# Salient Object Detection in Optical Remote Sensing Images Driven by Transformer

Gongyang Li, Zhen Bai, Zhi Liu, *Senior Member, IEEE*, Xinpeng Zhang, *Member, IEEE*, and Haibin Ling, *Fellow, IEEE*

*Abstract*—Existing methods for *Salient Object Detection in Optical Remote Sensing Images* (ORSI-SOD) mainly adopt *Convolutional Neural Networks* (CNNs) as the backbone, such as VGG and ResNet. Since CNNs can only extract features within certain receptive fields, most ORSI-SOD methods generally follow the local-to-contextual paradigm. In this paper, we propose a novel *Global Extraction Local Exploration Network* (GeleNet) for ORSI-SOD following the global-to-local paradigm. Specifically, GeleNet first adopts a transformer backbone to generate four-level feature embeddings with global long-range dependencies. Then, GeleNet employs a *Direction-aware Shuffle Weighted Spatial Attention Module* (D-SWSAM) and its simplified version (SWSAM) to enhance local interactions, and a *Knowledge Transfer Module* (KTM) to further enhance cross-level contextual interactions. D-SWSAM comprehensively perceives the orientation information in the lowest-level features through directional convolutions to adapt to various orientations of salient objects in ORSIs, and effectively enhances the details of salient objects with an improved attention mechanism. SWSAM discards the direction-aware part of D-SWSAM to focus on localizing salient objects in the highest-level features. KTM models the contextual correlation knowledge of two middle-level features of different scales based on the self-attention mechanism, and transfers the knowledge to the raw features to generate more discriminative features. Finally, a saliency predictor is used to generate the saliency map based on the outputs of the above three modules. Extensive experiments on three public datasets demonstrate that the proposed GeleNet outperforms relevant state-of-the-art methods. The code and results of our method are available at https://github.com/MathLee/GeleNet.

*Index Terms*—Salient object detection, optical remote sensing image, transformer, directional convolution, shuffle weighted spatial attention.

## I. INTRODUCTION

SALIENT Object Detection (SOD) focuses on finding and locating the most visually prominent objects/regions in a scene [1]–[3]. It is a common pre-processing step for many tasks in computer vision, such as quality assessment [4], [5], object segmentation [6]–[10], video compression [11], and object tracking [12]. Recently, SOD in *Optical Remote*

Gongyang Li, Zhen Bai, Zhi Liu, and Xinpeng Zhang are with Key Laboratory of Specialty Fiber Optics and Optical Access Networks, Joint International Research Laboratory of Specialty Fiber Optics and Advanced Communication, Shanghai Institute for Advanced Communication and Data Science, Shanghai University, Shanghai 200444, China, and School of Communication and Information Engineering, Shanghai University, Shanghai 200444, China. Gongyang Li and Zhi Liu are also with Wenzhou Institute of Shanghai University, Wenzhou 325000, China (email: ligongyang@shu.edu.cn; bz536476@163.com; liuzhisjtu@163.com; xzhang@shu.edu.cn).

Haibin Ling is with the Department of Computer Science, Stony Brook University, Stony Brook, NY 11794 USA (email: hling@cs.stonybrook.edu).

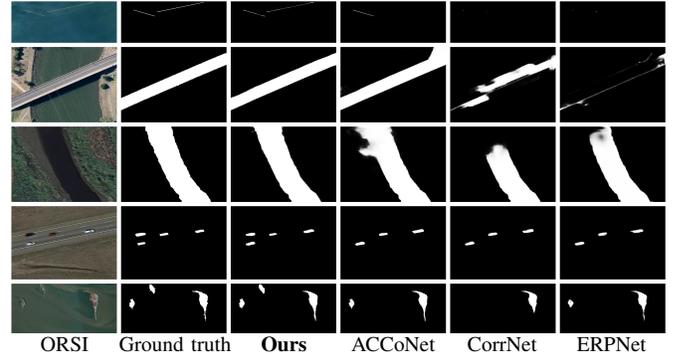*Corresponding authors: Zhi Liu and Xinpeng Zhang.*



Fig. 1. Saliency maps produced by our method and three state-of-the-art ORSI-SOD methods, including ACCoNet [16], CorrNet [17], and ERPNet [18]. Please zoom in for details, especially the first row.

*Sensing Images* (ORSI-SOD) [13]–[15], as an emerging topic, has attracted the attention of many researchers, and has been widely used in agriculture, forestry, environmental science, and security surveillance.

With the rapid development of deep learning, *Convolutional Neural Networks* (CNNs) [19] have dominated the field of computer vision with their powerful feature representation capabilities. Many effective CNN-based solutions for ORSI-SOD are proposed [15]–[18], [20]–[23]. While a few methods follow the local-to-local paradigm [18], [20], most methods adopt the local-to-contextual paradigm[1] [15]–[17], [21]–[25]. Both paradigms first use a CNN backbone, such as VGG [26] and ResNet [27], to extract basic feature embeddings. The local-to-local paradigm focuses on exploring valuable information in single-level feature embeddings. Differently, since the local-to-contextual one considers that CNNs only extract features within certain receptive fields, it focuses on designing specific modules to mine the contextual information between feature embeddings at multiple levels. The above paradigms promote the development of ORSI-SOD and achieve promising performance.

