

# PSNet: Parallel Symmetric Network for Video Salient Object Detection

Runmin Cong, *Member, IEEE*, Weiyu Song, Jianjun Lei, *Senior Member, IEEE*, Guanghui Yue, Yao Zhao, *Senior Member, IEEE*, and Sam Kwong, *Fellow, IEEE*

**Abstract**—For the video salient object detection (VSOD) task, how to excavate the information from the appearance modality and the motion modality has always been a topic of great concern. The two-stream structure, including an RGB appearance stream and an optical flow motion stream, has been widely used as a typical pipeline for VSOD tasks, but the existing methods usually only use motion features to unidirectionally guide appearance features or adaptively but blindly fuse two modality features. However, these methods underperform in diverse scenarios due to the uncomprehensive and unspecific learning schemes. In this paper, following a more secure modeling philosophy, we deeply investigate the importance of appearance modality and motion modality in a more comprehensive way and propose a VSOD network with up and down parallel symmetry, named PSNet. Two parallel branches with different dominant modalities are set to achieve complete video saliency decoding with the cooperation of the Gather Diffusion Reinforcement (GDR) module and Cross-modality Refinement and Complement (CRC) module. Finally, we use the Importance Perception Fusion (IPF) module to fuse the features from two parallel branches according to their different importance in different scenarios. Experiments on four dataset benchmarks demonstrate that our method achieves desirable and competitive performance. The code and results can be found from the link of [https://rmcong.github.io/proj\\_PSNet.html](https://rmcong.github.io/proj_PSNet.html).

**Index Terms**—Salient object detection, Video sequence, Parallel symmetric structure, Importance perception.

## I. INTRODUCTION

VIDEO salient object detection (VSOD) focuses on extracting the most attractive and motion related objects in a video sequence [1], [2], which has been used as a pre-processing step for a wide range of tasks, such as video understanding [3]–[6], video compression [7], video tracking [8], and video caption [9]. Due to the characteristic of video,

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

Weiyu Song and Yao Zhao are with the Institute of Information Science, Beijing Jiaotong University, Beijing 100044, China, and also with the Beijing Key Laboratory of Advanced Information Science and Network Technology, Beijing 100044, China (e-mail: wysong125@bjtu.edu.cn, yzhao@bjtu.edu.cn).

Jianjun Lei is with the School of Electrical and Information Engineering, Tianjin University, Tianjin 300072, China (e-mail: jjlei@tju.edu.cn).

Guanghui Yue is with the National, regional Key Technology Engineering Laboratory for Medical Ultrasound, Guangdong Key Laboratory for Biomedical Measurements and Ultrasound Imaging, School of Biomedical Engineering, Health Science Center, Shenzhen University, Shenzhen 518060, China (email: yueguanghui@szu.edu.cn).

Sam Kwong is with the Department of Computer Science, City University of Hong Kong, Hong Kong SAR, China, and also with the City University of Hong Kong Shenzhen Research Institute, Shenzhen 51800, China (e-mail: cssamk@cityu.edu.hk).

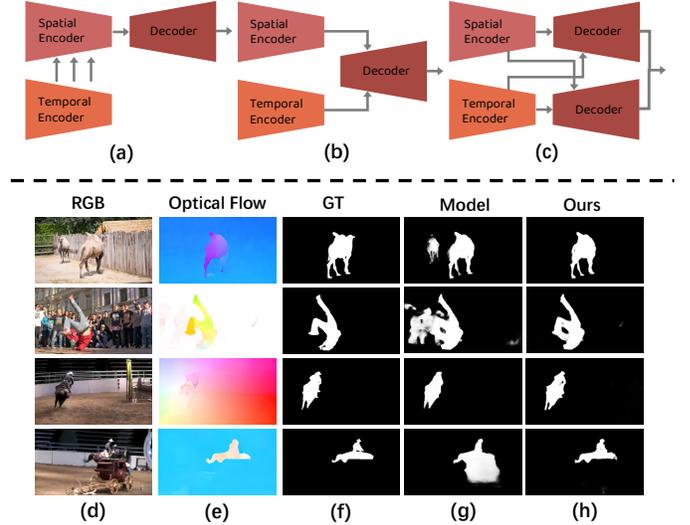


Fig. 1. Top: The structures of VSOD models between our method (c) and the other optical flow-based two-stream VSOD methods (a) (b). Bottom: The saliency results from different models in different scenes. (d) RGB images; (e) Optical flow images; (f) GT; (g) Saliency maps deduced by different methods, where the first row is generated by the MGA method [16], the second row is generated by our baseline model with addition fusion, and the last two rows are generated by the CAG method [18]; (h) Our model.

in addition to the appearance cue, the motion attribute plays an important role, which is different from the SOD task for static images. Entering the deep learning era, a variety of VSOD methods have been explored, which can be roughly divided into two categories, *e.g.*, single-stream methods using the temporal convolution/long short-term memory [10]–[14] and two-stream methods using the optical flow [15]–[18]. Even so, it is still very challenging for current VSOD methods to fully excavate and integrate the information from motion and appearance cues. For the optical flow-based two-stream VSOD model, how to achieve the information interaction according to the role of the two modalities is very important. In this paper, we first rethink and review the interaction mode in the optical flow-based two-stream VSOD structure, and the existing methods can be further categorized into two categories. One is the unidirectional guidance model, as shown in Fig. 1(a), in which the motion information mainly plays a supplementary role. For example, Li *et al.* [16] encouraged motion features to guide the appearance features in the designed VSOD model. As a result, the model pays too much attention to the spatial branch, while the advantage of the motion branch is weakened when dealing with some challenging scenes, such as the stationary

objects with salient appearance may be incorrectly preserved. (see the 1<sup>st</sup> row of Fig. 1). To alleviate the problems mentioned above, the undifferentiated and bidirectional fusion mechanism is proposed as another typical interaction mode, as shown in Fig. 1(b), which no longer distinguishes their primary and secondary roles. Fusing the two modality features by addition or concatenation is the simplest solution, but this way often fails to achieve the desired results, especially for some complex scenes (see the 2<sup>nd</sup> row of Fig. 1). In addition, some works [18] learn the weights to determine the contributions of spatial and temporal features, and then achieve adaptive fusion of two modality features. Although these methods appear to be quite intelligent and achieve relatively competitive performance, this black-box adaptive fusion strategy sometimes only trades off performance rather than maximizing gains when faced with different scenarios. As shown in 3<sup>rd</sup> and 4<sup>th</sup> rows of Fig. 1, they are different frames from different moments of the same video. Although they are similar scenes, the contribution of the two modality data to the final saliency detection is different. We can find that the appearance cues are more important than the motion cues in the 3<sup>rd</sup> row, where the dramatic moving of objects and the change of camera position lead to unclear and blur motion cues. While in the 4<sup>th</sup> row, motion cues can provide more effective guidance information compared with appearance cues that contain some irrepressible noise. According to these observations, when salient objects and backgrounds share similar appearances or background interference is disturbing, interlaced and wrong appearance cues could greatly contaminate the final detection results. But at this time, perhaps accurate motion cues will help us to segment the salient objects correctly. Alternatively, too slow or too fast object motion will blur the estimated optical flow map, thus failing to provide discriminative motion cues and affecting the final detection. In this case, satisfactory detection results can be obtained by exploiting the semantic information from distinctive appearance cues and features. In other words, the roles of the two modalities in different scenes or even similar scenes cannot be generalized, and the uncertainty of the scene makes it very difficult to model interaction fully adaptively. Instead of learning the importance of these two modalities regardless and fully adaptively, we propose a more secure modeling strategy, where the importance of appearance cues and motion cues will be comprehensively and explicitly taken into account to generate the saliency maps, as shown in Fig. 1(c). In our network, we design a top-bottom parallel symmetric structure, which sacrifices the full-automatic intelligence so that we can fuse features more comprehensively, considering the adaptability of the network to different scenarios. Since it struggles for the network to distinguish which modality is more important in one particular scenario, we design two branches with varying tendencies of importance for VSOD, taking one modality feature as a dominant role in each branch and then supplementing from another modality.

Under the parallel symmetric structure, we need to do two things, one is how to realize the utilization of the two modality information in each branch more clearly, and the other is how to integrate the information of the upper and lower

branches to generate the final result. For the first issue, we design the Gather Diffusion Reinforcement (GDR) module and Cross-modality Refinement and Complement (CRC) module to achieve dominate-modality feature reinforcement and cross-modality feature interaction, respectively. Considering that the high-level semantic information can reduce the interference of non-salient information in a single modality and multi-scale information can contribute to more comprehensive features, we design a GDR module to enhance the effectiveness of dominant features in each branch and improve the multi-scale correlation of the dominant features themselves. The outputs of the GDR module are then used for the CRC module in a top-down manner. The key ideas behind the design of the CRC module are as follows. Even if the data from one modality plays a dominant role, there is more or less useful information from the other modality. We divide this role into two types, one is the refinement role, which is mainly used to suppress the irrelevant redundancies in the dominant features, and the other is the complementary role, mainly used to compensate for potential information missing in dominant features. Therefore, we design the CRC module to achieve comprehensive information interaction in the case of explicit primary and secondary relations, which can play the most significant role in our proposed parallel symmetric framework. Although both our upper and lower branches are fully implemented in the VSOD task, the dominant modality they set is different. To obtain more robust and generalized final results, we need to integrate the two branches, which is the second problem we need to solve. Considering the different importance of the upper and lower branches in different scenarios, we introduce an Importance Perception Fusion (IPF) module for adaptive fusion. All designed modules are closely cooperated and integrated under our parallel symmetrical structure to achieve better detection performance. As shown in the 5<sup>th</sup> column of Fig. 1, our model can accurately locate salient objects in different types of scenes, with obvious advantages in detail representation and background suppression. The contributions of this paper can be summarized as:

- Considering the adaptability of the network to different scenarios and the uncertainty of the role of different modalities, we propose a parallel symmetric network (PSNet) for VSOD that simultaneously models the importance of two modality features in an explicit way.
- We propose a GDR module in each branch to perform multi-scale content enhancement for dominant features and design a CRC module to achieve cross-modality interaction, where the auxiliary features are applied to refine and supplement dominant features.
- Experimental results on four mainstream datasets demonstrate that our PSNet outperforms 25 state-of-the-art methods both quantitatively and qualitatively.

