

Adaptive Fusion for RGB-D Salient Object Detection

Ningning Wang, Xiaojin Gong*

College of Information Science & Electronic Engineering, Zhejiang University, Hangzhou, China.

{wangnn,gongxj}@zju.edu.cn

Abstract—RGB-D salient object detection aims to identify the most visually distinctive objects in a pair of color and depth images. Based upon an observation that most of the salient objects may stand out at least in one modality, this paper proposes an adaptive fusion scheme to fuse saliency predictions generated from two modalities. Specifically, we design a two-streamed convolutional neural network (CNN), each of which extracts features and predicts a saliency map from either RGB or depth modality. Then, a saliency fusion module learns a switch map that is used to adaptively fuse the predicted saliency maps. A loss function composed of saliency supervision, switch map supervision, and edge-preserving constraints is designed to make full supervision, and the entire network is trained in an end-to-end manner. Benefited from the adaptive fusion strategy and the edge-preserving constraint, our approach outperforms state-of-the-art methods on three publicly available datasets.

I. INTRODUCTION

Salient object detection aims to automatically identify the most attractive regions in images like human visual systems. It can serve as a useful pre-processing step for various computer vision applications such as image segmentation [1], person re-identification [2], object localization [3] and tracking [4], and therefore has received considerable attention. Although great progress has been made in this field, most works [5]–[11] focus merely on color images. When objects share similar appearances with their surroundings or present with complex background, the algorithms based on color images often fail to distinguish them as salient objects.

The above-mentioned challenges can be overcome to a large extent if depth information is available. In recent years, robust ranging sensors such as stereo cameras, RGB-D cameras, and lidars make it easy to collect paired color and depth images. RGB-D saliency detection has consequently been attracting more and more research interest. Published literatures made efforts on modeling depth-induced saliency detection [12]–[14] and fusing multi-modalities [15]–[19]. However, existing works performed fusion via either directly concatenating color and depth features, or element-wise multiplication/addition of predictions generated by the two modalities. Such fusion strategies are inadequate to combine complementary information from two modalities, leaving a room for performance improvement.

When observing objects in paired color and depth images, we roughly classify scenes into four categories: 1) Objects have distinguishable appearances in both modalities; 2) Objects have close depth values but distinguishable color appearances with backgrounds; 3) Objects share similar color

appearances with backgrounds but have different depth values; and 4) Scenes are cluttered in both color and depth images, as shown in Figure 1. For the first three scenarios, salient objects can be correctly detected at least in one modality when using state-of-the-art single-model based saliency detection methods. It implies that good results can be obtained for these scenarios if an algorithm could adaptively choose the predictions from one or the other modality.

Our work is inspired by the above-mentioned observation. In order to make fusion adaptive, We propose an end-to-end framework that consists of a two-streamed convolutional neural network (CNN) and a saliency fusion module to predict and fuse saliency predictions. Our main contributions lie in the following aspects:

- We design a two-streamed CNN to predict a saliency map from each modality separately. Each unimodal saliency prediction stream adopts a multi-scale feature aggregation strategy to make feature extraction and saliency prediction effective, while keeping the architecture simple.
- We propose a saliency fusion module that learns a switch map to adaptively fuse the predicted saliency maps. A pseudo ground truth switch map is constructed to supervise the learning so that the learned switch map can predict the weights for fusing RGB and depth saliency maps.
- The proposed approach is validated on three publicly available datasets, including NJUD [12], NLPR [15], and STEREO [20]. Experimental results show that our approach consistently outperforms state-of-the-art methods on all datasets.
- To make our work reproducible, we release our source code at https://github.com/Lucia-Ningning/Adaptive_Fusion_RGBD_Saliency_Detection.

II. RELATED WORK

A. RGB Saliency Detection

A great number of RGB salient object detection methods have been developed over the past decades. Traditional methods mainly rely on hand-crafted features and commit to mining effective priors such as center prior [5], contrast prior [6], [7], boundary and connectivity prior [21]. Owing to the deep learning revolution, CNN-based approaches have refreshed the previous records in recent years. Multi-scale features are first extracted by multiple CNNs and concatenated together, and then they are fed into a shallow network to predict