However, due to the characteristics of ORSI scenes, such as variation in object orientation, scale, and category, the above paradigms suffer from obvious limitations. The local-to-local paradigm ignores contextual information that is useful for handling the above scenes. The contextual information captured by the local-to-contextual paradigm is still based on convolution layers with limited receptive fields, which is also insufficient to handle challenging scenes of ORSIs. For

---

[1]Here, the first "local" in both paradigms specifically refers to using CNN backbones to extract features with limited receptive fields.

intuitive understanding, we show the saliency maps generated by typical methods for both paradigms in Fig. 1, where ACCoNet [16] and CorrNet [17] belong to the local-to-contextual paradigm, and ERPNet [18] belongs to the local-to-local paradigm. We find that these methods suffer from orientation insensitivity, incomplete detection, and missing salient objects.

Inspired by the above observations, in this paper, we propose to design new solutions following the global-to-local paradigm. Our main idea is to replace the CNN backbone with a transformer one that can establish global relationships (*i.e.*, changing the first "local" in two existing paradigms to "global") and to perform local enhancement on the extracted global features. With this idea, we build a novel *Global Extraction Local Exploration Network* (GeleNet) for ORSI-SOD with a transformer backbone. Transformers [28]–[31] are known to be good at modeling the global long-range dependencies between feature patches. This unique ability of transformer enables GeleNet to deal with complex scenes and changeable objects in ORSIs. Furthermore, in GeleNet, we focus on local and cross-level contextual interactions, which are beneficial for highlighting salient objects in ORSIs.

In particular, we adopt the popular PVT [32], [33] as the backbone of our GeleNet. To alleviate the orientation insensitivity issue of previous methods, we propose a *Direction-aware Shuffle Weighted Spatial Attention Module* (D-SWSAM), and assign it to the lowest-level features to adequately identify the orientation of objects through directional convolutions with four directions. D-SWSAM is also equipped with an improved attention mechanism to outline the details of salient objects. Since high-level features contain location information rather than orientation and texture information, we extract the corresponding part containing the improved attention mechanism from D-SWSAM, *i.e.*, SWSAM, and assign it to the highest-level features to determine the location of salient objects. The above modules can well enhance local interactions of intra-level features. In addition, we propose a *Knowledge Transfer Module* (KTM) for the remaining adjacent features to explore contextual interactions between inter-level features and transfer the specific knowledge of salient objects between adjacent features to the raw features. In this way, the proposed GeleNet can generate saliency maps with accurate orientations and complete objects, as illustrated in the third column of Fig. 1, and consistently outperforms compared methods on three datasets.

Our main contributions are summarized in three aspects:

- We propose a transformer-based ORSI-SOD solution, *GeleNet*, with the global-to-local paradigm, which is different from the local-to-contextual paradigm followed by most existing CNN-based methods. To the best of our knowledge, this is the first transformer-driven ORSI-SOD solution.
- We propose the D-SWSAM and its variant SWSAM to enhance local interactions of the extracted global feature embeddings. D-SWSAM can tackle the problem of objects with various orientations in ORSIs and enhance the details of salient objects in the lowest-level features,

while SWSAM can locate salient objects in the highest-level features.
- We propose the KTM to enhance contextual interactions of two middle-level features. In KTM, we model the contextual correlation knowledge of two types of combinations (*i.e.*, product and sum) of these features, and transfer the knowledge to the raw features to generate more discriminative features.

The rest of this paper is arranged as follows. In Sec. II, we review the related work. In Sec. III, we describe the details of the proposed GeleNet. In Sec. IV, we conduct comprehensive experiments and ablation studies. In Sec. V, we present the conclusion.

## II. RELATED WORK

### A. Salient Object Detection in Optical Remote Sensing Images

Salient object detection in optical remote sensing images plays an important role in understanding ORSIs. Recently, with the successive construction of the three datasets [14], [15], [22], numerous ORSI-SOD methods are proposed. Here we focus on CNN-based methods, which dominate this topic and achieve promising performance.

Existing CNN-based ORSI-SOD methods mainly follow two paradigms, *i.e.*, the local-to-local paradigm and the local-to-contextual paradigm. The local-to-local paradigm typically extracts feature embeddings containing local information through the CNN backbone, and then explores valuable information in single-level feature embeddings. For example, in [18], Zhou *et al*. extracted multi-level features through the CNN backbone, and performed edge extraction and feature fusion on each level of features in two parallel decoders. Li *et al*. [20] explored the complementarity of foreground, edge, background, and the global image-level content of single-level features, and aimed at generating complete salient objects. They focused on the extraction of various specific information on single-level features (*i.e.*, local features), ignoring the contextual interactions between local features at different levels.

The local-to-contextual paradigm, by contrast, explores contextual information between local feature embeddings at different levels, and is therefore popularly adopted by recent solutions. For example, Li *et al*. [15] extracted multi-level features from multiple inputs, and employed nested connections to aggregate them. Similarly, Zhou *et al*. [23] proposed a cascaded feature fusion module to fuse multi-level features from different branches. In [21], Huang *et al*. aggregated three high-level features to produce contextual semantic information to approximately locate salient objects. Li *et al*. [17] proposed a correlation module for continuous semantic features, generating an initial coarse saliency map for location guidance of low-level features. Tu *et al*. [22] proposed two decoders to aggregate three adjacent features twice with salient boundary features. Li *et al*. [16] designed a specific module for adjacent features, aiming at coordinating cross-scale interactions and mining valuable contextual information.

Despite great progress achieved by the local-to-contextual paradigm, the explored contextual interactions only mine

interactions between features at different levels through convolution-based modules. In this paper, inspired by the popular transformer [28]–[32], we propose the global-to-local paradigm that first models the global long-range dependencies between feature patches and then enhances the local and contextual interactions, and build a novel GeleNet for ORSI-SOD. Benefiting from the global view of the transformer and the local enhancement of our proposed modules, our GeleNet can better perceive salient objects with numerous scales, diverse types, and multiple numbers in ORSIs.