## II. RELATED WORK

### A. Salient Object Detection in Single Image and Image Group

For decades, single image-based SOD task has achieved extensive development [19]–[33], and has been widely used in many related fields [2], such as object segmentation [34],

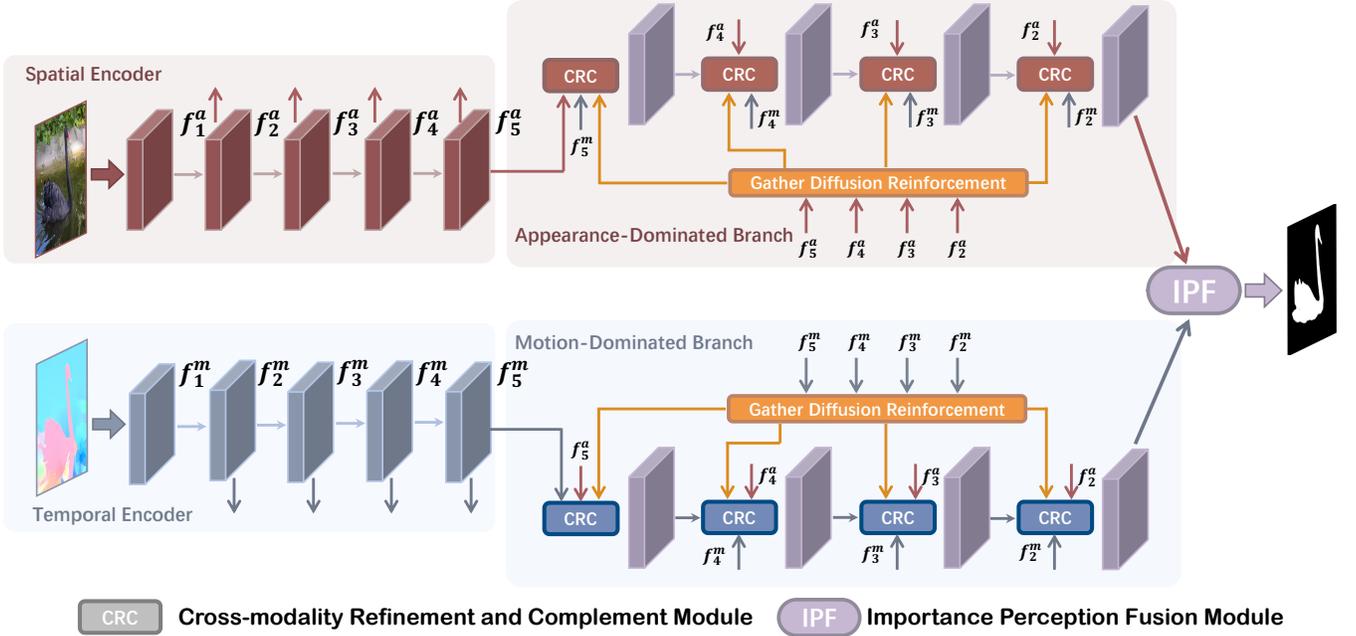


Fig. 2. The flowchart of the proposed Parallel Symmetric Network (PSNet) for video salient object detection. We first extract the multi-level features from RGB images and optical flow maps via spatial encoder and temporal encoder respectively, which are denoted as  $f_i^a$  and  $f_i^m$  ( $i = \{1, 2, \dots, 5\}$ ). Then, the appearance-dominated branch (top branch) and motion-dominated branch (bottom branch) are used to feature decoding. For each decoding, we use Gather Diffusion Reinforcement (GDR) module to perform cross-scale feature enhancement, and then use the Cross-modality Refinement and Complement (CRC) module to achieve cross-modality interaction with an explicit primary and secondary modality relationship. Finally, the Importance Perception Fusion (IPF) module is used to integrate the upper and lower branches by considering their different importance in different scenarios.

content enhancement [35]–[46], and quality assessment [47], [48]. Chen *et al.* [21] developed a method to make full use of global context. Liu *et al.* [22] introduced a network to selectively attend to informative context locations for each pixel. In addition, the salient boundaries have been introduced into the model to improve the representation and highlight the desirable boundaries [23]–[25]. Some methods integrated features in multiple layers of CNN to exploit the context information at different semantic levels [25], [26]. In some challenging and complex single image scenarios, some works seek help from other modality data (*e.g.*, depth map [49]–[55] and thermal map [56]). In addition, co-salient object detection (CoSOD) aims to detect salient objects from an image group containing several relevant images [57]–[66]. The difference between CoSOD and VSOD is that it does not have temporal consistency, and the co-salient object is generally only consistent in semantic categories, rather than the same object.

### B. Salient Object Detection in Video

The last decade has witnessed the considerable development of salient object detection in video sequences. Earlier VSOD methods mostly locate salient objects through hand-crafted features [67]–[70]. Tu *et al.* [67] detected the salient object in the video through two distinctive object detectors and refined the final spatiotemporal saliency result by measuring the foreground connectivity between two maps from two detectors. Chen *et al.* [68] divided the long-term video sequence into some short batches and proposed to detect saliency in a batch-wise way, where the low-rank coherency is introduced to

guarantee temporal smoothness. However, the performance of these methods is not satisfactory due to the limited feature representation capabilities. Recently, deep learning has demonstrated its power in VSOD tasks. Among them, some VSOD models adopt a single-stream structure that directly feeds the video sequences recursively into the network. For instance, Wang *et al.* [10] proposed the first work applying deep learning to the VSOD task. Li *et al.* [71] proposed a two-stage FCN-based model, where the first stage is responsible for detecting static saliency, and the second stage is utilized to detect spatiotemporal saliency with two consecutive frames. In general, this method models saliency in a relatively primitive way. With the development of the model, some subtle module designs are proposed. For example, Song *et al.* [13] used the designed Pyramid Dilated Bidirectional ConvLSTM to achieve deeper spatiotemporal feature extraction. Fan *et al.* [14] introduced a VSOD model based on ConvLSTM, which is applied to model spatiotemporal features in a fixed length of video frames. Moreover, a new VSOD dataset with human visual fixation to model the human saliency shifting is proposed as well. Chen *et al.* [72] focused on the results derived from previous SOTA models, which are applied as pseudo labels to fine-tune a new model, considering the motion quality. Chen *et al.* [12] presented a novel spatiotemporal modeling unit based on 3D convolution.

In addition, another typical VSOD pipeline is the two-stream structure, where the optical flow image generated by FlowNet2 [73] or other methods is directly fed into the network as another stream input. Current two-stream models can be divided into two categories. One is the uni-direction

guidance model, as shown in Fig. 1(a). Li *et al.* [16] used the two-stream model to extract two modality features, where the temporal branch is designed to affect the spatial branch for better salient results. Ren *et al.* [74] proposed to excite the video saliency branch with encoded optical flow features, and developed the semi-curriculum learning manner to learn saliency. The other is the undifferentiated and bidirectional fusion model, as shown in Fig. 1(b). Ji *et al.* [17] proposed to exploit the cross-modality features by considering a mutual restraint scheme. Chen *et al.* [18] tried to adaptively fuse features from motion and appearance via estimating confidence scores. Su *et al.* [15] dynamically learned the weight vector of two modality features and aggregated the corresponding features complementarily. Zhao *et al.* [75] used scribble labels to train the VSOD model, in which cross-modality fusion and temporal constraint are used to model spatiotemporal information.

It is worth mentioning that the VSOD task is highly related to the unsupervised video object segmentation (VOS) task. Lu *et al.* [76] proposed an unsupervised video object segmentation method, where a co-attention layer learns discriminative foreground information in video frame pairs. Zhou *et al.* [77] used an asymmetric motion-attentive transition to identify moving motion information and facilitate the representation of spatiotemporal cues in the zero-shot video object segmentation task. Wang *et al.* [78] built a fully connected graph to explore more representative and high-order relation information for zero-shot VOS. Cho *et al.* [79] regarded motion cues as optional in the unsupervised video object segmentation network, thereby designing a motion branch that can be adaptively turned on or off to participate or not in saliency detection.

The differences between our method and existing methods can be summarized in two major points. The previous two-stream structure is mainly a unidirectional guidance model or undifferentiated and bidirectional fusion model. But they may lead to insufficient information extraction due to wrong selection or ignoring the primary and secondary roles of different modalities (*i.e.*, appearance and motion). Hence, we propose a more comprehensive and more secure strategy for modeling cross-modality interaction in the VSOD task under the two-stream structure, including an appearance-dominated branch and a motion-dominated branch. Two branches each consider the fusion with opposite modality tendencies, owning a clear and specific modality guidance tendency. Furthermore, to implement our overall framework, we also design concrete models that differ from existing methods, where the GDR module, CRC module, and IPF module cooperate to fully mobilize the relationship between different modalities and different detection branches.

### III. METHODOLOGY

#### A. Overview of Proposed Network

As shown in Fig. 2, the proposed PSNet is a two-stream encoder-decoder network, following an up-down mirror-symmetrical structure. For the concise of the following description, we denote the current RGB frame as  $R_t$ , and the next RGB frame as  $R_{t+1}$ . These two adjacent images are input into FlowNet2 [73] to predict optical flow  $O_{t,t+1}$  in

an end-to-end way. With these inputs, the  $R_t$  and  $O_{t,t+1}$  are fed into the pre-trained ResNet50 backbone network that removes the last average pooling layer and the fully connected layer to obtain the encoder features of  $f_i^a$  and  $f_i^m$ , where  $i = \{1, 2, 3, 4, 5\}$  indicates the  $i^{th}$  layer. The parameters of the spatial encoder and temporal encoder are not shared in our model. In this network, we only use the features from the last four layers for the savings of computational costs. After that, both  $f_i^a$  and  $f_i^m$  are further input to the appearance-dominated branch and motion-dominated branch for feature decoding. As for the feature decoding process, we briefly illustrate the appearance-dominated branch as an example. First, all the dominant encoder features  $f_i^a$  from the last four layers (*i.e.*,  $i = \{1, 2, 3, 4, 5\}$ ) are embedded into the GDR module to achieve the dominate-modality feature reinforcement and generate the corresponding reinforced dominant features  $f_i^{a,r}$ . Following that, the reinforced dominant features  $f_i^{a,r}$ , the corresponding appearance and motion features of  $f_i^a$  and  $f_i^m$ , and the previous decoder features  $f_{i+1}^{a,d}$  are input to the CRC module, thereby completing the explicit cross-modality information interaction and obtaining the decoder features  $f_i^{a,d}$  of the current layer. Finally, we aggregate the outputs of the two decoder branches and generate the final saliency map through the IPF module.

#### B. Gather Diffusion Reinforcement Module

As mentioned earlier, each of our decoding branches has clear dominant and auxiliary modality partitions. In order to ensure the effectiveness and comprehensiveness of the dominant modality features as much as possible, we consider the following motivations for designing a GDR module to strengthen the dominant features of each layer. In the encoding stage, the features extracted by each layer are relatively independent and have their own characteristics. With the network going deeper, the high-level features may contain more location and abstract semantic information about the salient object. At the same time, the low-level features are prone to have more detailed information, such as textures and boundaries. Both high-level and low-level features are essential for salient object detection and integrating them can help to generate high-quality multi-level features. Based on this, the primary function of our GDR module is to correlate the relationship between encoder features at different scales to develop more comprehensive encoder features. The detailed architecture of the GDR module is shown in Fig. 3.