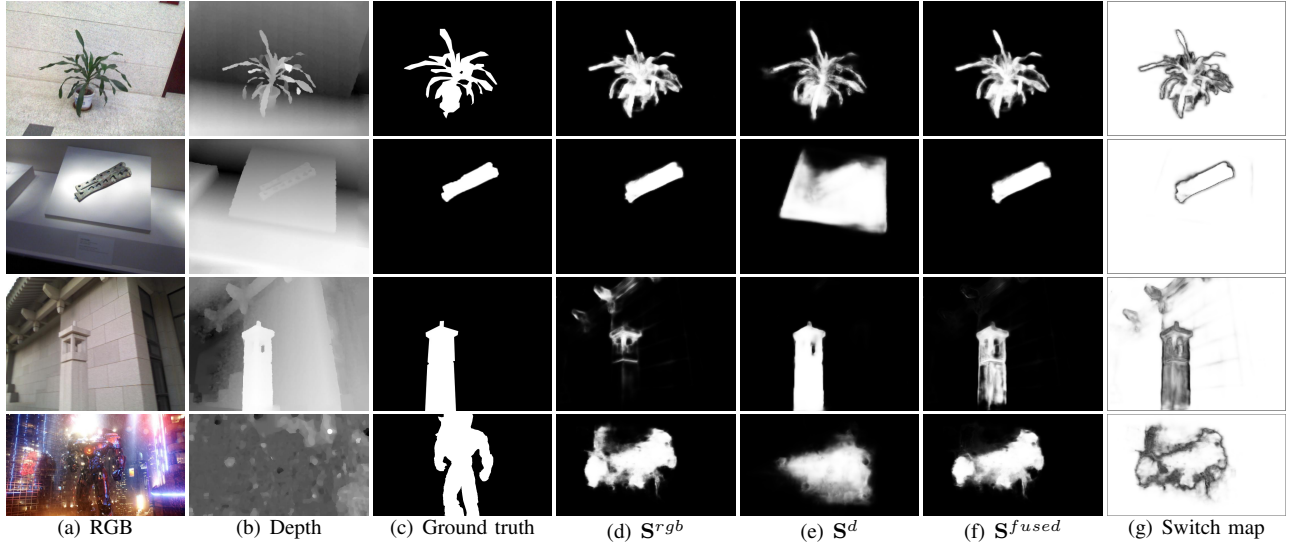


Fig. 1: Typical scenarios in RGB-D saliency object detection. Here, S^{rgb} denotes the result obtained by our RGB saliency prediction stream, S^d is the result from our depth saliency prediction stream, and S^{fused} is the final saliency detection result. Switch map is the map learned in our network for adaptive fusion.

saliency [8]. Whereafter, two-branched networks [10] [9] were designed to capture global and local context. In more recent years, deep hierarchical saliency networks (DHSNet) [22], short connections [11], and even more complicated structures such as Amulet [23] and agile Amulet [24] were developed to aggregate multi-scale features progressively and predict saliency within end-to-end frameworks. We adopt the progressive multi-scale feature aggregation strategy in our unimodal saliency prediction stream, but we keep the network as simple as possible.

B. RGB-D Saliency Detection

There are two major concerns existing in RGB-D saliency detection: 1) how to model the depth-induced saliency; and 2) how to fuse RGB and depth modalities for achieving better performance.

Regarding to the first problem, different features such as anisotropic center-surround difference [12] and local background enclosure (LBE) [13] were designed to evaluate saliency on depth maps. Global priors, including the normalized depth prior and the global-context surface orientation prior [14], were exploited as well. Although these features and priors are particularly effective for depth saliency detection, their performances are limited due to hand-crafted designs and multi-stage models.

For the second problem, existing approaches perform multi-modal fusion roughly at input, feature, or decision levels. For instance, Peng et al. [15] directly concatenated RGB and depth values and fed the 4-channel data into a multi-stage model to produce saliency maps. Qu et al. [16] extracted hand-crafted features from RGB-D superpixels and input them into a shallow CNN for feature combination and saliency regression. Han et al. [17] designed a two-streamed CNN to extract RGB and depth features separately and then fuse them with a joint

representation layer. Fusion in these methods is conducted with a single path. In order to enable sufficient fusion, a multi-scale multi-path fusion network [18] and a progressively complementarity-aware fusion network [19] were developed in more recent years. Although complicated architectures were designed, these methods perform fusion mainly rely on feature concatenation and element-wise addition/multiplication of prediction results. In contrast, we design an adaptive fusion scheme to fuse prediction results from RGB and depth modalities and achieve better detection performance.

III. THE PROPOSED METHOD

When a pair of RGB and depth images are given, we feed them into a two-streamed network for saliency detection. In each stream, features at different scales are progressively aggregated and a saliency map, S^{rgb} or S^d , is predicted based upon unimodal information. In addition, the last layer of RGB and depth features are concatenated to generate a switch map SW . The switch map further explicitly guides the fusion of S^{rgb} and S^d to produce the final saliency map S^{fused} . During training, all of the predicted saliency maps are supervised under the ground truth Y and the switch map is supervised with a pseudo ground truth constructed from Y and S^{rgb} . Figure 2 illustrates an overview of the proposed framework.

A. Unimodal Saliency Prediction Stream

A unimodal saliency prediction stream aims to predict a saliency map based upon a single modal information. Therefore, the design can be benefited from state-of-the-art RGB saliency detection methods. Our design adopts a multi-scale feature based saliency detection framework [11], [22]–[24], using an effective feature fusion strategy that progressively aggregates multi-scale features. In contrast to these methods, we keep our network structure as simple as possible.

