8)$$

where  $W$  and  $H$  are the width and height of the saliency map.

### 4.2 Implementation Details

We conduct our experiments using the Caffe [48] toolbox on a workstation with two GTX 1070 GPUs. The learning rate for the teacher network, the hierarchical cross-modal distillation network and the final RGB-D salient object detection network are  $1 \times 10^{-7}$ ,  $1 \times 10^{-6}$  and  $2 \times 10^{-9}$ , respectively.

TABLE 1: The Parameters of the Adaptation Layers and the Cross-level Transition Layers

Level	Adaptation 1	Adaptation 2	Transition layer
CA-Fuse 6	—	—	384, 1×1
CA-Fuse 5	384, 1×1	—	384, 1×1
CA-Fuse 4	384, 3×3	384, 3×3	256, 1×1
CA-Fuse 3	192, 3×3	192, 3×3	128, 1×1
CA-Fuse 2	128, 3×3	128, 3×3	—

For a fair comparison with the previous works based on the VGG network, we also adopt the VGG model as the backbone for both modalities and the detailed hierarchical cross-modal transfer architecture is illustrated in Fig. 2. The trunk inherits five convolutional blocks from the original VGG model. We add a new 512 13×13 convolutional layer for perceiving precedent features globally to enhance the localization ability. Then the strategy similar to [49] is leveraged to generate side outputs for each level. Specifically, the last layer in each convolutional block (e.g., Conv4\_3 and Conv2\_2) will be appended with one or two adaptation layers to the backbone. The details of the adaptation layers are shown in Table 1. We firstly train the teacher network with the architecture shown in the left of Fig.2. Concretely, the adapted features are used to infer level-specific saliency  $P_R^i$  via a 1×1 convolutional layer. Considering that it may be difficult for the first convolutional block to provide reliable cues, we do not involve it into inference for the teacher, student and the final RGB-D fusion network. Following Eq. (1),  $P_R^i$  will be combined with the predictions from deeper layers to generate the side-out  $\tilde{P}_R^i$  (refer to the BPDC module in Fig. 3 for implementation details). Accordingly, the loss function for the teacher network consists of the distance between the ground-truth mask and each side-out as well as the joint prediction combining all the side-outs as another constraint term:

$$L_{Teac} = \sum_{i=2}^K d(\tilde{P}_R^i, \mathbf{Y}) + d\left( \sum_{i=2}^6 \tilde{\mathbf{w}}_R^i \tilde{P}_R^i, \mathbf{Y} \right), \quad (9)$$

where  $\tilde{\mathbf{w}}_R^i$  is the weight for  $\tilde{P}_R^i$ .

The architecture of the student stream inherits the one of the teacher stream. When training the hierarchical cross-modal distillation network, the teacher stream is frozen. We adopt the cross-entropy loss for optimization when training the teacher network and the RGB-D fusion network:

$$d(\tilde{P}, \mathbf{Y}) = \mathbf{Y} \log \tilde{P} + (1 - \mathbf{Y}) \log(1 - \tilde{P}) \quad (10)$$

### 4.3 On the Hierarchical Cross-modal Distillation Schema

#### 4.3.1 Does the Student Network Learn Specific Cues?

The first question we want to investigate is whether the proposed hierarchical cross-modal distillation scheme can encourage the student network to learn specific cues to complement the source modality. Fig. 4 shows the individual inference from each level without combining with the predictions from deeper layers. It is not difficult to note that the side-outs from the student network present different