Given several features from different levels of the dominant branch, a Gather module is designed to exploit cross-layer and cross-scale information interaction. Specifically, considering that the features of different layers may contain some noise, especially in the low-level features, the coarse semantic mask predicted by the top encoder layer is used to filter out such noise, which is defined as:

$$mask_5 = \sigma(\mathcal{C}_{3 \times 3}(\mathcal{C}_{3 \times 3}(f_5))), \quad (1)$$

$$f_i^s = Up(mask_5) \otimes f_i, \quad (2)$$

where  $\mathcal{C}_{3 \times 3}$  is convolution layer with the kernel size of  $3 \times 3$ ,  $\sigma$  denotes sigmoid function,  $\otimes$  refers to element-wise

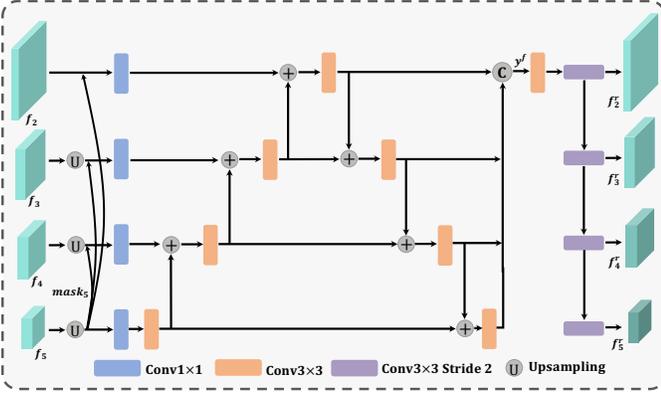


Fig. 3. The architecture of our proposed Gather Diffusion Reinforcement (GDR) module.

multiplication operation,  $f_i^s$  are the features after the semantic filtering, and  $Up$  is the upsampling operation. For simplicity, the superscript  $a$  or  $m$  of the encoder features  $f_i$  indicating the appearance branch or the motion branch is omitted.

Subsequently, inspired by [80], we use a recursive bidirectional structure to fuse multi-level features and establish relationships between features across scales in a hierarchical manner. That is, the fusion process is not limited to top-down but also explores bottom-up fusion to achieve more comprehensive multi-scale fusion. The top-down multi-scale feature interaction can be described as:

$$y_i = \begin{cases} \mathcal{C}_{3 \times 3}(\mathcal{C}_{1 \times 1}(f_i^s) + y_{i+1}), & i = \{2, 3, 4\} \\ \mathcal{C}_{3 \times 3}(\mathcal{C}_{1 \times 1}(f_i^s)), & i = 5 \end{cases} \quad (3)$$

where  $\mathcal{C}_{1 \times 1}$  is convolution layer with the kernel size of  $1 \times 1$ .

Then, the reverse operation is also performed to achieve a more comprehensive cross-scale interaction:

$$y'_i = \begin{cases} \mathcal{C}_{3 \times 3}(y_i + y'_{i-1}), & i = \{3, 4, 5\} \\ \mathcal{C}_{3 \times 3}(y_i), & i = 2 \end{cases} \quad (4)$$

As such, all interaction features  $\{y'_i \mid i = \{2, 3, 4, 5\}\}$  are fused together in the form of concatenation-convolution:

$$y^f = \mathcal{C}_{3 \times 3}(\text{Cat}[y'_2, y'_3, y'_4, y'_5]), \quad (5)$$

where  $\text{Cat}$  is channel-wise concatenation operation.

Finally, considering that each level of CRC needs a different scale of features from GDR, a straightforward way is to use a diffusion module to diffuse features. Here, we perform  $3 \times 3$  convolution with stride 2 on the fusion features and generate the reinforced features:

$$f_i^r = \begin{cases} \mathcal{C}_{3 \times 3 \text{ stride } 2}(y^f), & i = 2 \\ \mathcal{C}_{3 \times 3 \text{ stride } 2}(f_{i-1}^r), & i = \{3, 4, 5\} \end{cases} \quad (6)$$

where  $\mathcal{C}_{3 \times 3 \text{ stride } 2}$  denotes  $3 \times 3$  convolution with the stride of 2. The reinforced features in the appearance-dominated branch and motion-dominated branch can be distinguished as  $f_i^{a,r}$  and  $f_i^{m,r}$ , respectively. In fact, both DSS [81] and our GDR module adopt a structure similar to FPN [82], which is used for enriching the representation of multi-scale information. But our GDR module acts as a single-modality feature enhancement with cross-level, cross-scale information, and then passes them to the decoder. In the implementation,

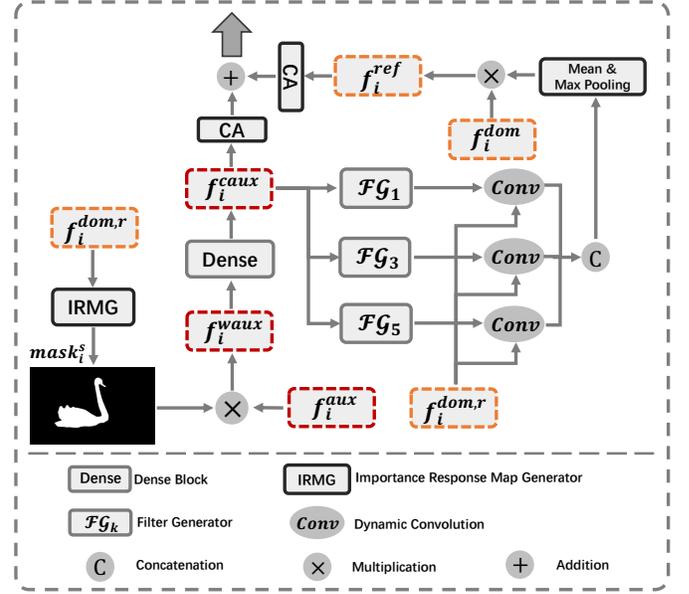


Fig. 4. The structure of our proposed Cross-modality Refinement and Complement (CRC) module.

for the FPN and short connection in [81], the multi-scale information is fused in a single direction (up-to-down). While for our GDR module, the interaction direction is not restricted to a single direction but follows a recursive way to achieve more comprehensive multi-scale information. In addition, the high-level features from the encoder are used to filter out the noise in low-level features for a more robust single-modality representation in our GDR module.

### C. Cross-modality Refinement and Complement Module

The cross-modality interaction has always been a core topic in video salient object detection. Fortunately, under our parallel symmetric architecture, each branch has a definite dominant modality and a corresponding auxiliary modality, so that we can design the interaction module more explicitly and clearly. For concise, we denote the dominant modality features as  $f_i^{dom}$ , and the auxiliary features as  $f_i^{aux}$ . The structure of CRC is depicted in Fig. 4. As mentioned earlier, for each dominated branch, the  $f_i^{dom}$  are more dominant than  $f_i^{aux}$ , but this does not mean that auxiliary features are entirely useless. In other words, there is still some helpful information in  $f_i^{aux}$ , which will contribute to the saliency feature learning. Therefore, starting from the auxiliary modality, we divide its role into two types of refinement and complement and design a CRC module to maximize the use of auxiliary information. Furthermore, the features of  $f_i^{dom}$  are fed into the GDR module to generate the reinforced dominant features  $f_i^{dom,r}$  for the current CRC module.

On the one hand, the supplement of the feature dimension is the most direct from the perspective of information interaction. But direct and indiscriminate integration of  $f_i^{aux}$  may introduce contamination noise into the dominant features. Therefore, we try to select auxiliary features from the perspective of dominant features and determine the feature components

that need to be supplemented. Specifically, two convolutions are employed on the  $f_i^{dom,r}$  to obtain the importance response map  $mask_i^s$ . Then, we use this map to weight the auxiliary features of  $f_i^{aux}$  and utilize the dense block and residual block to strengthen the features, thereby determining the auxiliary features that need to be supplemented to the reinforced dominant features. The above operations are presented as follows:

$$mask_i^s = \sigma \left( \mathcal{C}_{1 \times 1} \left( \mathcal{C}_{3 \times 3} \left( f_i^{dom,r} \right) \right) \right), \quad (7)$$

$$f_i^{waux} = mask_i^s \otimes f_i^{aux}, \quad (8)$$

$$f_i^{caux} = \mathcal{C}_{1 \times 1} \left( Dense \left( f_i^{waux} \right) + f_i^{waux} \right), \quad (9)$$

where  $f_i^{waux}$  are the weighted auxiliary features,  $f_i^{caux}$  denote the final complemented auxiliary features, and  $Dense$  represents the dense block in [83].

On the other hand, in addition to the complement of feature dimensions, auxiliary features can also be used to refine the irrelevant interference and misinformation in the dominant features. However, considering the difference and interference noise of the two modalities, as well as the large variation of the characteristics of the two modalities with the scene, we do not directly apply  $f_i^{caux}$  to refine  $f_i^{dom,r}$ , but introduce the dynamic convolution filters [84], [85] to adaptively generate the convolution kernel for different scenarios, so as to ensure the generalization and robustness of the network. With the complemented auxiliary features  $f_i^{caux}$ , we use the local dynamic convolution [84] with different dilated rates to generate multi-scale convolution kernels. Then, the generated dynamic convolution kernels are used to convolve the reinforced dominant features  $f_i^{dom,r}$ , achieving the goal of refining the details. The process can be expressed by:

$$k_1 = \mathcal{F}\mathcal{G}_1 \left( f_i^{caux} \right) \otimes f_i^{dom,r}, \quad (10)$$

$$k_3 = \mathcal{F}\mathcal{G}_3 \left( f_i^{caux} \right) \otimes f_i^{dom,r}, \quad (11)$$

$$k_5 = \mathcal{F}\mathcal{G}_5 \left( f_i^{caux} \right) \otimes f_i^{dom,r}, \quad (12)$$

$$f_i^{dy} = \mathcal{C}_{3 \times 3} \left( Cat \left[ k_1, k_3, k_5 \right] \right), \quad (13)$$

where  $\mathcal{F}\mathcal{G}_j$  presents the filter generator with the dilated rate of  $j$  by using two convolutions and reshaping operations, and  $\otimes$  indicates convolution operation.

Next, we employ  $f_i^{dy}$  to generate a refinement mask, and then revise and refine the features of  $f_i^{dom}$ :

$$C_i = Cat \left( maxpool \left( f_i^{dy} \right), avgpool \left( f_i^{dy} \right) \right), \quad (14)$$

$$mask_i^r = \sigma \left( \mathcal{C}_{3 \times 3} \left( C_i \right) \right), \quad (15)$$

$$f_i^{ref} = mask_i^r \otimes f_i^{dom}, \quad (16)$$

where  $f_i^{ref}$  denote the refined dominant features,  $mask_i^r$  is the generated refinement mask, and  $maxpool$  and  $avgpool$  denote max-pooling and average pooling respectively. Finally, we combine the complemented auxiliary features  $f_i^{caux}$  and refined dominant features  $f_i^{ref}$  after the channel compaction:

$$f_i^{rc} = CA \left( f_i^{caux} \right) + CA \left( f_i^{ref} \right), \quad (17)$$

where  $CA$  presents channel attention block [86]–[88].

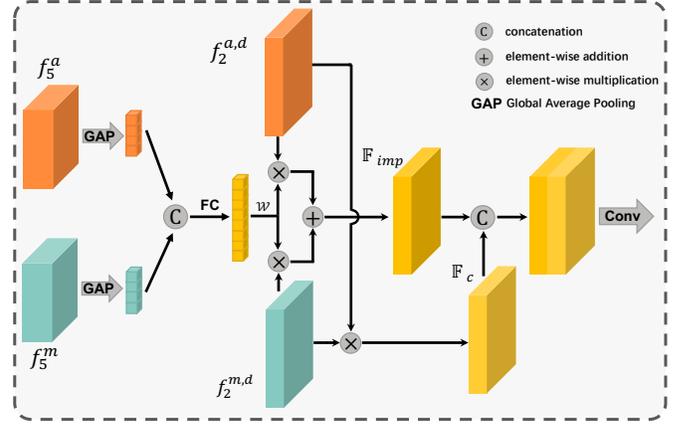


Fig. 5. The structure of our proposed Importance Perception Fusion (IPF) module.