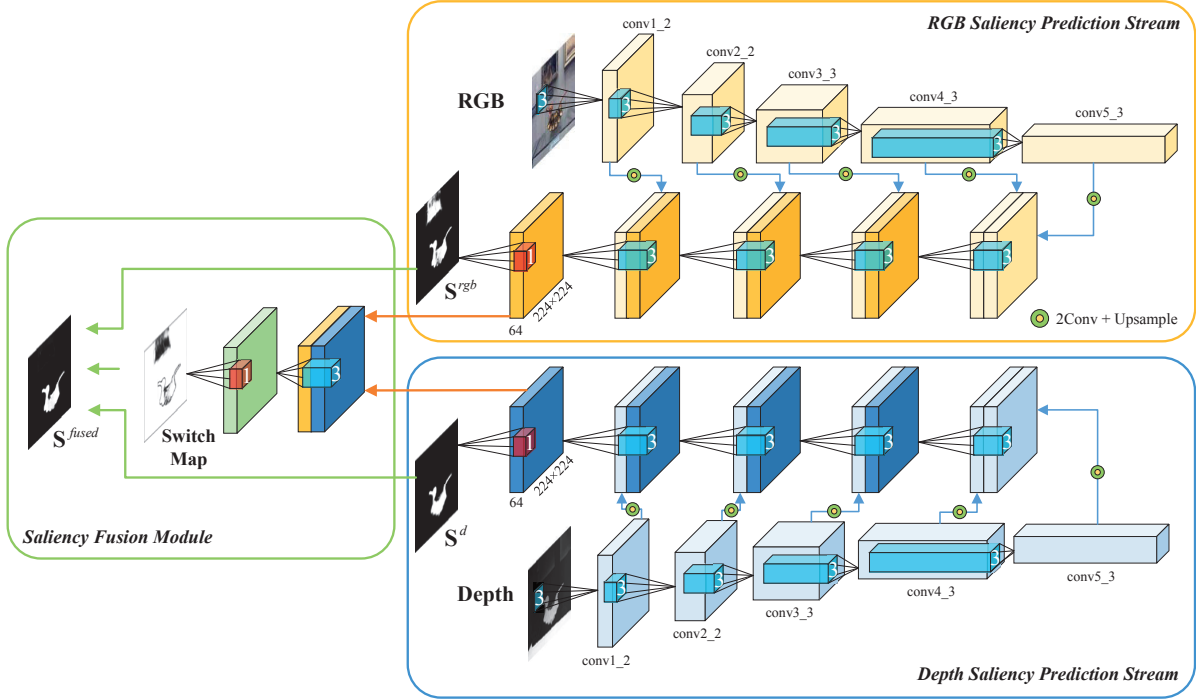


Fig. 2: The overview of our framework for RGB-D salient object detection.

Specifically, each stream is built upon the VGG-16 model [25] that contains 5 convolutional blocks. We drop the last pooling layer and the fully-connected layers to better fit for our task. Let us denote the outputs of each block, conv1_2, conv2_2, conv3_3, conv4_3, conv5_3, respectively, by $\mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3, \mathbf{A}_4, \mathbf{A}_5$. Each block also produces a side output \mathbf{F}_i by feeding \mathbf{A}_i into two extra convolutional layers and an up-sampling layer. The feature aggregation strategy progressively fuses the feature \mathbf{F}_i at scale i with the fused feature $\tilde{\mathbf{F}}_{i+1}$ from scale $i+1$. In the end, a saliency map \mathbf{S} is predicted based on the aggregated feature $\tilde{\mathbf{F}}_1$. Mathematically, we formulate the procedures of feature extraction and saliency prediction as follows:

$$\mathbf{F}_i = u(g(g(\mathbf{A}_i))) \quad 1 \leq i \leq 5, \quad (1)$$

$$\tilde{\mathbf{F}}_i = \begin{cases} g([\tilde{\mathbf{F}}_{i+1}, \mathbf{F}_i]) & 1 \leq i < 5 \\ \mathbf{F}_i & i = 5, \end{cases} \quad (2)$$

$$\mathbf{S} = h(\mathbf{W}_s * \tilde{\mathbf{F}}_1 + \mathbf{b}_s), \quad (3)$$

in which $g(\cdot)$ denotes the operations that consists of a 64-channel convolutional layer followed by a non-linear activation function. The kernel size of the convolution is 3×3 and the stride is 1. $u(\cdot)$ is an upsampling operation using bilinear interpolation. $[\cdot, \cdot]$ represents a channel-wise concatenation. \mathbf{W}_s and \mathbf{b}_s are, respectively, the parameters of the 1×1 kernel and the bias. $*$ represents the convolution operator and $h(\cdot)$ is the Sigmoid function.

This stream structure is applied to predict RGB saliency and depth saliency separately. The RGB saliency prediction stream takes a 3-channel color image as input while the depth stream inputs a 1-channel depth map. Except the inputs, these two streams share the same structure but with different parameter

values. In addition, it needs to be mentioned that we drop the superscription rgb or d in Eq.(1-3) for notational convenience.

B. Saliency Fusion Module

In contrast to previous RGB-D saliency detection works [15], [18] that fuse multi-modal predictions by element-wise addition or multiplication, we design a saliency fusion module that learns a switch map for adaptive fusion of the RGB saliency prediction \mathbf{S}^{rgb} and the depth saliency predictions \mathbf{S}^d . This module first concatenates the last layer features of two streams and then goes through a convolutional layer to learn a switch map \mathbf{SW} . In the end, a fused saliency map \mathbf{S}^{fused} is obtained. All operations in this module are formulated by:

$$\tilde{\mathbf{F}}^{sw} = g([\tilde{\mathbf{F}}_1^{rgb}, \tilde{\mathbf{F}}_1^d]), \quad (4)$$

$$\mathbf{SW} = h(\mathbf{W}_{sw} * \tilde{\mathbf{F}}^{sw} + \mathbf{b}_{sw}), \quad (5)$$

$$\mathbf{S}^{fused} = \mathbf{SW} \odot \mathbf{S}^{rgb} + (\mathbf{1} - \mathbf{SW}) \odot \mathbf{S}^d, \quad (6)$$

where $\tilde{\mathbf{F}}^{sw}$ represents the 64-channel feature fusing two modalities. $\tilde{\mathbf{F}}_1^{rgb}$ and $\tilde{\mathbf{F}}_1^d$ are, respectively, the features at the last layer of the color and depth streams. \mathbf{W}_{sw} and \mathbf{b}_{sw} are the parameters of the convolutional layer. \odot denotes the element-wise multiplication.