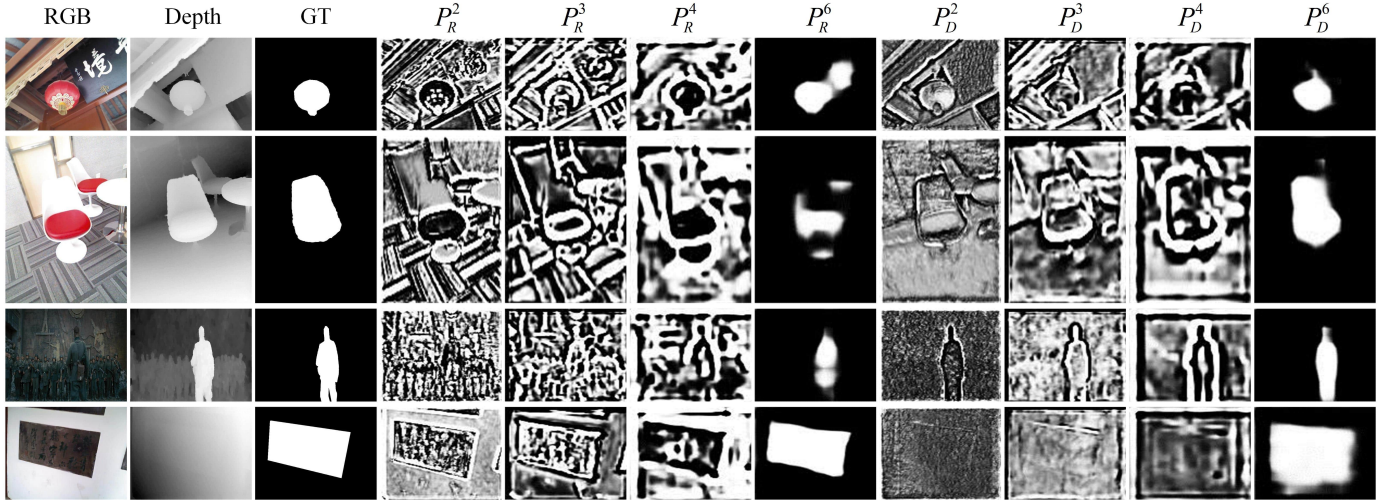


Fig. 4: The individual inference from the teacher and student streams in the hierarchical cross-modal distillation network.

TABLE 2: The Performance of Zero-shot Saliency Detection from Depth Images

Scheme	$F_\beta$			MAE		
	NLPR	NJUD	STEREO	NLPR	NJUD	STEREO
A	0.705	0.753	0.733	0.096	0.131	0.142
B	0.511	0.598	0.576	0.158	0.188	0.192
C	0.747	0.790	-	0.088	0.109	-

patterns in contrast to the ones generated by the teacher. More specifically, the shallow layers of the student network are only responsive to the depth variations while insensitive to the color/texture changes and the deeper layers, which are more responsible for locating the salient object, pay more attention to the object with distinguished depth. These differences demonstrate that the student network explores depth-specific saliency cues in each level effectively, which are complementary to the ones from the paired modality. Also, the cross-modal complementarity resides in multiple levels. These observations verify our motivations that a selector is in demand for highlighting the real complementary cues. Besides, a sufficient fusion scheme is also necessary considering the cross-modal complementarity in multiple layers.

### 4.3.2 For Zero-shot Saliency Detection from Depth Images

Our cross-modal transfer learning scheme allows zero-shot saliency detection for new modalities (e.g., depth). To verify this advantage, we combine the RGB salient object detection datasets including MSRA10K [9], ECSSD [50] and SED2 [51] to train the teacher network. Then we only use the RGB-D pairs from the RGB-D salient object detection training set to train the hierarchical cross-modal distillation network. We test the transferred student network with HHA (depth) images (noted as “A”). We also test the teacher network with the HHA images as inputs for comparison (noted as “B”). As shown in Table 2, the significant outperformance of the scheme “A” than “B” verifies the cross-modal discrep-

TABLE 3: Comparisons of Different Initialization Schemes for the Depth-induced Saliency Detection and RGB-D Saliency Detection Networks

Scheme	$F_\beta$			MAE		
	NLPR	NJUD	STEREO	NLPR	NJUD	STEREO
D-(A)	Fail to converge					
D-(B)	<u>0.747</u>	<u>0.780</u>	<u>0.732</u>	<u>0.080</u>	<u>0.096</u>	<u>0.127</u>
D-(C)	0.784	0.796	0.739	0.069	0.089	0.122
<i>fine</i> (D-(C))	<b>0.792</b>	<b>0.807</b>	<b>0.737</b>	<b>0.066</b>	<b>0.082</b>	<b>0.118</b>
RD-(A)	<u>0.842</u>	<u>0.854</u>	<u>0.868</u>	<u>0.054</u>	<u>0.063</u>	<u>0.066</u>
RD-(B)	0.861	0.860	0.877	0.049	0.061	0.061
RD-(C)	<b>0.872</b>	<b>0.871</b>	<b>0.880</b>	<b>0.046</b>	<b>0.057</b>	<b>0.060</b>

Notes:

- D-(B): Depth-induced saliency detection network initialized by task-adapted ImageNet pre-trained weights .
- *fine*(D-(C)): Depth-induced saliency detection network initialized by the proposed cross-modal distillation scheme.
- RD-(A): RGB-D saliency detection network initialized by ImageNet pre-trained weights.
- RD-(C): RGB-D saliency detection network initialized by the proposed cross-modal distillation scheme.

The comparison between them also demonstrates the notable improvement benefited from the distillation scheme. More details please refer to Section 4.3.3.

any and denotes the success of the proposed cross-modal transfer method. It is able to encourage modal-specific representations and inference, offering the promise in zero-shot saliency detection for new modalities.