With the final interaction features  $f_i^{rc}$ , we combine them with the decoder features of the previous layer  $f_{i+1}^{dom,d}$  to generate the final decoder features of the current layer:

$$f_i^{dom,d} = \begin{cases} \mathcal{C}_{3 \times 3} \left( Cat \left( f_i^{rc}, Up \left( f_{i+1}^{dom,d} \right) \right) \right), & i = \{2, 3, 4\} \\ \mathcal{C}_{3 \times 3} \left( Cat \left( f_i^{rc}, f_i^{dom,r} \right) \right), & i = 5 \end{cases} \quad (18)$$

where the decoder features of  $i^{th}$  layer in the appearance-dominated branch and motion-dominated branch are as  $f_i^{a,d}$  and  $f_i^{m,d}$ , respectively.

#### D. Importance Perception Fusion Module

Under the parallel symmetric framework, the appearance-dominated branch and motion-dominated branch generate comprehensive spatiotemporal features corresponding to different modality dominance with clear and well-defined roles. Although both branches can be regarded as complete VSOD branches, the dominant modality they set is different, and the learned features are also different. To obtain a more robust and generalized final result, inspired by the MF module in [89], we introduce an IPF module to achieve the branch fusion, considering the different importance of the upper and lower branches in different scenarios. Fig. 5 illustrates the flowchart of the IPF module. The input features of the IPF module can be divided into two parts. One is the last-layer features of the encoder ( $f_5$ ), which is used to perceive the different importance of the upper and lower branches in different scenes. In this way, an adaptive importance weight  $\mathcal{W} \in \mathbb{R}^{128}$  is learned. More specifically, the features from the 5<sup>th</sup> layer of appearance and motion encoders (i.e.,  $f_5^a$  and  $f_5^m$ ) are first fed into global average pooling and then concatenated together to learn a channel-wise weight:

$$\mathcal{W} = \sigma \left( FC \left( Cat \left( GAP \left( f_5^a \right), GAP \left( f_5^m \right) \right) \right) \right), \quad (19)$$

where  $GAP$  is global average pooling, and  $FC$  represents a full-connected layer.

Another input to the IPF module is the top-layer features of two decoders. The reason why we choose top-layer features  $f_2^d$  of the decoder is that the  $f_2^d$  has higher resolution

than  $f_3^d$  and  $f_4^d$ , and contains more comprehensive decoding information, which is more suitable for our purpose in IPF. More specifically, we use the weight  $\mathcal{W}$  to combine the output decoder features of two branches  $f_2^{a,d}$  and  $f_2^{m,d}$  into importance weighted features  $\mathbb{F}_{imp}$ :

$$\mathbb{F}_{imp} = \mathcal{W} \odot f_2^{a,d} + (1 - \mathcal{W}) \odot f_2^{m,d}, \quad (20)$$

where  $f_2^{a,d}$  and  $f_2^{m,d}$  are the decoder features of last layer in the corresponding branch, and  $\odot$  denotes element-wise multiplication with the broadcasting strategy.

Furthermore, the common response between two outputs from two branches is also important for the final saliency result. Thus, a simple but effective way is to perform multiplication to highlight the common part of the two branches:

$$\mathbb{F}_c = f_2^{a,d} \otimes f_2^{m,d}. \quad (21)$$

Finally, the common features  $\mathbb{F}_c$  and importance weighted features  $\mathbb{F}_{imp}$  are combined by concatenation operation to predict the final saliency map:

$$pre_s = \sigma(\mathcal{C}_{3 \times 3}(\mathcal{C}_{3 \times 3}(\text{Cat}(\mathbb{F}_c, \mathbb{F}_{imp}))))). \quad (22)$$

### E. Loss Function

The network is trained in a multiple supervision manner for the sake of faster convergence and better performance. First, for the final saliency results generated by the IPF module, we employ a joint loss function  $\mathcal{L}_{sal}$  to train our model, which is given by:

$$\mathcal{L}_{sal} = \mathcal{L}_{bce}(pre_s, GT) + \mathcal{L}_{ssim}(pre_s, GT), \quad (23)$$

where  $\mathcal{L}_{bce}$  is the binary cross-entropy loss, and  $\mathcal{L}_{ssim}$  is the structural similarity loss.

In addition, we add the side-output supervision on each branch. Taking the appearance-dominated branch as an example, firstly, the saliency map  $S_a$  predicted by the appearance-dominated branch will be trained by the  $\mathcal{L}_{sal}$  loss. Besides, the intermediate saliency results from GDR and CRC modules are also supervised by the ground truth. Specifically, the supervisions include: (1) the appearance backbone saliency map  $mask_5$  deduced from the 5<sup>th</sup> layer of appearance backbone, which is employed in GDR as a noise filter; (2) the importance response map  $mask_i^s$  deduced from each CRC in the appearance-dominated branch, in which  $i = \{2, 3, 4, 5\}$ . Therefore, for the appearance-dominated branch, the loss function can be formulated as:

$$\begin{aligned} \mathcal{L}_{appearance} = & \mathcal{L}_{sal}(S_a, GT) + \lambda_1 \mathcal{L}_{bce}(mask_5, GT) \\ & + \lambda_2 \sum_{i=2}^5 \mathcal{L}_{bce}(mask_i^s, GT) \end{aligned}, \quad (24)$$

where  $\lambda_1$  and  $\lambda_2$  are hyper-parameters for balancing the losses, which are empirically set to 0.6 and 0.4, respectively. To fit the size of the predicted map, all ground truth will be downsampled to the size of the corresponding predicted map. Similarly, we can get the loss function of the motion-dominated branch, denoted as  $\mathcal{L}_{motion}$ .

Finally, the total loss is defined as follows:

$$\mathcal{L}_{total} = \mathcal{L}_{sal} + \mathcal{L}_{appearance} + \mathcal{L}_{motion}. \quad (25)$$

## IV. EXPERIMENTS

### A. Datasets and Evaluation Metrics

We conduct experiments on four widely used public VSOD datasets in order to fully evaluate the effectiveness of our proposed method, *i.e.*, DAVIS [101], DAVSOD [14], SegV2 [102], and ViSal [103]. **DAVIS** [101] dataset consists of 50 clips of 480p and 720p videos with high-quality dense annotations, which is further split into 30 videos for training and 20 videos for testing. **DAVSOD** [14] dataset includes 226 clips of densely annotated videos, where the salient objects are annotated by dynamic eye-tracking. In DAVSOD, 80 clips of videos are for testing. **SegV2** [102] dataset is an early proposed dataset with 14 videos and 1065 annotated frames, including multiple objects that make it more challenging than others. **ViSal** [103] is a dataset containing 19 videos with 193 pixel-wise annotated frames. In this paper, the test will be carried out on the whole datasets of ViSal and SegV2, and the testing subset of DAVIS and DAVSOD datasets. To quantitatively evaluate the effectiveness of the proposed method, we introduce three evaluation metrics, including maximum F-measure ( $F_\beta$ ) [104], S-measure ( $S_m$ ) [105], and Mean Absolute Error (MAE) [106]. For these three metrics, except for the MAE score, larger values of the F-measure and S-measure indicate better performance.

### B. Implementation Details

We use the Pytorch toolbox to implement our network and train our model with an NVIDIA GTX3090 GPU. We also implement our network by using the MindSpore Lite tool<sup>1</sup>. Following the setting in [17], we use the stage-wise training protocol with image saliency datasets and video saliency datasets to train our model. In the first stage, we initialize our spatial backbone with a ResNet-50 [107]. Following [17], we remove the CRC and IPF modules, and pre-train this model on the training set of the DUTS dataset [108]. In this stage, the batch size and initial learning rate are set to 16 and 0.002, respectively. Moreover, the learning rate decays 0.1 times per 10 epochs. In the second stage, we use FlowNet2 [73] to generate the corresponding optical flow map for each frame of the DAVIS dataset and pre-train the temporal branch. The training settings are the same as stage 1. Next, in stage 3, we use the DAVIS dataset, including RGB images and optical flow maps, to fine-tune our whole PSNet. Concretely, we load the learned weights from stage 1 and stage 2 to the spatial branch and temporal branch, respectively. The number of batch sizes is set to 8. The learning rate is set to 0.0002 for finer learning and stops learning after 20 epochs. In each stage, we use the stochastic gradient descent (SGD) optimizer to train our model with a momentum of 0.9 and a weight decay of 0.0005. We resize all input images to  $384 \times 384$ . Furthermore, we apply a multi-scale training strategy with scales of  $\{0.75, 1, 1.25\}$ , random horizontal flipping, and random vertical flipping to enhance the generalizability and stability of our trained model.

<sup>1</sup><https://www.mindspore.cn/>

TABLE I  
 QUANTITATIVE RESULTS ON THE DAVIS, SEG V2, DAVSOD, AND VISAL DATASETS. THE TOP TWO SCORE WAS MARKED IN **BOLD** AND UNDERLINE, RESPECTIVELY.