The design of the switch map is motivated by the observation mentioned in Sec. I. That is, good detection results are achieved in most scenarios if the algorithm can automatically choose the predictions from either RGB or depth modality. To this end, we construct a pseudo ground truth switch map \mathbf{Y}^{sw} to guide the learning of \mathbf{SW} . It is defined by

$$\mathbf{Y}^{sw} = \mathbf{S}^{rgb} \odot \mathbf{Y} + (\mathbf{1} - \mathbf{S}^{rgb}) \odot (\mathbf{1} - \mathbf{Y}). \quad (7)$$

\mathbf{Y}^{sw} gets 1 if the RGB saliency prediction \mathbf{S}^{rgb} and the ground truth \mathbf{Y} are both salient or nonsalient, and 0 otherwise. It means that if \mathbf{S}^{rgb} correctly identifies salient objects, then we choose the prediction from the RGB stream as the final result; otherwise, the prediction from the depth stream is chosen.

In implementation, the switch map is a 1-channel image whose pixel values are assigned in $[0, 1]$. Therefore, instead of alternatively choosing the prediction from one or the other modality, the switch map plays a role to adaptively weigh the RGB and depth predictions, and therefore the fused saliency map is a weighted sum of the two predictions.

C. Loss Function

During training, a set of samples $\mathcal{C} = \{(\mathbf{X}_i, \mathbf{D}_i, \mathbf{Y}_i)\}_{i=1}^N$ are given, in which N is the total number of samples. $\mathbf{X}_i = \{x_{i,j}\}_{j=1}^T$ and $\mathbf{D}_i = \{d_{i,j}\}_{j=1}^T$ are a pair of RGB and depth images with T pixels. $\mathbf{Y}_i = \{y_{i,j}\}_{j=1}^T$ is the corresponding binary ground truth saliency map, with 1 denoting salient pixels and 0 for the background. Our network is trained to generate an edge-preserving saliency map by learning a switch map and fusing two unimodal saliency predictions. Therefore, the loss function is designed to contain three terms: a saliency loss \mathcal{L}_{sal} , a switch loss \mathcal{L}_{sw} , as well as an edge-preserving loss \mathcal{L}_{edge} . That is,

$$\mathcal{L} = \mathcal{L}_{sal} + \mathcal{L}_{sw} + \mathcal{L}_{edge}. \quad (8)$$

Saliency Loss. There are three saliency maps produced in our network: \mathbf{S}^{rgb} , \mathbf{S}^d , and \mathbf{S}^{fused} . We use the ground truth to supervise each of them. A standard cross-entropy loss is adopted to compute the difference between predicted results and the ground truth. Therefore, The saliency loss is defined by

$$\mathcal{L}_{sal} = \mathcal{L}_{sal}^{rgb} + \mathcal{L}_{sal}^d + \mathcal{L}_{sal}^{fused}, \quad (9)$$

where

$$\mathcal{L}_{sal}^m = - \sum_{i=1}^N \sum_{j=1}^T (y_{i,j} \log \mathbf{S}_{i,j}^m + (1 - y_{i,j}) \log(1 - \mathbf{S}_{i,j}^m)). \quad (10)$$

Here, the superscript m denotes a modality that may be rgb , d , or $fused$. $\mathbf{S}_{i,j}^m$ represents the probability predicted by the modality m for pixel j in the i -th image to be salient.

Switch Loss. The switch map is supervised by the pseudo ground truth \mathbf{Y}^{sw} constructed in Eq.(7). We use the cross-entropy loss to penalize the learning of the switch map as well. The loss is defined by:

$$\mathcal{L}_{sw} = - \sum_{i=1}^N \sum_{j=1}^T (y_{i,j}^{sw} \log \mathbf{SW}_{i,j} + (1 - y_{i,j}^{sw}) \log(1 - \mathbf{SW}_{i,j})) \quad (11)$$

where $y_{i,j}^{sw}$ is the j -th pixel of the pseudo ground truth switch map for the i -th image. $\mathbf{SW}_{i,j}$ represents the probability for the pixel to choose the RGB prediction $\mathbf{S}_{i,j}^{rgb}$.

Edge-preserving Loss. The edge-preserving property has been considered in previous RGB saliency detection

works [23], [26] to obtain sharp salient object boundaries and improve detection performance. In contrast to these works that used superpixel boundaries as constraints [26] or adopted short connections in network for boundary refinement [23], we formulate the edge-preserving constraint as a loss term supervising the fused saliency map. It is defined by

$$\mathcal{L}_{edge} = \frac{1}{N} \sum_{i=1}^N \|\partial_x(\mathbf{S}_i^{fused}) - \partial_x(\mathbf{Y}_i)\|_2^2 + \|\partial_y(\mathbf{S}_i^{fused}) - \partial_y(\mathbf{Y}_i)\|_2^2, \quad (12)$$

where $\partial_x(\cdot)$ and $\partial_y(\cdot)$ are gradients in horizontal and vertical direction respectively. This loss preserves edges by minimizing the differences between the edges in the fused saliency maps and those in the ground truth maps.