We also report the RGB-only saliency results on the target dataset as another baseline (using RGB images as inputs to feed the teacher network, denoted as C). The RGB detector trained on the source dataset obtains satisfactory performance on the RGB images, while our zero-shot saliency detector on depth images achieve comparable

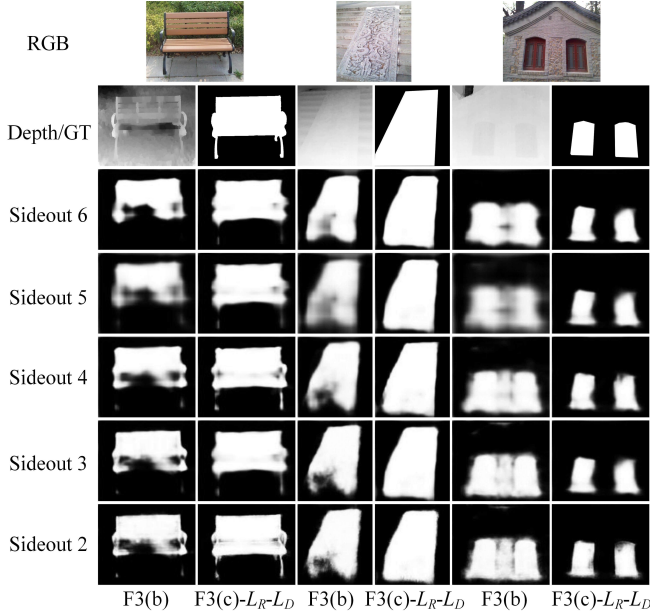


Fig. 5: Analyze the components in the CA-Fuse block visually.

detection performance (noted as A). Note that the STEREO dataset contains some same images in the source RGB saliency dataset. So we do not list the comparison results on this dataset.

### 4.3.3 Advantages as A Pre-training Scheme

In this section, we report the advantages of the proposed hierarchical cross-modal distillation scheme as a pre-training method for depth-induced saliency detection. We involve other two initialization strategies for comparison: 1) **D-(A)**: Random initialization; 2) **D-(B)**: Using the weights of the RGB-induced saliency detection network as initialization; 3) **D-(C)**: Using the proposed hierarchical cross-modal distillation schema. As shown in Table 3, we find that with ground-truth masks in each level as supervision, the convergence cannot be guaranteed when training the depth CNN from scratch, even we carefully tune the learning parameters. Compared to finetuning the RGB CNN, a huge improvement can be observed by using the hierarchical cross-modal distillation scheme, which demonstrates the efficacy of our cross-modal adaptation strategy. This improvement indicates that the side-outs from the teacher network serve as better guidance than the more correct ground-truth. We attribute this superiority to the progressive enhancement across the side-outs from the teacher. Compared to supervising each level with the same ground-truth, these side-outs become more direct and level-specific supervisions. The evolution of the side-outs demystifies level-specific contributions and cross-level collaborations explicitly. As a result, the goal of each level of the student is simplified as mimicking level-specific inference. For example, the goal of shallow layers is to learn low-level features for identifying object edges, which is a much easier task for them than predicting the completed saliency map. Additionally, finetuning the adapted student stream with ground-truth masks (denoted as “*fine*(D-(C))” allows further improvement.

TABLE 4: Analyze the Components in the CA-Fuse Block Quantitatively

Block	$F_\beta$			MAE		
	NLPR	NJUD	STEREO	NLPR	NJUD	STEREO
<b>F3(b)</b>	0.842	0.843	0.860	0.056	0.065	0.070
<b>F3(c)-<math>L_R-L_D</math></b>	0.867	0.866	0.878	0.048	0.059	0.062
<b>F3(c)</b>	<b>0.872</b>	<b>0.871</b>	<b>0.880</b>	<b>0.046</b>	<b>0.057</b>	<b>0.060</b>
<b>F3(c)-<math>\tilde{F}_{m,RD}^{m+1}</math></b>	0.867	0.862	0.878	0.047	0.059	0.060

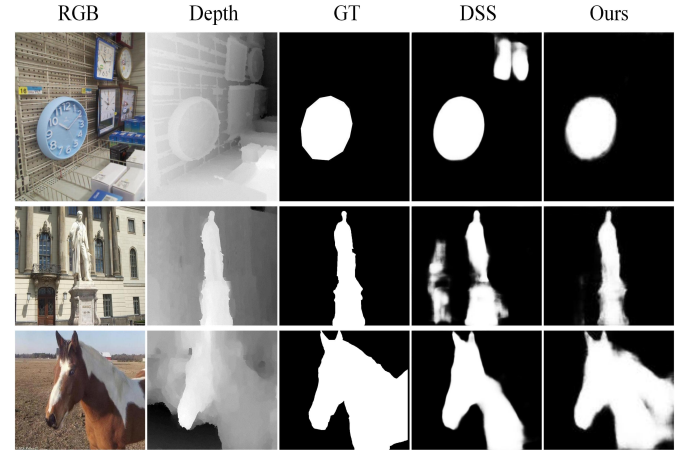


Fig. 6: Visual comparison to state-of-the-art RGB salient object detection method.