### B. Salient Object Detection with Transformer

Transformer was first proposed for *Natural Language Processing* (NLP) [28], which is good at modeling global long-range dependencies between word vectors. Following its success in NLP, researchers have extended it into computer vision and achieved remarkable progress on numerous tasks, especially on dense prediction tasks [31]–[33].

Here, we introduce some representative works on transformer-based SOD, involving SOD in *Natural Scene Images* (NSI-SOD) [34], [35], RGB-D/T SOD [36]–[38], co-saliency detection [39], and video SOD [39]. In general, transformer-based SOD methods can be roughly divided into three types depending on where the transformer is used. The first type of method adopted transformer as the feature extractor in the encoding phase. For instance, Liu *et al*. [36] used Swin Transformer [31] to extract basic features from RGB-D/T pairs, and aligned cross-modality features through attention mechanism to generate discriminative features. Liu *et al*. [34] achieved effective context modeling using the same backbone as [36] for NSI-SOD. The second type of method adopted transformer to develop modules in the decoding phase. Liu *et al*. [37] proposed a triplet transformer embedding module to enhance high-level features by learning long-range dependencies across layers. In [39], Su *et al*. proposed a unified transformer framework for co-saliency detection and video SOD, which is equipped with two transformer blocks to capture the long-range dependencies among a group of features from different images/frames. Fang *et al*. [38] proposed a multiple transformer module to learn the common information of cross-modality and cross-scale features. The last type of method utilized the pure transformer architecture to achieve SOD. Liu *et al*. [35] adopted T2T-ViT [30] as the backbone, and proposed a multi-task transformer decoder to jointly detect salient objects and boundaries.

The above transformer-based SOD methods achieve impressive results on specific SOD tasks. Therefore, we introduce the transformer into the ORSI-SOD task, and propose the first transformer-driven ORSI-SOD method, *i.e.*, GeleNet. Our method belongs to the first type of method, and adopts PVT [32], [33] as the backbone to extract long-range dependency features from input ORSIs.

### C. Attention Mechanism

Attention mechanism is widely used in computer vision and image analysis. In general, it includes channel attention [40], spatial attention [41], and self-attention [28], [42]. SENet [40] was a classic channel attention model, which explicitly represents dependencies between channels to adaptively recalibrate features. ECANet [43] developed an extremely lightweight channel attention module through a fast 1D convolution. Moreover, CBAM [41] additionally introduced spatial attention, and inferred attention maps along channel and spatial domain in turn for adaptive feature enhancement. Li *et al*. [44] proposed the *Spatial Group-wise Enhance* (SGE), which first splits features into several sub-features, then extracts specific semantics from each sub-feature, and finally adjusts the importance of semantics of each sub-feature by an attention factor. Zhang *et al*. [45] proposed a lightweight shuffle attention, which also first splits features into several groups, then performs channel attention and spatial attention in parallel, and finally introduces channel shuffle to allow information communication along channels.

Both SGE [44] and shuffle attention [45] consider only the attention of each sub-feature, but ignore the consistency of attention between different sub-features, which is not friendly to SOD. In addition, since the global features extracted by the transformer lack channel interaction, it is unreasonable for shuffle attention to put the shuffle operation at the end. Therefore, we propose an improved spatial attention module, namely SWSAM, which focuses on enhancing the channel interactions of global features and improving the effectiveness of spatial attention to highlight salient regions more accurately. Notably, we further integrate SWSAM and directional convolutions, and propose D-SWSAM to adapt to various orientations of salient objects in ORSIs. Moreover, we also propose a self-attention-based KTM to model and transfer the contextual knowledge to generate more discriminative features.

## III. PROPOSED METHOD

In this section, we elaborate on the proposed transformer-driven GeleNet. In Sec. III-A, we depict the network overview. In Sec. III-B and Sec. III-C, we introduce D-SWSAM and KTM, respectively. In Sec. III-D, we present the saliency predictor and loss function.

### A. Network Overview

As illustrated in Fig. 2, the proposed GeleNet follows the popular three-stage structure [46], [47] in SOD, including a feature extractor for basic feature generation, three modules (*i.e.*, D-SWSAM, KTM, and SWSAM) for feature interaction/enhancement, and a saliency predictor for saliency map generation.