D. Implementation Details

Our approach is implemented based upon TensorFlow [27]. We adopt the VGG-16 model [25] as the backbone for a fair comparison with previous works. All parameters except those in VGG-16 are initialized via Xavier [28]. Our entire network is trained in an end-to-end manner using the aforementioned loss function. The loss is optimized by the Adam optimizer [29] with a batch size of 8 and a learning rate of 10^{-4} . All input images are resized to the resolution of 224×224 for training and test. We conduct our experiments on a PC with a single NVIDIA GTX 1080Ti GPU. The test time for each RGB-D image pair takes only 0.03s.

IV. EXPERIMENTAL RESULTS

A. Datasets

To validate the proposed approach, we conduct a series of experiments on three publicly available datasets: NJUD [12], NLPR [15], and STEREO [20]. The NJUD dataset [12] contains 2003 binocular image pairs collected from Internet, 3D movies and photographs. NLPR [15] consists of 1000 images captured by Kinect, covering a variety of indoor and outdoor scenes under different illumination conditions. STEREO [20] provides the Web links for downloading stereoscopic images and a total of 797 pairs are gathered.

For a fair comparison to state-of-the-arts, we utilize the same data split as in [17]. The training set contains 1400 samples from the NJUD dataset and 650 samples from NLPR. 100 image pairs from NJUD and 50 image pairs from NLPR are sampled to form the validation set. The test set consists of the remaining data in these two datasets, together with the full STEREO dataset. In addition, we augment the training set by flipping all training samples horizontally.

B. Evaluation Metrics

We adopt the precision-recall (PR) curves, the F-measure score, and the mean absolute error (MAE) for performance evaluation. These metrics are widely used in saliency detection tasks. The PR curves are plotted by binarizing a predicted saliency map using 255 thresholds equally distributed in $[0, 1]$ and comparing the binarized map with the ground truth.

The F-measure is a weighted harmonic mean of the precision and recall, defined by

$$F_\beta = \frac{(1 + \beta^2) \cdot Precision \cdot Recall}{\beta^2 \cdot Precision + Recall}. \quad (13)$$

As done in previous works [17]–[19], β^2 is set to be 0.3 for emphasizing the importance of precision. We compare two kinds of F-measure scores, which are the maximum F-measure and the mean F-measure, respectively. The maximum F-measure is the highest score computed by the PR pairs in PR curves. The mean F-measure is computed by using an adaptive threshold that is set to be the sum of mean and standard deviation of each saliency map. The MAE [23] measures the saliency detection accuracy by

$$MAE = \frac{1}{T} \sum_{j=1}^T |S_j - Y_j|. \quad (14)$$

C. Ablation Study

We first conduct experiments to validate the effectiveness of the components in our proposed model. To this end, different settings are considered: 1) the full model, denoted by ‘AF’; 2) the model without edge-preserving loss, denoted by ‘AF-Edge’; 3) the model without switch map and edge-preserving loss, denoted by ‘AF-Edge-SW’; In this model, we concatenate the features from two streams and feed them into a 1×1 convolutional layer to predict the fused saliency map directly. 4) the model containing only the RGB saliency prediction stream, denoted by ‘ S^{rgb} ’; and 5) the one containing only the depth stream, denoted by ‘ S^d ’. Table I reports the mean F-measure scores for these models on three datasets.

TABLE I: The results for component analysis.

Settings	NJUD	NLPR	STEREO
S^{rgb}	0.854	0.857	0.874
S^d	0.800	0.754	0.770
AF-Edge-SW	0.872	0.865	0.879
AF-Edge	0.874	0.873	0.886
AF	0.878	0.881	0.891

The effectiveness of the saliency fusion module: The comparison of ‘AF-Edge-SW’ and ‘AF-Edge’ in Table I demonstrates the improvement in mean F-measure with our saliency fusion module. The results in Fig. 1 illustrate the fusion of S^{rgb} and S^d visually. When S^{rgb} correctly detects the salient objects, as the scenarios shown in the first two rows, our approach fuses more information from the RGB predictions by highlighting most regions in the switch maps. When objects share similar color appearances with backgrounds but have different depth values, as shown in the third row, our approach suppresses unreliable predictions in S^{rgb} by assigning low weights for these regions in the switch map. Thus, more information from S^d are fused. As expected, the proposed saliency fusion module can tackle these three types of scenarios correctly.

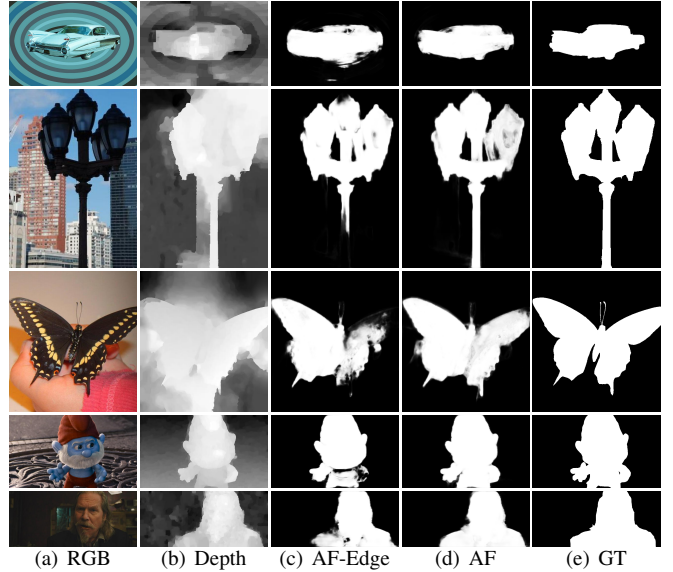


Fig. 3: Comparison of predictions with and without the edge-preserving loss.