We also report the performance of using the proposed cross-modal transfer learning method as pre-training for the final RGB-D salient object detection network. We involve other two strategies for comparison. **RD-(A)**: Both the RGB and the depth streams are initialized by the VGG model without respective finetuning with the RGB-D saliency datasets. This strategy is adopted in [16]; **RD-(B)**: Stage-wise training. It means finetuning the RGB stream with the VGG model as initialization firstly. Then we train the depth stream starting from the well-trained RGB weights. This strategy is widely adopted in the previous works such as [26], [33]; **RD-(C)**: Using the trained hierarchical cross-modal distillation network as initialization. With the three initialization schemes, we then train the RGB-D fusion network using the RGB-D pairs and ground-truth. The comparison in Table 3 showcases the outperformance of the proposed cross-modal transfer schema, suggesting its success in learning better modal-specific representations.

### 4.4 On the CA-Fuse Block

In this section, we analyze the components in the CA-Fuse block. We first study the importance of introducing cross-modal residual functions. Fig. 5 illustrates the side outputs from each level with different designs shown in Fig. 3. The columns indexed as “F3(b)” show that the saliency maps inferred in the top-down pattern can be basically refined from coarse to fine with the help of the added supervision in each level and the cross-level combinations. However, the salient objects are not uniformly highlighted and some background regions are failed to be identified,