Concretely, we use the *Pyramid Vision Transformer* (PVT) [33] as the backbone, whose input size is set to $3 \times 352 \times 352$. PVT consists of four transformer encoder blocks denoted as $T^i$ ($i \in \{1, 2, 3, 4\}$), and can generate four-level basic global features denoted as $\hat{\boldsymbol{f}}_t^i \in \mathbb{R}^{c_i \times h_i \times w_i}$, where $c_i \in \{64, 128, 320, 512\}$, and $h_i/w_i = \frac{352}{2^{i+1}}$. To improve the computational efficiency, we unify the channel number of $\hat{\boldsymbol{f}}_t^i$ ($i \in \{1, 3, 4\}$) to 32 by the channel normalization (*i.e.*, a convolution layer), generating $\boldsymbol{f}_t^i \in \mathbb{R}^{c \times h_i \times w_i}$, where $c$ is 32. Notably, for $\hat{\boldsymbol{f}}_t^2$, we not only reduce its channel number to 32, but also adjust its resolution from $44 \times 44$ to $22 \times 22$ for
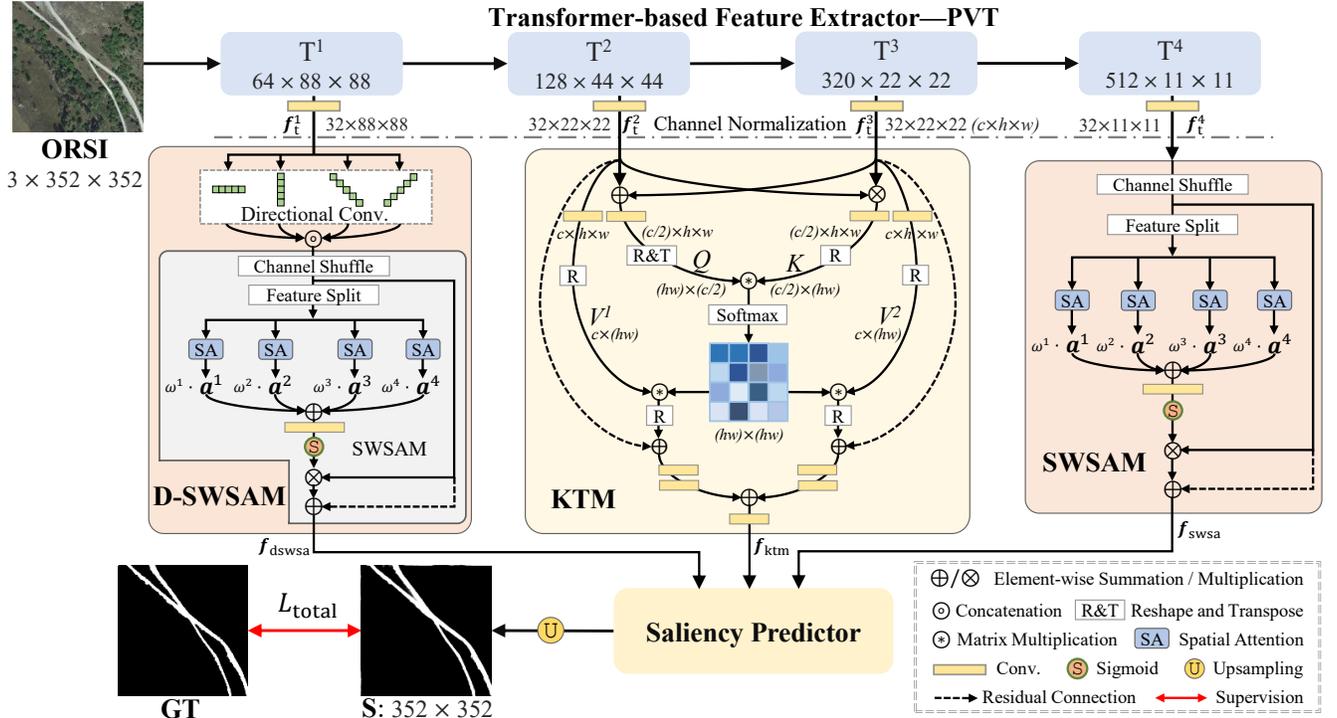
Fig. 2. Pipeline of the proposed transformer-driven GeleNet, which consists of a feature extractor, three modules, and a saliency predictor. First, we adopt a transformer-based feature extractor PVT-v2-b2 [33] to extract four-level basic feature embeddings with global long-range dependencies. Then, we employ the *Direction-aware Shuffle Weighted Spatial Attention Module* (D-SWSAM), the *Knowledge Transfer Module*, and the variant of D-SWSAM (*i.e.*, SWSAM) to deal with the corresponding features, respectively. Specifically, in D-SWSAM, we perform four directional convolutions with different directions (*i.e.*, horizontal, vertical, leading diagonal, and reverse diagonal) on the lowest-level features to extract specific orientation information, and then use SWSAM to outline the details regions. We also adopt SWSAM to enhance the location of salient objects in the highest-level features. In KTM, we model the contextual correlation knowledge of two types of combinations (*i.e.*, product and sum) of two middle-level features, and transfer the knowledge to the raw features to generate more discriminative features. Finally, we use a saliency predictor to generate a saliency map from the outputs of the above modules.

subsequent processing in KTM, generating $f_t^2 \in \mathbb{R}^{32 \times 22 \times 22}$. For the lowest-level features $f_t^1$ and the highest-level features $f_t^4$, we adopt an improved spatial attention mechanism for local enhancement. According to the characteristics of features at different levels, we adopt D-SWSAM for $f_t^1$ to extract orientation information and achieve local detail enhancement, generating $f_{\text{dswsa}}$. While we adopt SWSAM for $f_t^4$ to achieve local location enhancement, generating $f_{\text{swsa}}$. Moreover, we adopt KTM to activate cross-level contextual interactions of $f_t^2$ and $f_t^3$, generating discriminative features $f_{\text{ktm}}$. Taking advantage of PVT and these three novel modules, we infer salient objects in the saliency predictor, which is a variant of the effective partial decoder [48].

### B. Direction-aware Shuffle Weighted Spatial Attention Module

Since the basic features extracted by PVT is with global long-range dependencies, we want to explore their local enhancements to complement their global information and adapt to complex scenes in ORSIs. To be precise, we hope to consistently highlight salient regions in features across different channels, which is important for SOD. Traditional spatial attention [41] is known to be effective way to achieve this goal, however, it generates the spatial attention map in a global manner. Specifically, it performs global max pooling and global average pooling on all channels, which may produce an insufficient spatial attention map. Differently, SGE [44],

as a grouping attention, splits features into several subsets and generates a specific spatial attention map from each sub-feature for individual enhancement. While considering only the attention of each sub-feature, SGE ignores the consistency of attention between different sub-features, resulting in the lack of consistency in the group-enhanced features, which is not friendly to SOD. Inspired by [41], [44], we propose an effective grouping spatial attention mechanism for SOD, *i.e.*, the *Shuffle Weighted Spatial Attention Module* (SWSAM), which first generates the local spatial attention map from each sub-feature, and then adopts the weighted fusion operation to produce the final spatial attention map for consistent enhancement.