The effectiveness of the edge-preserving loss: With the edge-preserving loss, ‘AF’ achieves superior performance to ‘AF-Edge’ as reported in Table I. The results in Fig. 3 illustrate that the saliency maps predicted by ‘AF’ can reduce the blur effect around objects’ boundaries when the objects have similar appearances with the background. In addition, the salient objects are detected more coherently and completely with the edge-preserving constraint. The superiority in both quantitative and qualitative comparisons proves the effectiveness of this loss.

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

We further compare our full model with two traditional methods including GP [14] and LBE [13], together with three CNN-based RGB-D saliency detection networks, including CTMF [17], MPC1 [18] and PCA [19]. The quantitative comparisons are reported in Table II, Fig. 4, and Fig. 5. Qualitative comparisons are demonstrated in Fig. 6.

Quantitative Comparison: As shown in Table II, Fig. 4, and Fig. 5, the proposed method outperforms other state-of-the-art methods in terms of all evaluation metrics. Table II and Fig. 4 show that all deep learning based approaches outperform traditional methods by a great margin; and end-to-end frameworks, including PCA [19] and our approach, are superior to multi-stage methods such as CTMF [17] and MPC1 [18]. Moreover, benefited from our fusion scheme and edge-preserving loss, the proposed method consistently improves the F-measure and MAE achieved by PCA on all three datasets, especially on NLPR where accurate depth data are collected by Kinect. The results indicate that our model can fuse depth information with RGB data more effectively.

Qualitative Comparison: Fig. 6 provides a visual comparison between our model and other approaches. For the typical scenarios that share the similar appearance with the

TABLE II: Comparison of maximum F-measure and MAE.

Methods	NJUD		NLPR		STEREO	
	F_β	MAE	F_β	MAE	F_β	MAE
GP	0.773	0.1679	0.764	0.1108	0.783	0.1564
LBE	0.718	0.2381	0.687	0.2191	0.698	0.2421
CTMF	0.857	0.0847	0.841	0.0554	0.853	0.0849
MPCI	0.868	0.0789	0.841	0.0585	0.861	0.0780
PCA	0.887	0.0592	0.864	0.0433	0.884	0.0592
AF	0.899	0.0534	0.899	0.0327	0.904	0.0462

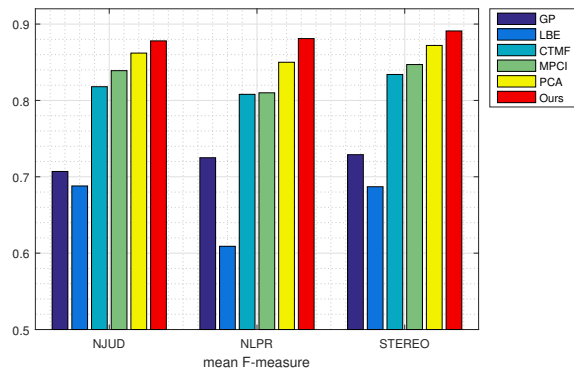


Fig. 4: Comparison of mean F-measure.

background, as shown in the first two rows, the proposed method can better capture effective information in depth data and localize the salient objects accurately. The depth distributions in the third and fourth rows are indistinguishable for the salient objects. Other methods fail to highlight complete and uniform salient objects while our fusion strategy can avoid such depth confusions to a great extent. Moreover, benefited from the edge-preserving loss, the proposed method preserves rich details and sharp boundaries in comparison with the others as demonstrated in the last two rows.

Failed Cases: The proposed approach is capable of detecting salient objects as long as the objects stand out in one modality. When objects are not distinguishable in both modalities, our approach fails as expected. Fig. 7 demonstrates two typical examples. As shown in the figure, such scenarios are challenging to all existing methods.

V. CONCLUSION

In this paper, we have presented a novel end-to-end framework for RGB-D salient object detection. Instead of directly concatenating RGB and depth features or element-wisely multiplying/adding saliency predictions, we introduce a switch map that is adaptively learned to fuse the effective information from RGB and depth predictions. An edge-preserving loss is also designed for correcting blurry boundaries and further improving spatial coherence. The experiments have demonstrated that the proposed method consistently outperforms other state-of-the-art methods on different datasets.

ACKNOWLEDGEMENT

This work was supported in part by Major Scientific Project of Zhejiang Lab (No. 2018DD0ZX01) and the Natural Sci-

ence Foundation of Zhejiang Province, China under Grant LY17F010007.