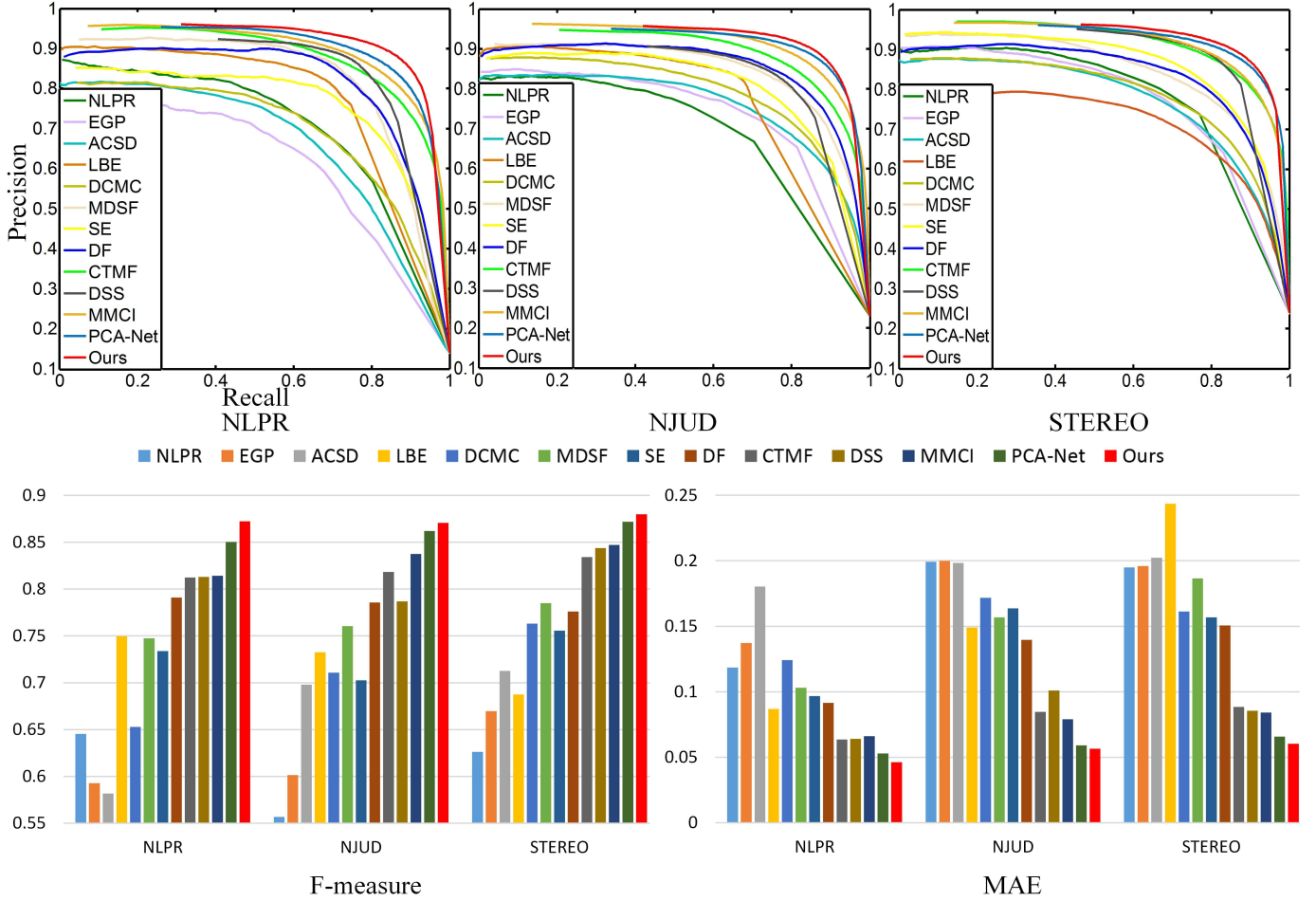


Fig. 7: Quantitative comparison to state-of-the-art RGB and RGB-D salient object detection methods. The MDSF only reports the results on the NLPR and NJUD datasets.

suggesting that directly concatenating cross-modal features is incapable of capturing the complementary information sufficiently. Then we improve the “F3(b)” block by adding cross-modal residual connections and this variant is denoted as “F3(c)- $L_R-L_D$ ”. Benefiting from the cross-modal residual functions, the complementary cues from both modalities are incorporated more easily, resulting in more informative multi-modal fusion. The comparison between “F3(b)” and “F3(c)- $L_R-L_D$ ” in Table 4 verifies the large performance gains from the cross-modal residual connections. Moreover, the comparison between “F3(c)- $L_R-L_D$ ” and “F3(c)” verifies the benefits of adding supervisions on the RGB and depth branches ( $L_R$  and  $L_D$ ), which further boost the emergence of the complementary cues from the paired modality. Another question we want to study is whether it is beneficial to transmit cross-level features to the adjacent shallower layer? To answer this question, we remove the  $\tilde{F}_{m, RD}^{m+1}$  in Fig. 3(c). Accordingly,  $\tilde{F}_R^m$  and  $\tilde{F}_D^m$  will be concatenated for joint inference. We denote this variant as “F3(c)- $\tilde{F}_{m, RD}^{m+1}$ ”. The quantitative comparison in Table 4 reports the noticeable gains by transmitting the cross-level features. We attribute this improvement to the richer RGB-D representations due to combining cross-level features.

#### 4.5 Comparison to State-of-the-art Methods

We compare our model (the variant “RD-(C)”) to 11 state-of-the-art RGB-D salient object detection methods: NLPR [4], EGP [52], ACSD [21], DCMC [22], LBE [17], MDSF [23], SE [53], DF [25], CTMF [26], MMCI [33] and our preliminary work PCA-Net [16], among which DF [25], CTMF [26], MMCI [33] and PCA-Net [16] are CNN-based methods. We also compare our method with a state-of-the-art RGB salient object detection model DSS [49] to verify the benefits of the synchronized depth data. Fig. 6 presents the comparison to the RGB-induced saliency visually. It can be noted that when the salient object and the background are with similar appearance or the background is seriously cluttered or the salient object is non-uniform, it is difficult to locate the salient object correctly and highlight the salient regions uniformly by relying on RGB inputs only. In these scenes, our model effectively incorporates the complementary cues from the paired depth data to overcome these deficiencies to identify the real salient object and highlight the salient regions consistently. The quantitative comparison in Fig. 7 shows that our proposed method outperforms others significantly. Compared to other RGB-D salient object detection methods, the proposed one holds distinguished advantages