In addition, salient objects in ORSIs usually have various orientations, as shown in Fig. 1 and Fig. 2, which often bring troubles to existing methods using the traditional convolutions. To solve this issue, we specifically introduce directional convolutions with different directions [49] into SWSAM, and propose D-SWSAM to explicitly extract orientation information of salient objects and achieve local enhancement. Moreover, we arrange D-SWSAM to deal with $f_t^1$. The detailed structure of D-SWSAM is presented in the left part of Fig. 2. In the following, we elaborate D-SWSAM in three parts, *i.e.*, the directional convolution unit, the channel shuffle and feature split, and the weighted spatial attention, of which the latter two parts constitute SWSAM.

*1) Directional Convolution Unit.* The directional convolution unit takes into account the four basic directions, and is composed of four directional convolution layers with horizontal (h), vertical (v), leading diagonal (ld), and reverse diagonal (rd) directions [49]. We parallelize these four directional convolution layers to simultaneously mine different orientation information of $\boldsymbol{f}_{\mathrm{t}}^{1}$, and concatenate the results for information integration. We formulate the above process as follows:

$$\boldsymbol{f}_{\mathrm{ori}} = \mathrm{conv}_{\mathrm{h}}(\boldsymbol{f}_{\mathrm{t}}^{1}) \odot \mathrm{conv}_{\mathrm{v}}(\boldsymbol{f}_{\mathrm{t}}^{1}) \odot \mathrm{conv}_{\mathrm{ld}}(\boldsymbol{f}_{\mathrm{t}}^{1}) \odot \mathrm{conv}_{\mathrm{rd}}(\boldsymbol{f}_{\mathrm{t}}^{1}), \tag{1}$$

where $\boldsymbol{f}_{\mathrm{ori}} \in \mathbb{R}^{32 \times 88 \times 88}$ denotes the output orientation features, $\odot$ is the concatenation, $\mathrm{conv}_{\mathrm{x}}(\cdot)$ is the directional convolution layer with the direction $\mathrm{x} \in \{\mathrm{h}, \mathrm{v}, \mathrm{ld}, \mathrm{rd}\}$. Considering the input feature size and computational efficiency, here we set the kernel size and the output channel of four directional convolutions to 5 and 8, respectively. To show the extracted orientation information intuitively, we expand $\boldsymbol{f}_{\mathrm{ori}}$ across channel as $[\boldsymbol{f}_{\mathrm{h}}^{1}, ..., \boldsymbol{f}_{\mathrm{h}}^{8}, \boldsymbol{f}_{\mathrm{v}}^{1}, ..., \boldsymbol{f}_{\mathrm{v}}^{8}, \boldsymbol{f}_{\mathrm{ld}}^{1}, ..., \boldsymbol{f}_{\mathrm{ld}}^{8}, \boldsymbol{f}_{\mathrm{rd}}^{1}, ..., \boldsymbol{f}_{\mathrm{rd}}^{8}]$, where $\boldsymbol{f}_{\mathrm{x}} \in \mathbb{R}^{1 \times 88 \times 88}$ is a single-channel feature and we omit its superscript, and each directional convolution layer generates an eight-channel feature.

*2) Channel Shuffle and Feature Split.* Inspired by ShuffleNet [50] and shuffle attention [45], which shuffle features to achieve information communication along channels, we shuffle $\boldsymbol{f}_{\mathrm{ori}}$ with four groups to evenly disperse the orientation information, achieving $\boldsymbol{f}_{\mathrm{shuf}} \in \mathbb{R}^{32 \times 88 \times 88}$, which can be expanded as $[\boldsymbol{f}_{\mathrm{h}}^{1}, \boldsymbol{f}_{\mathrm{v}}^{1}, \boldsymbol{f}_{\mathrm{ld}}^{1}, \boldsymbol{f}_{\mathrm{rd}}^{1}, ..., \boldsymbol{f}_{\mathrm{h}}^{4}, \boldsymbol{f}_{\mathrm{v}}^{4}, \boldsymbol{f}_{\mathrm{ld}}^{4}, \boldsymbol{f}_{\mathrm{rd}}^{4}, ..., \boldsymbol{f}_{\mathrm{h}}^{8}, \boldsymbol{f}_{\mathrm{v}}^{8}, \boldsymbol{f}_{\mathrm{ld}}^{8}, \boldsymbol{f}_{\mathrm{rd}}^{8}]$.