Methods	Years	DAVIS			SegV2			ViSal			DAVSOD		
		$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$
<b>Deep Learning Static Saliency Object Detection</b>													
EGNet [23]	2019	0.768	0.829	0.056	0.774	0.845	0.024	0.941	0.946	0.015	0.604	0.719	0.101
CPD [90]	2019	0.778	0.859	0.032	0.778	0.841	0.023	0.941	0.942	0.016	0.608	0.724	0.092
ITSD [91]	2020	0.835	0.876	0.033	0.807	0.787	0.027	–	–	–	0.651	0.747	0.094
<b>Traditional Video Saliency Object Detection</b>													
MSTM [92]	2016	0.429	0.583	0.165	0.526	0.643	0.114	0.673	0.749	0.095	0.344	0.532	0.211
SGSP [93]	2017	0.655	0.592	0.138	0.673	0.681	0.124	0.677	0.706	0.165	0.426	0.577	0.207
STBP [94]	2016	0.544	0.677	0.096	0.640	0.735	0.061	0.622	0.629	0.163	0.410	0.568	0.160
FDOS [67]	2017	0.701	0.784	0.061	0.683	0.765	0.045	0.767	0.801	0.063	0.456	0.582	0.157
SCOM [95]	2018	0.783	0.832	0.048	0.764	0.815	0.030	0.831	0.762	0.122	0.464	0.599	0.220
SFLR [68]	2017	0.727	0.790	0.056	0.745	0.804	0.037	0.779	0.814	0.062	0.478	0.624	0.132
<b>Deep Learning Video Saliency Object Detection</b>													
SCNN [96]	2018	0.714	0.783	0.064	–	–	–	0.831	0.847	0.071	0.532	0.674	0.128
DLVS [10]	2018	0.708	0.794	0.061	–	–	–	0.852	0.881	0.048	0.521	0.657	0.129
FGRN [71]	2018	0.783	0.838	0.043	–	–	–	0.848	0.861	0.045	0.573	0.693	0.098
MBNM [97]	2018	0.861	0.887	0.031	0.716	0.809	0.026	0.883	0.898	0.020	0.520	0.637	0.159
PDB [13]	2018	0.861	0.887	0.028	0.800	0.864	0.024	0.888	0.907	0.032	0.572	0.698	0.116
RCR [98]	2019	0.855	0.882	0.027	0.781	0.842	0.035	0.906	0.922	0.026	0.653	0.741	0.087
SSAV [14]	2019	0.861	0.893	0.026	0.801	0.851	0.023	0.939	0.943	0.020	0.603	0.724	0.092
MGA [16]	2019	0.892	0.910	0.023	0.821	0.865	0.030	0.933	0.936	0.017	0.640	0.738	0.084
PCSA [11]	2020	0.880	0.902	0.022	0.810	0.865	0.025	0.940	0.946	0.017	0.655	0.741	0.086
CASNet [99]	2020	0.860	0.873	0.032	0.847	0.820	0.029	–	–	–	–	–	–
DSNet [15]	2020	0.891	0.914	<u>0.018</u>	0.832	0.875	0.028	0.950	0.949	<u>0.013</u>	–	–	–
STVS [12]	2021	0.865	0.892	<u>0.018</u>	<b>0.860</b>	<b>0.891</b>	<u>0.017</u>	<u>0.952</u>	<u>0.952</u>	<u>0.013</u>	0.651	0.746	0.086
WVSOD [75]	2021	0.793	0.846	0.038	0.762	0.819	0.033	0.875	0.883	0.035	0.593	0.694	0.115
TransVOS [100]	2021	0.869	0.885	<u>0.018</u>	0.800	0.816	0.024	0.928	0.917	0.021	–	–	–
CAG [18]	2021	0.898	0.906	<u>0.018</u>	0.826	0.865	0.027	0.950	0.950	<u>0.013</u>	0.670	0.762	<b>0.072</b>
FSNet [17]	2021	<b>0.907</b>	<b>0.920</b>	0.020	0.805	0.870	0.024	–	–	–	<b>0.685</b>	<b>0.773</b>	<b>0.072</b>
PSNet	–	<b>0.907</b>	<u>0.919</u>	<b>0.016</b>	<u>0.852</u>	<u>0.889</u>	<b>0.016</b>	<b>0.955</b>	<b>0.954</b>	<b>0.012</b>	<u>0.678</u>	<u>0.765</u>	0.074

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

Our proposed method is compared with 25 state-of-the-art methods, including three static SOD methods (EGNet [23], CPD [90], ITSD [91]), six traditional VSOD methods (MSTM [92], STBP [94], SGSP [93], SCOM [95], FDOS [67], SFLR [68]), and sixteen deep learning-based VSOD methods (SCNN [96], DLVS [10], FGRN [71], MBNM [97], PDB [13], RCR [98], SSAV [14], MGA [16], PCSA [11], CASNet [99], WVSOD [75], DSNet [15], TransVOS<sup>2</sup> [100], STVS [12], CAG [18], FSNet [17]). For fair comparisons, all the saliency maps are provided by authors or tested by the released code under the default parameters.

1) **Quantitative Evaluation:** The quantitative results max-F, S-measure, and MAE results are listed in Table I. The static SOD methods perform well in some simple datasets (*e.g.*, ViSal) and even outperform VSOD methods, mainly because the image appearance cues in these datasets dominate most scenes. Quantitatively, the static SOD method CPD [90] wins the percentage gain of 0.9% in  $F_\beta$  and 31.3% in MAE against the MGA [16] method on the ViSal dataset. However, this advantage will no longer exist in the face of complex video scenes, such as the DAVSOD dataset, whose performance is far lower than the VSOD methods. For traditional VSOD methods, due to the limitation of only using hand-crafted

features, their performance is even lower than that of static SOD methods. Taking the best traditional VSOD method SCOM [95] as an example, it achieves comparable performance with some deep learning-based static SOD methods (*e.g.*, EGNet [23] and CPD [90]) on the DAVIS dataset. However, its performance is 50% lower than deep learning-based static SOD methods on the DAVSOD dataset. Seeing Table I, it is observed that our method achieves competitive performance on these four datasets, basically ranking in the top two. Specifically, our method outperforms all other models on the ViSal dataset, which achieves the percentage gain of 7.6% in terms of MAE score compared with the **second best** method (*i.e.*, CAG [18]). In addition, compared with the **second best** model on the DAVIS dataset, the percentage gain reaches 11.1% for the MAE score. Our method is slightly inferior to the FSNet method [17] on the DAVSOD dataset, but achieves comparable performance on the DAVIS dataset, and has a clear performance advantage on the SegV2 dataset. From the analysis of the model size, the size of our method (67.9 M) is only about 80% of the size of FSNet (83.4 M). Overall, our method still has certain advantages in terms of performance and model size.

The training time of PSNet is about 20 hours for all three stages of training, and the testing speed of our PSNet reaches 19 FPS with the model size of 67.9 M. Compared with optical-flow-based methods, such as the MGA (47 FPS and 91.5

<sup>2</sup>TransVOS is a semi-supervised VOS method.

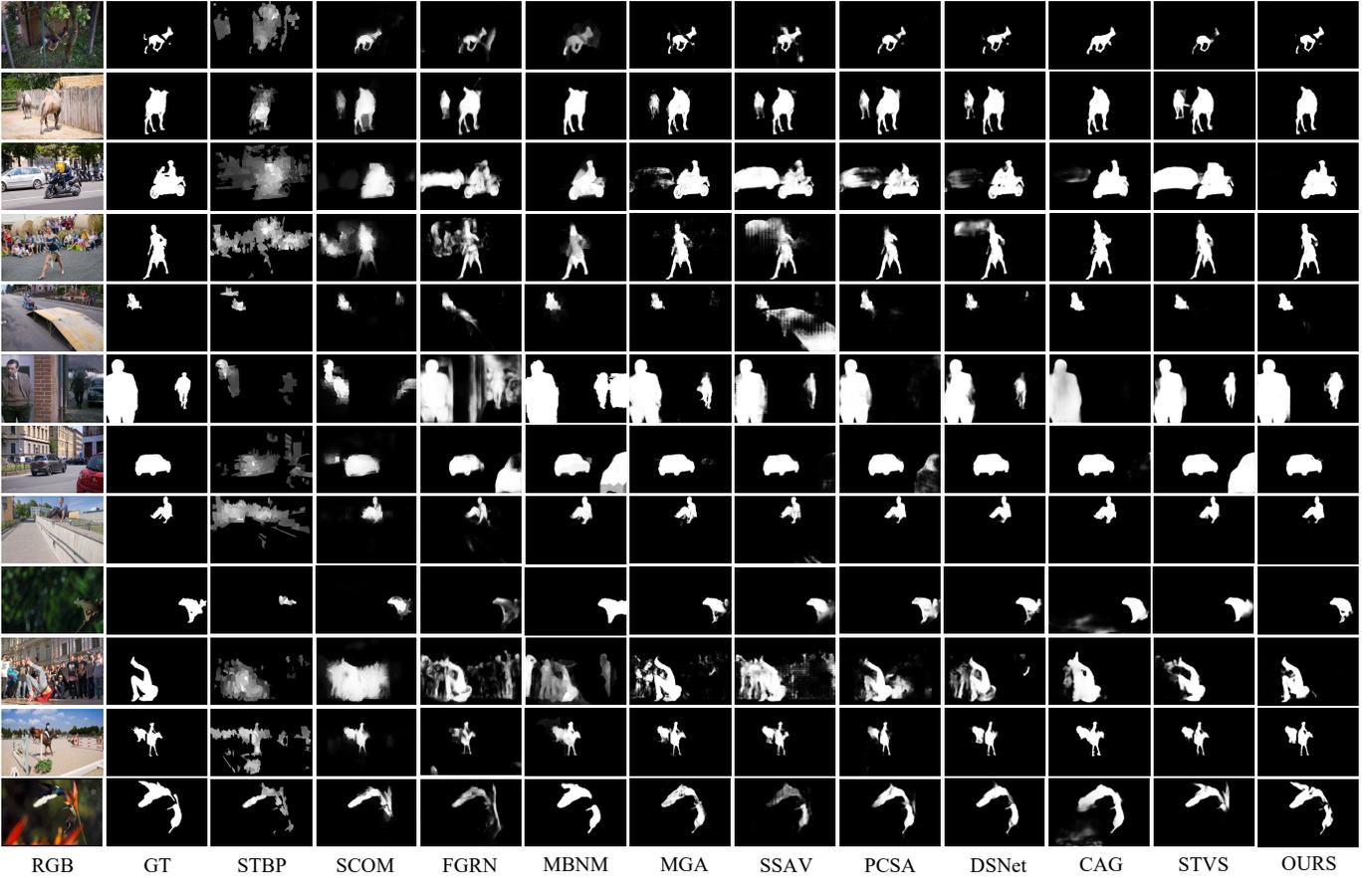


Fig. 6. The visualization results of different video salient object detection methods.

M) [16] and CAG (29 FPS and 55.3 M) [18], although the performance of our algorithm achieves the best result, our testing time and model size are not optimal, which is related to the inclusion of operations such as dynamic convolution in our network design. That is, at present, our model is still far from the real-time effect. Therefore, in the future, we can consider lightweight alternative modules to further improve the efficiency of model testing.

2) **Qualitative Evaluation:** To further illustrate the advantages of our proposed method, we provide some qualitative saliency results in Fig. 6. Compared with other methods, our method achieves superior results with complete object structure, precise saliency location, and sharp boundaries. As can be seen, the traditional VSOD methods cannot achieve desirable results due to their limitations and deficiencies in feature representation, such as the 3<sup>rd</sup> and 4<sup>th</sup> columns. By contrast, deep learning-based methods achieve more competitive results, especially our proposed method is capable of addressing scenes with small objects and complex backgrounds. Taking the 2<sup>nd</sup> and 3<sup>rd</sup> rows as an example, these two sequences contain some tough challenges, in which the background exists some moving but non-salient objects. However, most methods, such as STVS [12] and PCSA [11], cannot completely suppress the distracting background objects in such complicated scenes. Thanks to the design of our network, our model can completely suppress such background disturbances that consider

motion modality as the dominant feature in such scenes and reduce the interference of wrong appearance cues. Meanwhile, in the 10<sup>th</sup> row, the scene is more complex, where the man dancing in front of the audience is our salient object. However, the motion of the foreground objects changes very quickly, and the audience gathered in the back will not only form a relatively strong disturbance in appearance but also have a certain movement of their own, which will undoubtedly make things worse. Therefore, basically all comparison algorithms struggle to handle this scene well, especially the background areas. By contrast, our method can more fully exploit the roles of different modalities through two symmetrical parallel branches, resulting in satisfactory saliency results.

#### D. Ablation Study

In this section, some experiments are conducted to verify the effectiveness of our proposed pipeline and key modules.

1) **Verification of the GDR and CRC modules:** We first conduct several experiments to demonstrate the effectiveness of the proposed GDR and CRC modules. Therefore, we keep the IPF module in this part of the experiment. In order to construct our baseline model, the GDR and CRC modules are removed from the two branches to construct our baseline model. Due to the different importance tendencies of the two modality features in these two branches, we retain the importance sensor in CRC to generate the importance response

TABLE II  
THE ABLATION VERIFICATION OF CRC AND GDR MODULES ON THE DAVIS AND DAVSOD DATASETS.