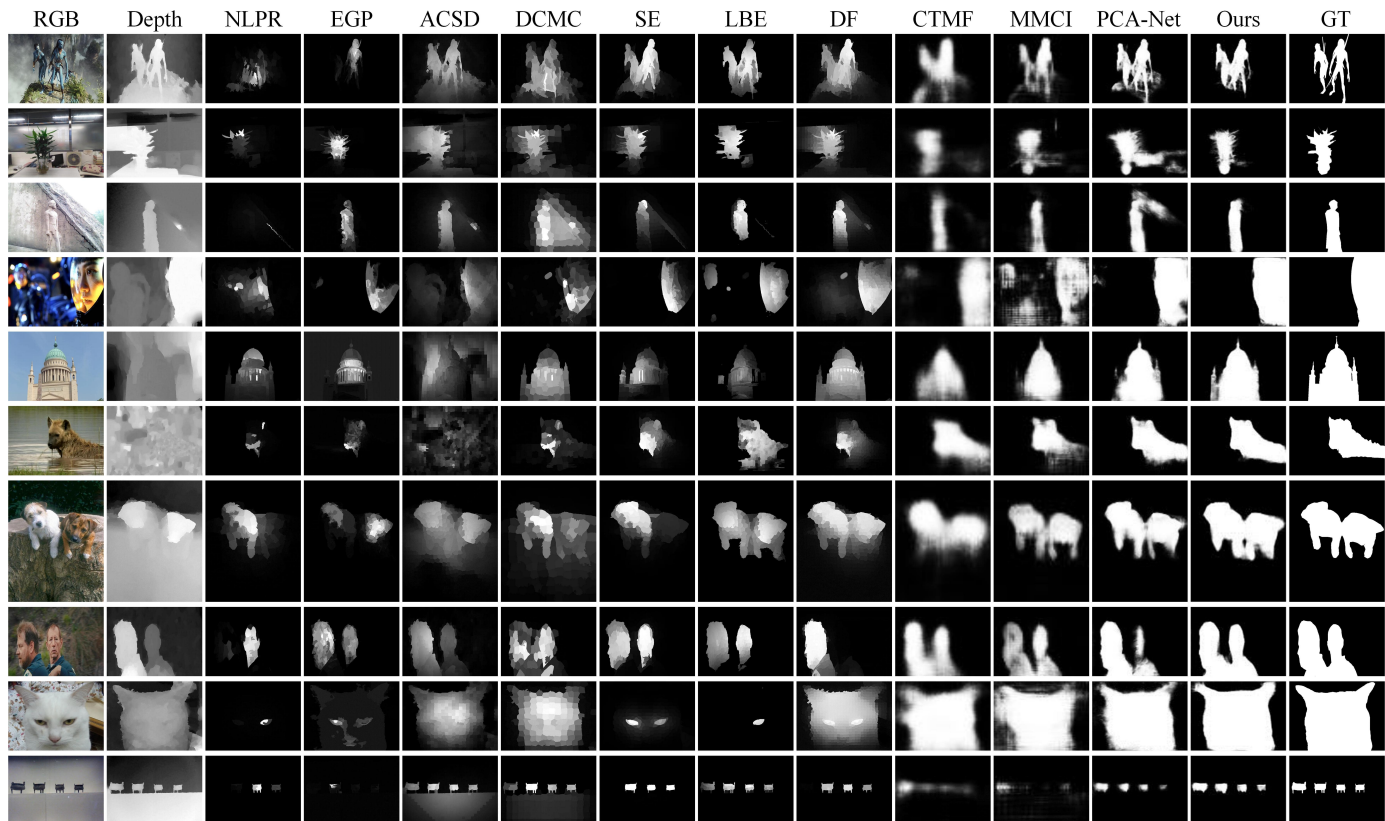


Fig. 8: Visual comparison to state-of-the-art RGB-D salient object detection methods.

in learning, selecting and fusing cross-modal complements. The methods [4], [17], [21], [22], [23], [52], [53] based on handcrafted RGB-D features are easy to be confused by complex background and intra-variant salient objects due to the lack of high-level global contexts. Previous CNN-based methods [25], [26] that combine cross-modal features only in a single level are incapable of capturing the cross-modal complementarity residing in high-level contexts and low-level spatial cues simultaneously. The “early fusion” schema adopted in [25] results in inconsistent highlighting of salient regions and the “late fusion” strategy used in [26] leads to severely blurred saliency maps. Although the work [33] remedies this shortcoming by designing two branches for global reasoning and local capturing respectively, it only leverages the last fully connected layer and an intermediate convolutional layer for joint inference. The final results are combined by directly summing the results from two branches, which is unlikely to combine local spatial cues and global contexts robustly. In contrast, our preliminary work [16] involves the information in all layers via a top-down path, which is able to progressively select and fuse the complements from each modal/level and refine the saliency maps gradually. By further using the hierarchical cross-modal distillation schema proposed in this extended work, the salient object is better located. Besides, the saliency maps are more uniform and carry better details than the ones generated by [16], implying the advantages of the proposed cross-modal transfer scheme in learning better modal-specific representations. In various challenging

scenes shown in Fig. 8, such as the background is complex (the 1<sup>st</sup>-2<sup>nd</sup> rows); the salient object and background have indistinguishable appearance or depth (the 3<sup>rd</sup>-4<sup>th</sup> and 5<sup>th</sup>-6<sup>th</sup> rows, resp.); the appearance or depth in the salient objects is non-uniform (the 7<sup>th</sup> and 8<sup>th</sup> rows, resp.); large/small salient objects (the 9<sup>th</sup> and 10<sup>th</sup> rows, resp.); multiple separated salient objects (the 10<sup>th</sup> row). In these cases, our proposed model can learn rich representations from each modality, select the complementary cues and fuse them informatively for successful joint inference.

## 5 CONCLUSION

In this paper, we propose a comprehensive view and a systematic solution for RGB-D salient object detection. The philosophy in designing an RGB-D system is generalized as three keys: modal-specific representations learning, complementary information selection and cross-modal complements fusion. Accordingly, we propose a new cross-modal transfer learning scheme, an explicit cross-modal complementarity selector and a sufficient cross-modal cross-level fusion pattern. The proposed solution solves the problems of zero-shot detection and multi-modal fusion jointly. We believe the insights provided from this work will allow us to learn better representations from new unlabeled modalities and more sufficient fusion for other multi-modal systems.