Then, we split $\boldsymbol{f}_{\mathrm{shuf}}$ into four feature subsets along channel, generating $\{\boldsymbol{f}_{\mathrm{s-shuf}}^{1}, \boldsymbol{f}_{\mathrm{s-shuf}}^{2}, \boldsymbol{f}_{\mathrm{s-shuf}}^{3}, \boldsymbol{f}_{\mathrm{s-shuf}}^{4}\} \in \mathbb{R}^{8 \times 88 \times 88}$, where $\boldsymbol{f}_{\mathrm{s-shuf}}^{n}$ ($n \in \{1, 2, 3, 4\}$) can be expanded as $[\boldsymbol{f}_{\mathrm{h}}^{2n-1}, \boldsymbol{f}_{\mathrm{v}}^{2n-1}, \boldsymbol{f}_{\mathrm{ld}}^{2n-1}, \boldsymbol{f}_{\mathrm{rd}}^{2n-1}, \boldsymbol{f}_{\mathrm{h}}^{2n}, \boldsymbol{f}_{\mathrm{v}}^{2n}, \boldsymbol{f}_{\mathrm{ld}}^{2n}, \boldsymbol{f}_{\mathrm{rd}}^{2n}]$. The above operations activate the interaction between features of different orientations, so that each sub-feature evenly contains orientation information in four directions, which is conducive to generating an accurate spatial attention map for each sub-feature.

*3) Weighted Spatial Attention.* We then apply the traditional spatial attention [41] to the above sub-features $\boldsymbol{f}_{\mathrm{s-shuf}}^{n}$, generating corresponding spatial attention maps $\boldsymbol{a}^{n} \in (0,1)^{1 \times 88 \times 88}$ as follows:

$$\boldsymbol{a}^{n} = \mathrm{SA}(\boldsymbol{f}_{\mathrm{s-shuf}}^{n}), \tag{2}$$

where $\mathrm{SA}(\cdot)$ is the spatial attention operation. These four spatial attention maps can extract salient regions in local sub-features comprehensively without neglecting salient regions in the original complete $\boldsymbol{f}_{\mathrm{ori}}$.

Next, we design a learnable attention fusion approach, that is, set a learnable parameter $w^{n} \in [0,1]$ for each spatial attention map $\boldsymbol{a}^{n}$ and aggregate them as follows:

$$\boldsymbol{a}_{\mathrm{ori}} = \mathrm{sigmoid}(\mathrm{conv}(\sum_{n=1}^{4} w^{n} \cdot \boldsymbol{a}^{n})), \tag{3}$$

where $\boldsymbol{a}_{\mathrm{ori}} \in (0,1)^{1 \times 88 \times 88}$ is the aggregated spatial attention map, $w^{n}$ is initialized as 0.25 and gradually converges to appropriate weights, $\sum_{n=1}^{4} w^{n} = 1$, $\mathrm{conv}(\cdot)$ is the normal convolution layer, and $\mathrm{sigmoid}(\cdot)$ is the sigmoid activation function. In this way, we can obtain a comprehensive and

orientation-sensitive spatial attention map $\boldsymbol{a}_{\mathrm{ori}}$. We adopt $\boldsymbol{a}_{\mathrm{ori}}$ to achieve consistent detail enhancement, generating the output feature of D-SWSAM $\boldsymbol{f}_{\mathrm{dswsa}} \in \mathbb{R}^{32 \times 88 \times 88}$ as follows:

$$\boldsymbol{f}_{\mathrm{dswsa}} = (\boldsymbol{a}_{\mathrm{ori}} \otimes \boldsymbol{f}_{\mathrm{shuf}}) \oplus \boldsymbol{f}_{\mathrm{shuf}}, \tag{4}$$

where $\otimes$ is the element-wise multiplication and $\oplus$ is the element-wise summation. Notably, here we perform detail enhancement on $\boldsymbol{f}_{\mathrm{shuf}}$ rather than $\boldsymbol{f}_{\mathrm{ori}}$, which continues to maintain valid channel interactions.

*4) Applying SWSAM for Location Enhancement.* As shown in Fig. 2, instead of D-SWSAM, we apply SWSAM on the highest-level features $\boldsymbol{f}_{\mathrm{t}}^{4}$ for location enhancement. This is because $\boldsymbol{f}_{\mathrm{t}}^{4}$ mainly contains location information, rather than detail information such as orientation information and texture information, which means that the directional convolution unit in D-SWSAM is superfluous. Therefore, we abandon this unit. In addition, $\boldsymbol{f}_{\mathrm{t}}^{4}$ is extracted using PVT which focuses on modeling the long-range dependencies between feature patches and inevitably ignores feature interactions between channels. So we maintain the channel shuffle operation in SWSAM to explicitly increase the channel interaction. In this way, we can obtain the output feature of SWSAM $\boldsymbol{f}_{\mathrm{swsa}} \in \mathbb{R}^{32 \times 11 \times 11}$.

In summary, our D-SWSAM and SWSAM are designed according to specific characteristics of extracted global features of ORSIs to better enhance local interactions. We believe our D-SWSAM can effectively assist GeleNet to adapt to salient objects with various orientations in ORSIs, and our SWSAM can assist GeleNet to accurately locate all salient objects in ORSIs.

### C. Knowledge Transfer Module

For the lowest-level and highest-level features, we design special modules to process them to achieve local interactions according to their respective characteristics. However, it is insufficient to consider only local enhancement, we enhance cross-level contextual interactions on two middle-level features (*i.e.*, $\boldsymbol{f}_{\mathrm{t}}^{2}$ and $\boldsymbol{f}_{\mathrm{t}}^{3}$) to explore the discriminative information of salient objects. Inspired by the self-attention mechanism [28], [42], we propose a knowledge transfer module to achieve the goal. The detailed structure of KTM is presented in the middle part of Fig. 2. In the following, we introduce the two KTM components, *i.e.*, the contextual correlation knowledge modeling and the knowledge transfer.