REFERENCES

- [1] B. Lai and X. Gong, "Saliency guided dictionary learning for weakly-supervised image parsing," in *Proc. IEEE CVPR*, Jun. 2016, pp. 3630–3639.
- [2] S. Bi, G. Li, and Y. Yu, "Person re-identification using multiple experts with random subspaces," *J. Image Graph.*, vol. 2, no. 2, 2014.
- [3] R. G. Cinbis, J. Verbeek, and C. Schmid, "Multi-fold mil training for weakly supervised object localization," in *Proc. IEEE CVPR*, Jun. 2014, pp. 2409–2416.
- [4] Y. Wu, Y. Sui, and G. Wang, "Vision-based real-time aerial object localization and tracking for uav sensing system," *IEEE Access*, vol. 5, pp. 23 969–23 978, 2017.
- [5] T. Judd, K. Ehinger, F. Durand, and A. Torralba, "Learning to predict where humans look," in *Proc. IEEE ICCV*, Sept. 2009, pp. 2106–2113.
- [6] T. Liu, Z. Yuan, J. Sun, J. Wang, N. Zheng, X. Tang, and H.-Y. Shum, "Learning to detect a salient object," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 2, pp. 353–367, Feb. 2011.
- [7] 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, Mar. 2015.
- [8] G. Li and Y. Yu, "Visual saliency based on multiscale deep features," in *Proc. IEEE CVPR*, Jun. 2015, pp. 5455–5463.
- [9] L. Wang, H. Lu, X. Ruan, and M.-H. Yang, "Deep networks for saliency detection via local estimation and global search," in *Proc. IEEE CVPR*, Jun. 2015, pp. 3183–3192.
- [10] G. Li and Y. Yu, "Deep contrast learning for salient object detection," in *Proc. IEEE CVPR*, Jun. 2016, pp. 478–487.
- [11] 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 CVPR*, Jul. 2017, pp. 5300–5309.
- [12] R. Ju, L. Ge, W. Geng, T. Ren, and G. Wu, "Depth saliency based on anisotropic center-surround difference," in *Proc. IEEE ICIP*, Oct. 2014, pp. 1115–1119.
- [13] D. Feng, N. Barnes, S. You, and C. McCarthy, "Local background enclosure for rgb-d salient object detection," in *Proc. IEEE CVPR*, Jun. 2016, pp. 2343–2350.
- [14] J. Ren, X. Gong, L. Yu, W. Zhou, and M. Ying Yang, "Exploiting global priors for rgb-d saliency detection," in *Proc. IEEE CVPRW*, Jun. 2015, pp. 25–32.
- [15] H. Peng, B. Li, W. Xiong, W. Hu, and R. Ji, "Rgb-d salient object detection: a benchmark and algorithms," in *Proc. ECCV*, Sept. 2014, pp. 92–109.
- [16] 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, May 2017.
- [17] 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 T. Cybern.*, Nov. 2017.
- [18] 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.*, vol. 86, pp. 376–385, 2019.
- [19] H. Chen and Y. Li, "Progressively complementarity-aware fusion network for rgb-d salient object detection," in *Proc. IEEE CVPR*, 2018, pp. 3051–3060.
- [20] Y. Niu, Y. Geng, X. Li, and F. Liu, "Leveraging stereopsis for saliency analysis," in *Proc. IEEE CVPR*, Jun. 2012, pp. 454–461.
- [21] W. Zhu, S. Liang, Y. Wei, and J. Sun, "Saliency optimization from robust background detection," in *Proc. IEEE CVPR*, Jun. 2014, pp. 2814–2821.
- [22] N. Liu and J. Han, "Dhsnet: Deep hierarchical saliency network for salient object detection," in *Proc. IEEE CVPR*, Jun. 2016, pp. 678–686.
- [23] P. Zhang, D. Wang, H. Lu, H. Wang, and X. Ruan, "Amulet: Aggregating multi-level convolutional features for salient object detection," in *Proc. IEEE ICCV*, Oct. 2017, pp. 202–211.
- [24] P. Zhang, L. Wang, D. Wang, H. Lu, and C. Shen, "Agile amulet: Real-time salient object detection with contextual attention," *arXiv preprint arXiv:1802.06960*, 2018.
- [25] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [26] X. Wang, H. Ma, X. Chen, and S. You, "Edge preserving and multi-scale contextual neural network for salient object detection," *IEEE Trans. Image Process.*, vol. 27, no. 1, pp. 1057–1149, Jan 2018.

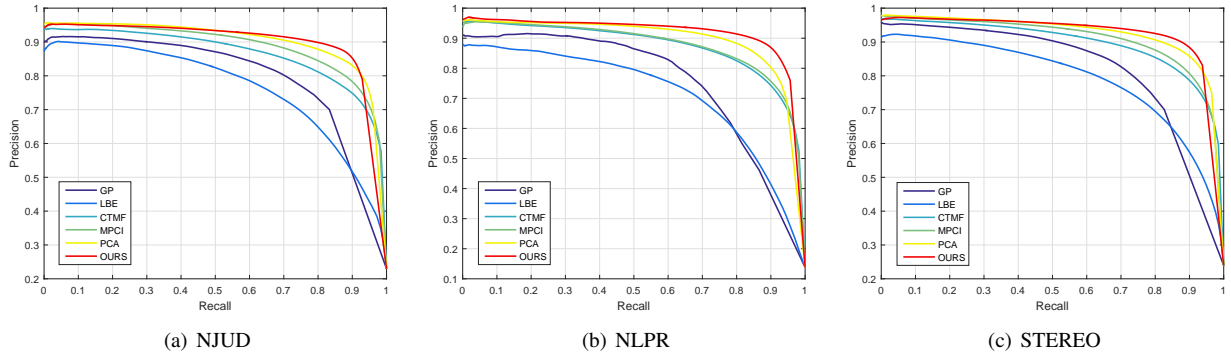


Fig. 5: Comparison of PR curves.