## ACKNOWLEDGMENTS

This work was supported by the Research Grants Council of Hong Kong (Project No CityU 11205015 and CityU

11255716).

## REFERENCES

- [1] J. Han, L. Shao, D. Xu, and J. Shotton, "Enhanced computer vision with microsoft kinect sensor: A review," *IEEE Trans. Cybern.*, vol. 43, no. 5, pp. 1318–1334, 2013.
- [2] S. Gupta, P. Arbeláez, R. Girshick, and J. Malik, "Indoor scene understanding with rgb-d images: Bottom-up segmentation, object detection and semantic segmentation," *Int. J. Comput. Vis.*, vol. 112, no. 2, pp. 133–149, 2015.
- [3] M. Camplani, S. L. Hannuna, M. Mirmehdi, D. Damen, A. Paiement, L. Tao, and T. Burghardt, "Real-time rgb-d tracking with depth scaling kernelised correlation filters and occlusion handling," in *Proc. British Mach. Vis. Conf.*, 2015, pp. 145–1.
- [4] H. Peng, B. Li, W. Xiong, W. Hu, and R. Ji, "Rgb-d salient object detection: a benchmark and algorithms," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 92–109.
- [5] L. Shao and M. Brady, "Specific object retrieval based on salient regions," *Pattern Recognit.*, vol. 39, no. 10, pp. 1932–1948, 2006.
- [6] V. Mahadevan, N. Vasconcelos *et al.*, "Biologically inspired object tracking using center-surround saliency mechanisms," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 3, pp. 541–554, 2013.
- [7] J. Harel, C. Koch, and P. Perona, "Graph-based visual saliency," in *Proc. Adv. Neural Inf. Process. Syst.*, 2007, pp. 545–552.
- [8] J. Yang and M.-H. Yang, "Top-down visual saliency via joint crf and dictionary learning," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 3, pp. 576–588, 2017.
- [9] M.-M. Cheng, N. J. Mitra, X. Huang, P. H. Torr, and S.-M. Hu, "Global contrast based salient region detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 3, pp. 569–582, 2015.
- [10] A. Borji, M.-M. Cheng, H. Jiang, and J. Li, "Salient object detection: A benchmark," *IEEE Trans. Image Process.*, vol. 24, no. 12, pp. 5706–5722, 2015.
- [11] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik, "Learning rich features from rgb-d images for object detection and segmentation," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 345–360.
- [12] H. Fu, D. Xu, S. Lin, and J. Liu, "Object-based rgb-d image co-segmentation with mutex constraint," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2015, pp. 4428–4436.
- [13] H. Fu, D. Xu, and S. Lin, "Object-based multiple foreground segmentation in rgb-d video," *IEEE Trans. Image Process.*, vol. 26, no. 3, pp. 1418–1427, 2017.
- [14] R. Cong, J. Lei, H. Fu, W. Lin, Q. Huang, X. Cao, and C. Hou, "An iterative co-saliency framework for rgb-d images," *IEEE Trans. Cybern.*, 2017.
- [15] R. Cong, J. Lei, H. Fu, Q. Huang, X. Cao, and C. Hou, "Co-saliency detection for rgb-d images based on multi-constraint feature matching and cross label propagation," *IEEE Trans. Image Process.*, vol. 27, no. 2, pp. 568–579, 2018.
- [16] H. Chen and Y. Li, "Progressively complementarity-aware fusion network for rgb-d salient object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2018, pp. 3051–3060.
- [17] D. Feng, N. Barnes, S. You, and C. McCarthy, "Local background enclosure for rgb-d salient object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2016, pp. 2343–2350.
- [18] C. Lang, T. V. Nguyen, H. Katti, K. Yadati, M. Kankanhalli, and S. Yan, "Depth matters: Influence of depth cues on visual saliency," in *Proc. Eur. Conf. Comput. Vis.*, 2012, pp. 101–115.
- [19] K. Desingh, K. M. Krishna, D. Rajan, and C. Jawahar, "Depth really matters: Improving visual salient region detection with depth," in *Proc. British Mach. Vis. Conf.*, 2013.
- [20] A. Ciptadi, T. Hermans, and J. M. Rehg, "An in depth view of saliency," in *Proc. British Mach. Vis. Conf.*, 2013.
- [21] R. Ju, L. Ge, W. Geng, T. Ren, and G. Wu, "Depth saliency based on anisotropic center-surround difference," in *Proc. IEEE Int. Conf. Image Process.*, 2014, pp. 1115–1119.
- [22] 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," *Signal Process. Lett.*, vol. 23, no. 6, pp. 819–823, 2016.
- [23] 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.
- [24] Y. Niu, Y. Geng, X. Li, and F. Liu, "Leveraging stereopsis for saliency analysis," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2012, pp. 454–461.
- [25] L. Qu, S. He, J. Zhang, J. Tian, Y. Tang, and Q. Yang, "Rgb-d salient object detection via deep fusion," *IEEE Trans. Image Process.*, vol. 26, no. 5, pp. 2274–2285, 2017.