*1) Contextual Correlation Knowledge Modeling.* In SOD, the product of two features can reveal the significant information co-existing in both features, which is conducive to collaboratively identifying objects. The sum of two features can comprehensively capture the information contained in both features without omission, which is conducive to elaborating objects. In particular for our framework, the product and sum of $\boldsymbol{f}_{\mathrm{t}}^{2}$ and $\boldsymbol{f}_{\mathrm{t}}^{3}$ are complementary to a certain extent. Therefore, we adopt self-attention [28], [42] to model the contextual correlation knowledge between the product and sum of $\boldsymbol{f}_{\mathrm{t}}^{2}$ and $\boldsymbol{f}_{\mathrm{t}}^{3}$.

As stated in Sec. III-A, we have unified the size of $\boldsymbol{f}_{\mathrm{t}}^{2}$ and $\boldsymbol{f}_{\mathrm{t}}^{3}$ to $32 \times 22 \times 22$. For convenience, we denote the size of $\boldsymbol{f}_{\mathrm{t}}^{2}$ and $\boldsymbol{f}_{\mathrm{t}}^{3}$ to $c \times h \times w$, as shown in Fig. 2. Here, we denote

the product and sum of $\boldsymbol{f}_t^2$ and $\boldsymbol{f}_t^3$ as $\boldsymbol{f}_{\text{pro}} \in \mathbb{R}^{c \times h \times w}$ and $\boldsymbol{f}_{\text{sum}} \in \mathbb{R}^{c \times h \times w}$, respectively. As the KTM illustrated in Fig. 2, to reduce the computational cost, we perform a convolution layer with the channel number of $c/2$ on $\boldsymbol{f}_{\text{pro}}$ and $\boldsymbol{f}_{\text{sum}}$ to generate two new features $\{\tilde{\boldsymbol{f}}_{\text{pro}}, \tilde{\boldsymbol{f}}_{\text{sum}}\} \in \mathbb{R}^{(c/2) \times h \times w}$. Then, we reshape and transpose $\tilde{\boldsymbol{f}}_{\text{sum}}$ to obtain $\boldsymbol{f}_Q \in \mathbb{R}^{(hw) \times (c/2)}$, and reshape $\tilde{\boldsymbol{f}}_{\text{pro}}$ to obtain $\boldsymbol{f}_K \in \mathbb{R}^{(c/2) \times (hw)}$. After that we model the contextual correlation knowledge $\mathbf{C} \in \mathbb{R}^{(hw) \times (hw)}$ between $\boldsymbol{f}_Q$ and $\boldsymbol{f}_K$ as follows:

$$\mathbf{C} = \text{softmax}(\boldsymbol{f}_Q \circledast \boldsymbol{f}_K), \tag{5}$$

where $\text{softmax}(\cdot)$ is the softmax activation function and $\circledast$ is the matrix multiplication. In this way, we model the pixel-to-pixel dependencies between the co-existing significant information of $\boldsymbol{f}_{\text{pro}}$ and the comprehensive information of $\boldsymbol{f}_{\text{sum}}$, which are effective to avoid missing salient regions/objects in ORSIs.

*2) Knowledge Transfer.* Meanwhile, we use a convolution layer on $\boldsymbol{f}_t^2$ and $\boldsymbol{f}_t^3$ to generate two new features $\{\tilde{\boldsymbol{f}}_t^2, \tilde{\boldsymbol{f}}_t^3\} \in \mathbb{R}^{c \times h \times w}$, and then reshape them to obtain $\{\boldsymbol{f}_{V^1}, \boldsymbol{f}_{V^2}\} \in \mathbb{R}^{c \times (hw)}$. After that we transfer the modeled knowledge $\mathbf{C}$ to $\boldsymbol{f}_{V^1}$ and $\boldsymbol{f}_{V^2}$ to generate the informative transferred features $\{\boldsymbol{f}_{\text{tsf}}^1, \boldsymbol{f}_{\text{tsf}}^2\} \in \mathbb{R}^{c \times h \times w}$ as follows:

$$\begin{aligned}
\boldsymbol{f}_{\text{tsf}}^1 &= \text{R}(\boldsymbol{f}_{V^1} \circledast \text{T}(\mathbf{C})), \\
\boldsymbol{f}_{\text{tsf}}^2 &= \text{R}(\boldsymbol{f}_{V^2} \circledast \text{T}(\mathbf{C})),
\end{aligned} \tag{6}$$

where $\text{R}(\cdot)$ and $\text{T}(\cdot)$ mean reshape and transpose, respectively. Following [42], we introduce a trainable weight to adaptively fuse $\boldsymbol{f}_{\text{tsf}}^1$ and raw $\boldsymbol{f}_t^2$ through residual connection, and do the same for $\boldsymbol{f}_{\text{tsf}}^2$ and raw $\boldsymbol{f}_t^3$, generating $\{\tilde{\boldsymbol{f}}_{\text{tsf}}^1, \tilde{\boldsymbol{f}}_{\text{tsf}}^2\} \in \mathbb{R}^{c \times h \times w}$. Finally, we adopt an element-wise summation and a convolution layer to integrate the cross-level $\tilde{\boldsymbol{f}}_{\text{tsf}}^1$ and $\tilde{\boldsymbol{f}}_{\text{tsf}}^2$, generating the discriminative output feature of KTM $\boldsymbol{f}_{\text{ktm}} \in \mathbb{R}^{c \times h \times w}$.

In summary, $\boldsymbol{f}_{\text{ktm}}$ inherits the properties of two combinations of $\boldsymbol{f}_t^2$ and $\boldsymbol{f}_t^3$, so it has the ability to simultaneously identify and elaborate salient objects. In addition, compared to $\boldsymbol{f}_{\text{dswsa}}$ and $\boldsymbol{f}_{\text{swsa}}$, $\boldsymbol{f}_{\text{ktm}}$ is more contextual, which is beneficial for our GeleNet to combine with local enhanced features (*i.e.*, $\boldsymbol{f}_{\text{dswsa}}$ and $\boldsymbol{f}_{\text{swsa}}$) for better salient object inference.