	DAVIS			DAVSOD		
	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$
B	0.891	0.904	0.020	0.658	0.746	0.082
B+GDR	0.895	0.907	0.019	0.663	0.750	0.080
B+CRC	0.904	0.913	0.017	0.671	0.759	0.075
B+CRC+GDR	0.907	0.919	0.016	0.678	0.765	0.074

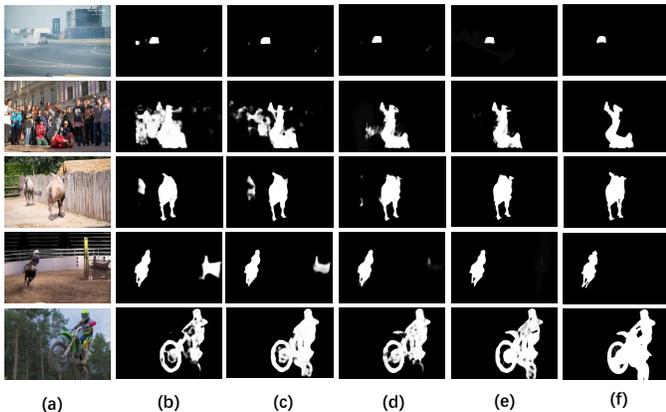


Fig. 7. Some visual comparisons of different ablation settings. (a) RGB images; (b) B; (c) B+GDR; (d) B+CRC; (e) B+GDR+CRC; (f) GT.

map. Then, we activate the auxiliary features by multiple them with the importance response map. And the activated auxiliary features are further concatenated with the dominant features in a particular branch. Finally, we accordingly construct a baseline model for our verification (denoted as ‘B’ in Table II). We gradually add CRC and GDR modules into the baseline model for ablation experiments, and the quantitative and qualitative results are shown in Table II and Fig. 7.

Firstly, we add the GDR module to the baseline model (denoted as ‘B + GDR’) to demonstrate the effectiveness of the proposed GDR module. In ‘B + GDR’, the basic setting is similar to ‘B’, but the auxiliary features are activated by the enhanced features from the GDR module. As reported in the second row of Table II, compared with the baseline model on the DAVSOD dataset, the MAE score is improved from 0.082 to 0.080, with a percentage gain of 2.4%. From Fig. 7(c), after introducing the GDR module, we can see that some background noise can be suppressed slightly, such as the left camel in the third image. In addition, we also add the CRC module to the baseline model (denoted as ‘B + CRC’) to verify the effectiveness of the CRC module. We can see that introducing the CRC module can boost performance compared with the baseline. Quantitatively, on the DAVIS dataset, introducing the CRC module achieves performance gains of 1.5% in terms of  $F_\beta$ , and 15.0% for the MAE score. As can be seen in Fig. 7(d), the model with the CRC module has better background suppression ability, such as the items on the right in the fourth image being effectively suppressed. Thus, these observations verify that

TABLE III  
THE IPF MODULE VERIFICATION ON DAVIS AND DAVSOD DATASETS. APPEARANCE REPRESENTS APPEARANCE-DOMINATED BRANCH. MOTION REPRESENTS MOTION-DOMINATED BRANCH.

	DAVIS			DAVSOD		
	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$	$F_\beta \uparrow$	$S_m \uparrow$	MAE $\downarrow$
Parallel-A	0.905	0.917	0.017	0.661	0.749	0.079
Parallel-C	0.900	0.915	0.019	0.671	0.761	0.076
Parallel-F	0.901	0.915	0.018	0.668	0.757	0.076
Parallel-IPF	0.907	0.919	0.016	0.678	0.765	0.074
Appearance	0.898	0.911	0.018	0.662	0.753	0.081
Motion	0.896	0.910	0.019	0.648	0.744	0.079

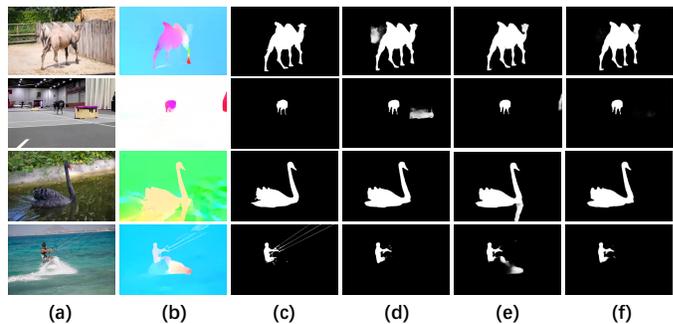


Fig. 8. Some visual comparisons on the output saliency results of Appearance-Dominated Branch, Motion-Dominated Branch, and our IPF. (a) RGB images; (b) Optical flow maps; (c) GT; (d) Saliency results deduced from Appearance-Dominated Branch; (e) Saliency results deduced from Motion-Dominated Branch; (f) Saliency results deduced from our proposed PSNet.

the CRC module can effectively extract useful complementary information from auxiliary modality features and further refine our more important modality features. Finally, both the GDR and CRC module are introduced into the baseline model to form our full model, which is denoted as ‘B + CRC + GDR’. Compared with other ablation settings in Table II, our full model achieves the best performance. From Fig. 7(e), we can see that our method achieves a more complete structure, more accurate localization, and clearer background.

2) *Verification of the IPF module*: To further verify the effectiveness of our proposed IPF module, we conduct the following ablation models.

- ‘Parallel-A’ denotes that a simple element-wise addition is used to fuse output features from two branches.
- ‘Parallel-C’ denotes that concatenation operation is used for fusing two output features.
- ‘Parallel-F’ denotes that channel-wise attention is used to adaptively fuse the output features of two branches.
- ‘Parallel-IPF’ denotes our proposed Importance Perception Fusion module, which uses high-level features in the backbone to sense the importance of two modality data.

We retain the multiplication operation for the above experiments to get the common response features between two branch output features and combine the common features with the fused features. As shown in Table III, our designed IPF module obtains a more robust and generalized final result by considering the different importance of the upper and lower

branches in different scenarios. Compared with the ‘Parallel-F’ mode, the percentage gain of the MAE score reaches 11.1% on the DAVIS dataset and 2.6% on the DAVSOD dataset, respectively. In addition, we also report the results for the upper and lower branches (*i.e.*, Appearance-Dominated Branch and Motion-Dominated Branch) under the full-model architecture with the IPF module, as shown in the last two rows of Table III. It can be seen that the results of any single branch cannot reach the results with the IPF module, which also illustrates the effectiveness and necessity of our IPF module design. Moreover, in order to further understand the effectiveness of the proposed IPF, some visualization results are shown in Fig. 8. As illustrated in Fig. 8, the motion-dominated branch achieves better saliency results in the top two rows of scenes. While in the last two rows, the appearance-dominated branch achieves better saliency results. And for all these scenes depicted in Fig. 8, our proposed PSNet with the IPF module achieves robust and stable saliency results.

## V. CONCLUSION

In this paper, we presented a parallel symmetric network (PSNet) for video salient object detection. Noticing that the importance between appearance cues and motion cues is different under different scenes, we propose to detect saliency via two parallel symmetric branches (*i.e.*, appearance-dominated branch and motion-dominated branch) in an explicitly discriminative way. These two branches have the same structure but regard different modality data as dominant features. Especially, the GDR module is proposed to highlight the multi-scale and multi-layer information, and the CRC module is designed to extract useful information from less important modality data and refine dominant modality data. We also introduce the IPF module to sense the importance weights of two modality data and fuse them adaptively. Extensive quantitative evaluations and visualization on four benchmark datasets demonstrate that our model achieves promising performance.