- [26] 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.*, 2017.
- [27] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2012, pp. 1097–1105.
- [28] A. Wang, J. Cai, J. Lu, and T.-J. Cham, "Mmss: Multi-modal sharable and specific feature learning for rgb-d object recognition," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015, pp. 1125–1133.
- [29] S.-J. Park, K.-S. Hong, and S. Lee, "Rdfnet: Rgb-d multi-level residual feature fusion for indoor semantic segmentation," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017.
- [30] X. Xu, Y. Li, G. Wu, and J. Luo, "Multi-modal deep feature learning for rgb-d object detection," *Pattern Recognit.*, vol. 72, pp. 300–313, 2017.
- [31] H. Zhu, J.-B. Weibel, and S. Lu, "Discriminative multi-modal feature fusion for rgb-d indoor scene recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2016, pp. 2969–2976.
- [32] X. Fan, Z. Liu, and G. Sun, "Salient region detection for stereoscopic images," in *Proc. IEEE Int. Conf. Digital Signal Process.*, 2014, pp. 454–458.
- [33] 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 Recognit.*, 2018.
- [34] Y. Cheng, R. Cai, Z. Li, X. Zhao, and K. Huang, "Locality-sensitive deconvolution networks with gated fusion for rgb-d indoor semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, vol. 3, 2017.
- [35] D. Lin, G. Chen, D. Cohen-Or, P.-A. Heng, and H. Huang, "Cascaded feature network for semantic segmentation of rgb-d images," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017, pp. 1320–1328.
- [36] G. Lin, A. Milan, C. Shen, and I. Reid, "Refinenet: Multi-path refinement networks for high-resolution semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, vol. 1, no. 2, 2017, p. 3.
- [37] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [38] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," *arXiv preprint arXiv:1412.6550*, 2014.
- [39] Q. Li, S. Jin, and J. Yan, "Mimicking very efficient network for object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2017, pp. 7341–7349.
- [40] Z. Huang and N. Wang, "Like what you like: Knowledge distill via neuron selectivity transfer," *arXiv preprint arXiv:1707.01219*, 2017.
- [41] C. M. Christoudias, R. Urtasun, M. Salzmann, and T. Darrell, "Learning to recognize objects from unseen modalities," in *Proc. Eur. Conf. Comput. Vis.*, 2010, pp. 677–691.
- [42] R. Socher, M. Ganjoo, C. D. Manning, and A. Ng, "Zero-shot learning through cross-modal transfer," in *Proc. Adv. Neural Inf. Process. Syst.*, 2013, pp. 935–943.
- [43] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, "Multimodal deep learning," in *Proc. Int. Conf. Mach. Learn.*, 2011, pp. 689–696.
- [44] J. Hoffman, S. Gupta, and T. Darrell, "Learning with side information through modality hallucination," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2016, pp. 826–834.
- [45] N. C. Garcia, P. Morerio, and V. Murino, "Modality distillation with multiple stream networks for action recognition," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 103–118.
- [46] S. Gupta, J. Hoffman, and J. Malik, "Cross modal distillation for supervision transfer," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2016, pp. 2827–2836.
- [47] K. Lenc and A. Vedaldi, "Understanding image representations by measuring their equivariance and equivalence," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2015, pp. 991–999.
- [48] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proc. ACM Int. Conf. Multimedia*, 2014, pp. 675–678.

- [49] Q. Hou, M.-M. Cheng, X. Hu, A. Borji, Z. Tu, and P. Torr, "Deeply supervised salient object detection with short connections," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2017, pp. 5300–5309.
- [50] Q. Yan, L. Xu, J. Shi, and J. Jia, "Hierarchical saliency detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2013, pp. 1155–1162.
- [51] S. Alpert, M. Galun, R. Basri, and A. Brandt, "Image segmentation by probabilistic bottom-up aggregation and cue integration," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2007, pp. 1–8.
- [52] J. Ren, X. Gong, L. Yu, W. Zhou, and M. Ying Yang, "Exploiting global priors for rgb-d saliency detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, 2015, pp. 25–32.
- [53] J. Guo, T. Ren, and J. Bei, "Salient object detection for rgb-d image via saliency evolution," in *Proc. IEEE Int. Conf. Multimedia and Expo*, 2016, pp. 1–6.