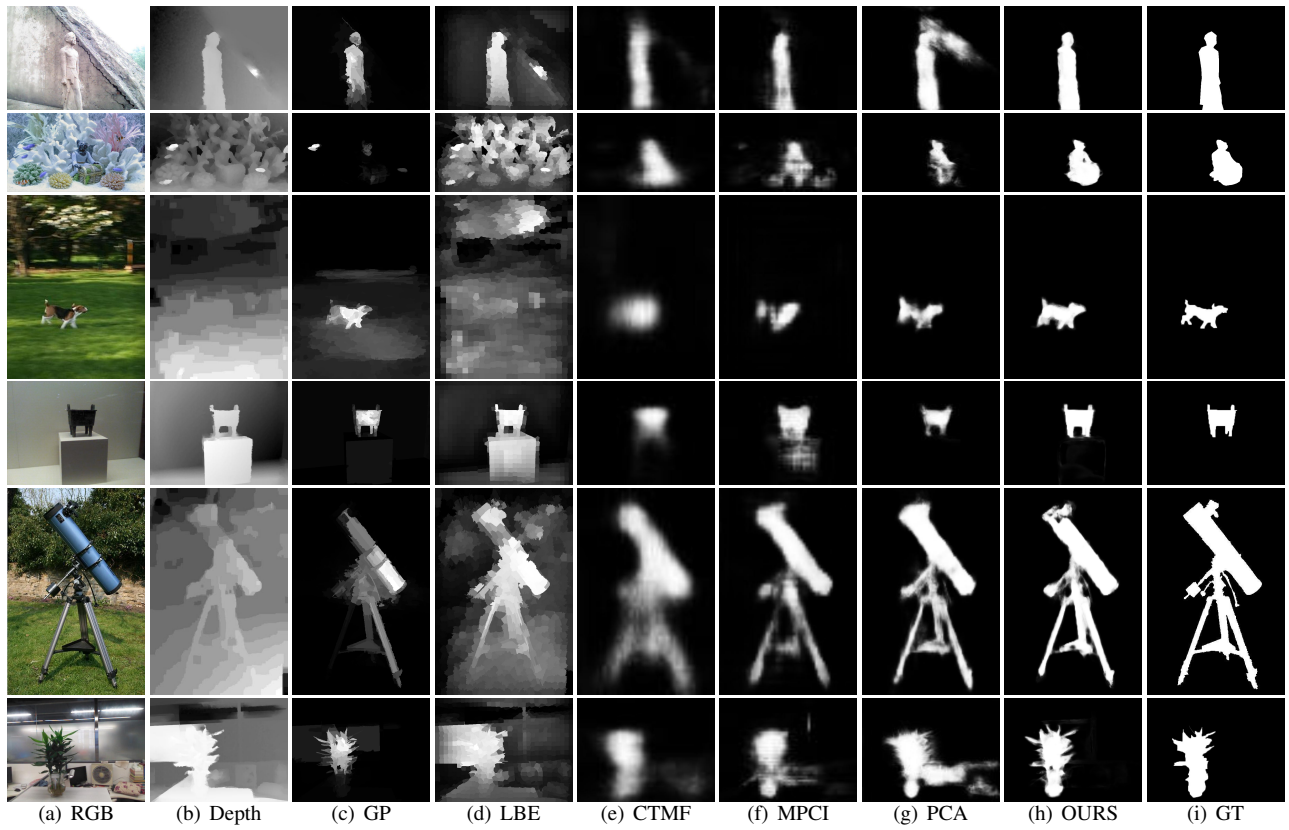


Fig. 6: Visual comparison of saliency maps.

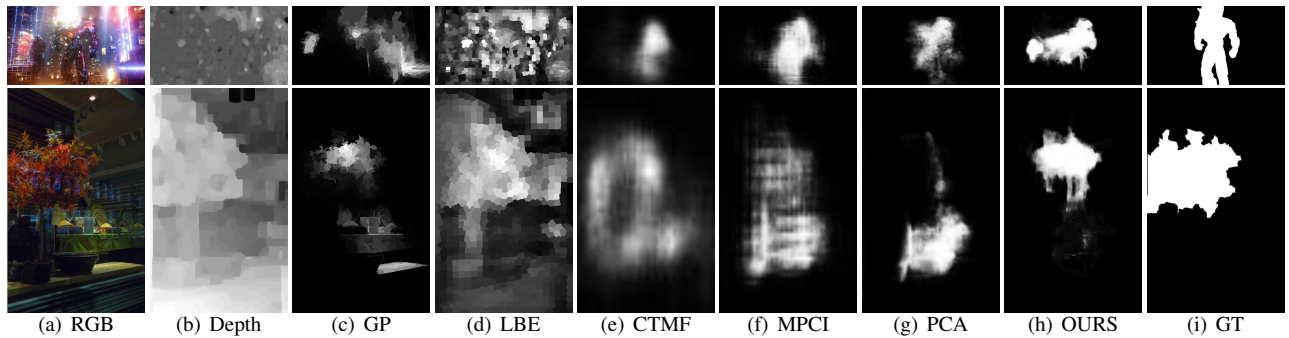


Fig. 7: Failed cases.

- [27] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. Corrado, A. Davis, J. Dean, M. Devin *et al.*, “Tensorflow: Large-scale machine learning on heterogeneous distributed systems,” *arXiv preprint arXiv:1603.04467*, 2016.
- [28] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *Proc. AISTATS*, May 2010, pp. 249–256.
- [29] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.