## REFERENCES

- [1] W. Wang, Q. Lai, H. Fu, J. Shen, H. Ling, and R. Yang, “Salient object detection in the deep learning era: An in-depth survey,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 6, pp. 3239–3259, 2022. **1**
- [2] R. Cong, J. Lei, H. Fu, M.-M. Cheng, W. Lin, and Q. Huang, “Review of visual saliency detection with comprehensive information,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 29, no. 10, pp. 2941–2959, 2019. **1, 2**
- [3] W. Wang, J. Shen, R. Yang, and F. Porikli, “Saliency-aware video object segmentation,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 1, pp. 20–33, 2018. **1**
- [4] W. Wang, J. Shen, X. Lu, S. C. H. Hoi, and H. Ling, “Paying attention to video object pattern understanding,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 7, pp. 2413–2428, 2021. **1**
- [5] P. Wen, R. Yang, Q. Xu, C. Qian, Q. Huang, R. Cong, and J. Si, “DMVOS: Discriminative matching for real-time video object segmentation,” in *Proc. ACM MM*, 2020, pp. 2048–2056. **1**
- [6] W. Wang, J. Shen, Y. Yu, and K.-L. Ma, “Stereoscopic thumbnail creation via efficient stereo saliency detection,” *IEEE Trans. Vis. Comput. Graphics*, vol. 23, no. 8, pp. 2014–2027, 2016. **1**
- [7] L. Itti, “Automatic foveation for video compression using a neurobiological model of visual attention,” *IEEE Trans. Image Process.*, vol. 13, no. 10, pp. 1304–1318, 2004. **1**
- [8] Z. Zhou, W. Pei, X. Li, H. Wang, F. Zheng, and Z. He, “Saliency-associated object tracking,” in *Proc. ICCV*, 2021, pp. 9846–9855. **1**
- [9] S. Li, Z. Tao, K. Li, and Y. Fu, “Visual to text: Survey of image and video captioning,” *IEEE Trans. Emerg. Topics in Comput. Intell.*, vol. 3, no. 4, pp. 297–312, 2019. **1**
- [10] W. Wang, J. Shen, and L. Shao, “Video salient object detection via fully convolutional networks,” *IEEE Trans. Image Process.*, vol. 27, no. 1, pp. 38–49, 2018. **1, 3, 8**
- [11] Y. Gu, L. Wang, Z. Wang, Y. Liu, M.-M. Cheng, and S.-P. Lu, “Pyramid constrained self-attention network for fast video salient object detection,” in *Proc. AAAI*, vol. 34, no. 07, 2020, pp. 10 869–10 876. **1, 8, 9**
- [12] C. Chen, G. Wang, C. Peng, Y. Fang, D. Zhang, and H. Qin, “Exploring rich and efficient spatial temporal interactions for real-time video salient object detection,” *IEEE Trans. Image Process.*, vol. 30, pp. 3995–4007, 2021. **1, 3, 8, 9**
- [13] H. Song, W. Wang, S. Zhao, J. Shen, and K.-M. Lam, “Pyramid dilated deeper ConvLSTM for video salient object detection,” in *Proc. ECCV*, 2018, pp. 715–731. **1, 3, 8**
- [14] D.-P. Fan, W. Wang, M.-M. Cheng, and J. Shen, “Shifting more attention to video salient object detection,” in *Proc. CVPR*, 2019, pp. 8554–8564. **1, 3, 7, 8**
- [15] Y. Su, W. Wang, J. Liu, P. Jing, and X. Yang, “DS-Net: Dynamic spatiotemporal network for video salient object detection,” *arXiv preprint arXiv:2012.04886*, 2020. **1, 4, 8**
- [16] H. Li, G. Chen, G. Li, and Y. Yu, “Motion guided attention for video salient object detection,” in *Proc. ICCV*, 2019, pp. 7274–7283. **1, 4, 8, 9**
- [17] G.-P. Ji, K. Fu, Z. Wu, D.-P. Fan, J. Shen, and L. Shao, “Full-duplex strategy for video object segmentation,” in *Proc. ICCV*, 2021, pp. 4922–4933. **1, 4, 7, 8**
- [18] P. Chen, J. Lai, G. Wang, and H. Zhou, “Confidence-guided adaptive gate and dual differential enhancement for video salient object detection,” in *Proc. ICME*, 2021, pp. 1–6. **1, 2, 4, 8, 9**
- [19] W. Wang and J. Shen, “Deep visual attention prediction,” *IEEE Trans. Image Process.*, vol. 27, no. 5, pp. 2368–2378, 2018. **2**
- [20] W. Wang, J. Shen, X. Dong, A. Borji, and R. Yang, “Inferring salient objects from human fixations,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 8, pp. 1913–1927, 2020. **2**
- [21] Z. Chen, Q. Xu, R. Cong, and Q. Huang, “Global context-aware progressive aggregation network for salient object detection,” in *Proc. AAAI*, 2020, pp. 10 599–10 606. **2, 3**
- [22] N. Liu, J. Han, and M.-H. Yang, “PiCANet: Learning pixel-wise contextual attention for saliency detection,” in *Proc. CVPR*, 2018, pp. 3089–3098. **2, 3**
- [23] J.-X. Zhao, J.-J. Liu, D.-P. Fan, Y. Cao, J. Yang, and M.-M. Cheng, “EGNet: Edge guidance network for salient object detection,” in *Proc. ICCV*, 2019, pp. 8779–8788. **2, 3, 8**
- [24] X. Qin, Z. Zhang, C. Huang, C. Gao, M. Dehghan, and M. Jagersand, “BASNet: Boundary-aware salient object detection,” in *Proc. CVPR*, 2019, pp. 7479–7489. **2, 3**
- [25] 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. 121–134, 2017. **2, 3**
- [26] Y. Pang, X. Zhao, L. Zhang, and H. Lu, “Multi-scale interactive network for salient object detection,” in *Proc. CVPR*, 2020, pp. 9413–9422. **2, 3**
- [27] C. Fang, H. Tian, D. Zhang, Q. Zhang, J. Han, and J. Han, “Densely nested top-down flows for salient object detection,” *Science China Information Sciences*, vol. 65, no. 8, pp. 1–14, 2022. **2**
- [28] Q. Zhang *et al.*, “Dense attention fluid network for salient object detection in optical remote sensing images,” *IEEE Trans. Image Process.*, vol. 30, pp. 1305–1317, 2021. **2**
- [29] R. Cong, Y. Zhang, L. Fang, J. Li, Y. Zhao, and S. Kwong, “RRNet: Relational reasoning network with parallel multi-scale attention for salient object detection in optical remote sensing images,” *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1558–0644, 2022. **2**
- [30] R. Cong, Q. Qin, C. Zhang, Q. Jiang, S. Wang, Y. Zhao, and S. Kwong, “A weakly supervised learning framework for salient object detection via hybrid labels,” *IEEE Trans. Circuits Syst. Video Technol.*, early access, doi: 10.1109/TCSVT.2022.3205182. **2**
- [31] X. Zhou, K. Shen, L. Weng, R. Cong, B. Zheng, J. Zhang, and C. Yan, “Edge-guided recurrent positioning network for salient object detection in optical remote sensing images,” *IEEE Trans. Cybern.*, early access, doi: 10.1109/TCYB.2022.3163152. **2**
- [32] C. Li, R. Cong, J. Hou, S. Zhang, Y. Qian, and S. Kwong, “Nested network with two-stream pyramid for salient object detection in optical remote sensing images,” *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 11, pp. 9156–9166, 2019. **2**

- [33] D. Zhang, J. Han, Y. Zhang, and D. Xu, "Synthesizing supervision for learning deep saliency network without human annotation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 7, pp. 1755–1769, 2019. [2](#)
- [34] H. Li, S. Cong, Runminand Kwong, C. Chen, Q. Xu, and C. Li, "Stereo superpixel: An iterative framework based on parallax consistency and collaborative optimization," *Information Sciences*, vol. 556, pp. 209–222, 2021. [2](#)
- [35] C. Li, J. Guo, B. Wang, R. Cong, Y. Zhang, and J. Wang, "Single underwater image enhancement based on color cast removal and visibility restoration," *J. Electronic Imaging*, vol. 25, no. 3, p. 033012, 2016. [3](#)
- [36] C. Li, S. Anwar, J. Hou, R. Cong, C. Guo, and W. Ren, "Underwater image enhancement via medium transmission-guided multi-color space embedding," *IEEE Trans. Image Process.*, vol. 30, pp. 4985–5000, 2021. [3](#)
- [37] J. Hu, Q. Jiang, R. Cong, W. Gao, and F. Shao, "Two-branch deep neural network for underwater image enhancement in HSV color space," *IEEE Signal Process. Lett.*, vol. 28, pp. 2152–2156, 2021. [3](#)
- [38] C. Guo, C. Li, J. Guo, C. C. Loy, J. Hou, S. Kwong, and R. Cong, "Zero-reference deep curve estimation for low-light image enhancement," in *Proc. CVPR*, 2020, pp. 1777–1786. [3](#)
- [39] C. Li, C. Guo, J. Guo, P. Han, H. Fu, and R. Cong, "PDR-Net: Perception-inspired single image dehazing network with refinement," *IEEE Trans. Multimedia*, vol. 22, no. 3, pp. 704–716, 2020. [3](#)
- [40] Q. Jiang, Y. Mao, R. Cong, W. Ren, C. Huang, and F. Shao, "Unsupervised decomposition and correction network for low-light image enhancement," *IEEE Trans. Intell. Transp. Syst.*, early access, doi: 10.1109/TITS.2022.3165176. [3](#)
- [41] F. Li, Y. Wu, H. Bai, W. Lin, R. Cong, and Y. Zhao, "Learning detail-structure alternative optimization for blind super-resolution," *IEEE Trans. Multimedia*, early access, doi: 10.1109/TMM.2022.3152090. [3](#)
- [42] Q. Tang, R. Cong, R. Sheng, L. He, D. Zhang, Y. Zhao, and S. Kwong, "Bridgenet: A joint learning network of depth map super-resolution and monocular depth estimation," in *Proc. ACM MM*, 2021, pp. 2148–2157. [3](#)
- [43] L. He, H. Zhu, F. Li, H. Bai, R. Cong, C. Zhang, C. Lin, M. Liu, and Y. Zhao, "Towards fast and accurate real-world depth super-resolution: Benchmark dataset and baseline," in *Proc. IEEE CVPR*, 2021, pp. 9229–9238. [3](#)
- [44] C. Guo, C. Li, J. Guo, R. Cong, H. Fu, and P. Han, "Hierarchical features driven residual learning for depth map super-resolution," *IEEE Trans. Image Process.*, vol. 28, no. 5, pp. 2545–2557, 2019. [3](#)
- [45] F. Li, R. Cong, H. Bai, and Y. He, "Deep interleaved network for image super-resolution with asymmetric co-attention," in *Proc. IJCAI*, 2020, pp. 534–543. [3](#)
- [46] Z. Wang, F. Li, R. Cong, H. Bai, and Y. Zhao, "Adaptive feature fusion network based on boosted attention mechanism for single image dehazing," *Multim. Tools Appl.*, vol. 81, no. 8, pp. 11 325–11 339, 2022. [3](#)
- [47] N. Yang, Q. Zhong, K. Li, R. Cong, Y. Zhao, and S. Kwong, "A reference-free underwater image quality assessment metric in frequency domain," *Signal Process. Image Commun.*, vol. 94, p. 116218, 2021. [3](#)
- [48] Q. Jiang, Y. Gu, C. Li, R. Cong, and F. Shao, "Underwater image enhancement quality evaluation: Benchmark dataset and objective metric," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 32, no. 9, pp. 5959–5974, 2022. [3](#)
- [49] C. Zhang, R. Cong, Q. Lin, L. Ma, F. Li, Y. Zhao, and S. Kwong, "Cross-modality discrepant interaction network for RGB-D salient object detection," in *Proc. ACM MM*, 2021, pp. 2094–2102. [3](#)
- [50] C. Li, R. Cong, S. Kwong, J. Hou, H. Fu, G. Zhu, D. Zhang, and Q. Huang, "ASIF-Net: Attention steered interweave fusion network for RGB-D salient object detection," *IEEE Trans. Cybern.*, vol. 50, no. 1, pp. 88–100, 2021. [3](#)
- [51] H. Wen, C. Yan, X. Zhou, R. Cong, Y. Sun, B. Zheng, J. Zhang, Y. Bao, and G. Ding, "Dynamic selective network for RGB-D salient object detection," *IEEE Trans. Image Process.*, vol. 30, pp. 9179–9192, 2021. [3](#)
- [52] C. Li, R. Cong, Y. Piao, Q. Xu, and C. C. Loy, "RGB-D salient object detection with cross-modality modulation and selection," in *Proc. ECCV*, 2020, pp. 225–241. [3](#)
- [53] Y. Mao, Q. Jiang, R. Cong, W. Gao, F. Shao, and S. Kwong, "Cross-modality fusion and progressive integration network for saliency prediction on stereoscopic 3D images," *IEEE Trans. Multimedia*, vol. 24, pp. 2435–2448, 2022. [3](#)
- [54] R. Cong, J. Lei, H. Fu, J. Hou, Q. Huang, and S. Kwong, "Going from RGB to RGBD saliency: A depth-guided transformation model," *IEEE Trans. Cybern.*, vol. 50, no. 8, pp. 3627–3639, 2020. [3](#)
- [55] R. Cong, Q. Lin, C. Zhang, c. Li, X. Cao, Q. Huang, and Y. Zhao, "CIR-Net: Cross-modality interaction and refinement for RGB-D salient object detection," *IEEE Trans. Image Process.*, 2022. [3](#)
- [56] R. Cong, K. Zhang, C. Zhang, F. Zheng, Y. Zhao, Q. Huang, and S. Kwong, "Does thermal really always matter for RGB-T salient object detection?" *IEEE Trans. Multimedia*, 2022. [3](#)
- [57] D. Zhang, J. Han, C. Li, and J. Wang, "Co-saliency detection via looking deep and wide," in *Proc. CVPR*, 2015. [3](#)
- [58] Q. Zhang, R. Cong, J. Hou, C. Li, and Y. Zhao, "CoADNet: Collaborative aggregation-and-distribution networks for co-salient object detection," in *Proc. NeurIPS*, 2020, pp. 6959–6970. [3](#)
- [59] R. Cong, N. Yang, C. Li, H. Fu, Y. Zhao, Q. Huang, and S. Kwong, "Global-and-local collaborative learning for co-salient object detection," *IEEE Trans. Cybern.*, early access, doi: 10.1109/TCYB.2022.3169431. [3](#)
- [60] R. Cong, J. Lei, H. Fu, Q. Huang, X. Cao, and N. Ling, "HSCS: Hierarchical sparsity based co-saliency detection for RGBD images," *IEEE Trans. Multimedia*, vol. 21, no. 7, pp. 1660–1671, 2019. [3](#)
- [61] Y. Zhang, L. Li, R. Cong, X. Guo, H. Xu, and J. Zhang, "Co-saliency detection via hierarchical consistency measure," in *Proc. IEEE ICME*, 2018, pp. 1–6. [3](#)
- [62] R. Cong, J. Lei, H. Fu, W. Lin, Q. Huang, X. Cao, and C. Hou, "An iterative co-saliency framework for RGBD images," *IEEE Trans. Cybern.*, vol. 49, no. 1, pp. 233–246, 2019. [3](#)
- [63] R. Cong, J. Lei, H. Fu, Q. Huang, X. Cao, and C. Hou, "Co-saliency detection for RGBD images based on multi-constraint feature matching and cross label propagation," *IEEE Trans. Image Process.*, vol. 27, no. 2, pp. 568–579, 2018. [3](#)
- [64] K. Zhang, T. Li, S. Shen, B. Liu, J. Chen, and Q. Liu, "Adaptive graph convolutional network with attention graph clustering for co-saliency detection," in *Proc. CVPR*, 2020, pp. 9050–9059. [3](#)
- [65] B. Jiang, X. Jiang, A. Zhou, J. Tang, and B. Luo, "A unified multiple graph learning and convolutional network model for co-saliency estimation," in *Proc. ACM MM*, 2019, pp. 1375–1382. [3](#)
- [66] D.-P. Fan, T. Li, Z. Lin, G.-P. Ji, D. Zhang, M.-M. Cheng, H. Fu, and J. Shen, "Re-thinking co-salient object detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 8, pp. 4339–4354, 2022. [3](#)
- [67] Z. Tu, Z. Guo, W. Xie, M. Yan, R. C. Veltkamp, B. Li, and J. Yuan, "Fusing disparate object signatures for salient object detection in video," *Pattern Recognit.*, vol. 72, pp. 285–299, 2017. [3, 8](#)
- [68] C. Chen, S. Li, Y. Wang, H. Qin, and A. Hao, "Video saliency detection via spatial-temporal fusion and low-rank coherency diffusion," *IEEE Trans. Image Process.*, vol. 26, no. 7, pp. 3156–3170, 2017. [3, 8](#)
- [69] F. Guo, W. Wang, J. Shen, L. Shao, J. Yang, D. Tao, and Y. Y. Tang, "Video saliency detection using object proposals," *IEEE Trans. Cybern.*, vol. 48, no. 11, pp. 3159–3170, 2017. [3](#)
- [70] R. Cong, J. Lei, H. Fu, F. Porikli, Q. Huang, and C. Hou, "Video saliency detection via sparsity-based reconstruction and propagation," *IEEE Trans. Image Process.*, vol. 28, no. 10, pp. 4819–4931, 2019. [3](#)
- [71] G. Li, Y. Xie, T. Wei, K. Wang, and L. Lin, "Flow guided recurrent neural encoder for video salient object detection," in *Proc. CVPR*, 2018, pp. 3243–3252. [3, 8](#)
- [72] C. Chen, J. Song, C. Peng, G. Wang, and Y. Fang, "A novel video salient object detection method via semisupervised motion quality perception," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 32, no. 5, pp. 2732–2745, 2021. [3](#)
- [73] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox, "FlowNet 2.0: Evolution of optical flow estimation with deep networks," in *Proc. CVPR*, 2017, pp. 2462–2470. [3, 4, 7](#)
- [74] S. Ren, C. Han, X. Yang, G. Han, and S. He, "TENet: Triple excitation network for video salient object detection," in *Proc. ECCV*, 2020, pp. 212–228. [4](#)
- [75] W. Zhao, J. Zhang, L. Li, N. Barnes, N. Liu, and J. Han, "Weakly supervised video salient object detection," in *Proc. CVPR*, June 2021, pp. 16 826–16 835. [4, 8](#)
- [76] X. Lu, W. Wang, C. Ma, J. Shen, L. Shao, and F. Porikli, "See more, know more: Unsupervised video object segmentation with co-attention siamese networks," in *Proc. CVPR*, 2019, pp. 3623–3632. [4](#)
- [77] T. Zhou, J. Li, S. Wang, R. Tao, and J. Shen, "MATNet: Motion-attentive transition network for zero-shot video object segmentation," *IEEE Trans. Image Process.*, vol. 29, pp. 8326–8338, 2020. [4](#)
- [78] W. Wang, X. Lu, J. Shen, D. J. Crandall, and L. Shao, "Zero-shot video object segmentation via attentive graph neural networks," in *Proc. ICCV*, 2019, pp. 9235–9244. [4](#)

- [79] S. Cho, M. Lee, S. Lee, C. Park, D. Kim, and S. Lee, "Treating motion as option to reduce motion dependency in unsupervised video object segmentation," *arXiv preprint arXiv:2209.03138*, 2022. 4
- [80] Y. Hong, H. Pan, W. Sun, and Y. Jia, "Deep dual-resolution networks for real-time and accurate semantic segmentation of road scenes," *arXiv preprint arXiv:2101.06085*, 2021. 5
- [81] Q. Hou, M. Cheng, X. Hu, A. Borji, Z. Tu, and P. H. S. Torr, "Deeply supervised salient object detection with short connections," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 4, pp. 815–828, 2019. 5
- [82] T. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie, "Feature pyramid networks for object detection," in *Proc. CVPR*, 2017, pp. 936–944. 5
- [83] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. CVPR*, 2017, pp. 4700–4708. 6
- [84] X. Jia, B. De Brabandere, T. Tuytelaars, and L. V. Gool, "Dynamic filter networks," in *Proc. NeurIPS*, 2016, pp. 667–675. 6
- [85] J. He, Z. Deng, and Y. Qiao, "Dynamic multi-scale filters for semantic segmentation," in *Proc. ICCV*, 2019, pp. 3562–3572. 6
- [86] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proc. CVPR*, 2018, pp. 7132–7141. 6
- [87] R. Cong, H. Yang, Y. Zhang, N. Yang, H. Li, X. Zhang, R. Li, Z. Chen, Y. Zhao, and S. Kwong, "Boundary guided semantic learning for real-time COVID-19 lung infection segmentation system," *IEEE Trans. Consum. Electron.*, early access, doi: 10.1109/TCE.2022.3205376. 6
- [88] R. Cong, H. Yang, Q. Jiang, W. Gao, H. Li, C. Wang, Y. Zhao, and S. Kwong, "BCS-Net: Boundary, context, and semantic for automatic COVID-19 lung infection segmentation from CT images," *IEEE Trans. Instrum. Meas.*, vol. 71, pp. 1–11, 2022. 6
- [89] Z. Chen, R. Cong, Q. Xu, and Q. Huang, "DPANet: Depth potentiality-aware gated attention network for RGB-D salient object detection," *IEEE Trans. Image Process.*, vol. 30, pp. 7012–7024, 2021. 6
- [90] Z. Wu, L. Su, and Q. Huang, "Cascaded partial decoder for fast and accurate salient object detection," in *Proc. CVPR*, 2019, pp. 3907–3916. 8
- [91] H. Zhou, X. Xie, J.-H. Lai, Z. Chen, and L. Yang, "Interactive two-stream decoder for accurate and fast saliency detection," in *Proc. CVPR*, 2020, pp. 9141–9150. 8
- [92] W.-C. Tu, S. He, Q. Yang, and S.-Y. Chien, "Real-time salient object detection with a minimum spanning tree," in *Proc. CVPR*, 2016, pp. 2334–2342. 8
- [93] Z. Liu, J. Li, L. Ye, G. Sun, and L. Shen, "Saliency detection for unconstrained videos using superpixel-level graph and spatiotemporal propagation," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 27, no. 12, pp. 2527–2542, 2016. 8
- [94] T. Xi, W. Zhao, H. Wang, and W. Lin, "Salient object detection with spatiotemporal background priors for video," *IEEE Trans. Image Process.*, vol. 26, no. 7, pp. 3425–3436, 2016. 8
- [95] Y. Chen, W. Zou, Y. Tang, X. Li, C. Xu, and N. Komodakis, "SCOM: Spatiotemporal constrained optimization for salient object detection," *IEEE Trans. Image Process.*, vol. 27, no. 7, pp. 3345–3357, 2018. 8
- [96] Y. Tang, W. Zou, Z. Jin, Y. Chen, Y. Hua, and X. Li, "Weakly supervised salient object detection with spatiotemporal cascade neural networks," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 29, no. 7, pp. 1973–1984, 2018. 8
- [97] S. Li, B. Seybold, A. Vorobyov, X. Lei, and C.-C. J. Kuo, "Unsupervised video object segmentation with motion-based bilateral networks," in *Proc. ECCV*, 2018, pp. 207–223. 8
- [98] P. Yan, G. Li, Y. Xie, Z. Li, C. Wang, T. Chen, and L. Lin, "Semi-supervised video salient object detection using pseudo-labels," in *Proc. ICCV*, 2019, pp. 7284–7293. 8
- [99] Y. Ji, H. Zhang, Z. Jie, L. Ma, and Q. J. Wu, "CASNet: A cross-attention siamese network for video salient object detection," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 32, no. 6, pp. 2676–2690, 2020. 8
- [100] J. Mei, M. Wang, Y. Lin, Y. Yuan, and Y. Liu, "TransVOS: Video object segmentation with transformers," *arXiv preprint arXiv:2106.00588*, 2021. 8
- [101] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung, "A benchmark dataset and evaluation methodology for video object segmentation," in *Proc. CVPR*, 2016, pp. 724–732. 7
- [102] F. Li, T. Kim, A. Humayun, D. Tsai, and J. M. Rehg, "Video segmentation by tracking many figure-ground segments," in *Proc. ICCV*, 2013, pp. 2192–2199. 7
- [103] J. Li, C. Xia, and X. Chen, "A benchmark dataset and saliency-guided stacked autoencoders for video-based salient object detection," *IEEE Trans. Image Process.*, vol. 27, no. 1, pp. 349–364, 2017. 7
- [104] R. Achanta, S. Hemami, F. Estrada, and S. Susstrunk, "Frequency-tuned salient region detection," in *Proc. CVPR*, 2009, pp. 1597–1604. 7
- [105] D.-P. Fan, M.-M. Cheng, Y. Liu, T. Li, and A. Borji, "Structure-measure: A new way to evaluate foreground maps," in *Proc. ICCV*, 2017, pp. 4548–4557. 7
- [106] C. Li, R. Cong, C. Guo, H. Li, C. Zhang, F. Zheng, and Y. Zhao, "A parallel down-up fusion network for salient object detection in optical remote sensing images," *Neurocomputing*, vol. 415, pp. 411–420, 2020. 7
- [107] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. CVPR*, 2016, pp. 770–778. 7
- [108] L. Wang, H. Lu, Y. Wang, M. Feng, D. Wang, B. Yin, and X. Ruan, "Learning to detect salient objects with image-level supervision," in *Proc. CVPR*, 2017, pp. 136–145. 7