### D. Saliency Predictor

To make better use of the informative output features of D-SWSAM, KTM and SWSAM, *i.e.*, $\boldsymbol{f}_{\text{dswsa}}$, $\boldsymbol{f}_{\text{ktm}}$ and $\boldsymbol{f}_{\text{swsa}}$, we adopt the effective partial decoder [48] as our saliency predictor to generate the saliency map. Normally, the resolutions of input features in the original partial decoder are $1\times$, $2\times$, and $4\times$. However, the resolutions of input features of our saliency predictor are $32 \times 11 \times 11$ ($\boldsymbol{f}_{\text{swsa}}$), $32 \times 22 \times 22$ ($\boldsymbol{f}_{\text{ktm}}$), and $32 \times 88 \times 88$ ($\boldsymbol{f}_{\text{dswsa}}$). Therefore, we make a small modification to the original partial encoder, *i.e.*, modify the upsampling rate, to adapt to the resolutions of our input features. In this way, our saliency predictor can generate an initial saliency map $\mathbf{s} \in [0,1]^{1 \times 88 \times 88}$. We restore its resolution to the same resolution as the input ORSI by a $4\times$ upsampling operation, and obtain the final saliency map $\mathbf{S} \in [0,1]^{1 \times 352 \times 352}$.

During the training phase, we train the proposed GeleNet with a hybrid loss function [67], [68], including the

intersection-over-union (IoU) loss and the binary cross-entropy (BCE) loss. We formulate the total loss function $L_{\text{total}}$ as follows:

$$L_{\text{total}} = \ell_{iou}(\mathbf{S}, \mathbf{G}) + \ell_{bce}(\mathbf{S}, \mathbf{G}), \tag{7}$$

where $\ell_{iou}(\cdot)$ and $\ell_{bce}(\cdot)$ are IoU loss and BCE loss, respectively, and $\mathbf{G} \in \{0,1\}^{1 \times 352 \times 352}$ is the ground truth (GT).

## IV. Experiments

### A. Experimental Setup

*1) Datasets.* We conduct experiments on the ORSSD [15], EORSSD [14], and ORSI-4199 [22] datasets. The ORSSD dataset is the first public dataset for ORSI-SOD, and contains 800 images and corresponding pixel-level GTs, of which 600 images are used for training and 200 images for testing. The EORSSD dataset contains 2,000 images and corresponding GTs, of which 1,400 images are used for training and 600 images for testing. The ORSI-4199 dataset is the biggest dataset for ORSI-SOD, and contains 4,199 images and corresponding GTs, of which 2,000 images are used for training and 2,199 images for testing. Following [14], [17], [23], we train our GeleNet on the training set of each dataset and test it on the test set of each dataset.

*2) Network Implementation.* All experiments are conducted on PyTorch [69] with an NVIDIA Titan X GPU (12GB memory). To balance the effectiveness and efficiency, we adopt PVT-v2-b2 [33] as the backbone, and initialize it with the pre-trained parameters. Newly added layers are all initialized with the "Kaiming" method [70]. We adopt rotation and a combination of flipping and rotation for data augmentation, and resize the input image and GT to $352 \times 352$. Our GeleNet is trained using Adam optimizer [71] for 45 epochs with a batch size of 8 and a base learning rate of $1e^{-4}$ which will decay to 1/10 every 30 epochs.

*3) Evaluation Metrics.* We adopt some widely used evaluation metrics to quantitatively evaluate the performance of our method and all compared methods on three datasets, including S-measure ($S_\alpha$, $\alpha = 0.5$) [72], F-measure ($F_\beta$, $\beta^2 = 0.3$) [73], E-measure ($E_\xi$) [74], mean absolute error (MAE, $\mathcal{M}$), precision-recall (PR) curve, and F-measure curve. Here we adopt the evaluation tool (Matlab version)[2] for convenient evaluation.

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

We compare our GeleNet with state-of-the-art NSI-SOD and ORSI-SOD methods on the EORSSD and ORSSD datasets, including R3Net [51], PoolNet [52], EGNet [53], GCPA [54], MINet [55], ITSD [56], GateNet [57], CSNet [58], SAM-Net [59], HVPNet [60], SUCA [61], PA-KRN [62], VST [35], DPORTNet [63], DNTD [64], ICON [65] with PVT backbone, LVNet [15], DAFNet [14], SARNet [21], MJRBM [22], EMFINet [23], ERPNet [18], ACCoNet [16], CorrNet [17], MCCNet [20], and HFANet [66]. The saliency maps for the above methods are obtained from authors and public benchmarks[3],[4] [14], [15], or by running public source codes. For the

---

[2]https://github.com/MathLee/MatlabEvaluationTools
[3]https://li-chongyi.github.io/proj_optical_saliency.html
[4]https://github.com/rmcong/DAFNet_TIP20

TABLE I
QUANTITATIVE COMPARISONS WITH STATE-OF-THE-ART NSI-SOD AND ORSI-SOD METHODS ON EORSSD AND ORSSD DATASETS. ↓ INDICATES THAT THE LOWER THE BETTER, WHILE ↑ THE OPPOSITE. WE MARK THE TOP TWO RESULTS IN RED AND BLUE, RESPECTIVELY.

| Methods | Type | EORSSD [14] | | | | | | | | ORSSD [15] | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $S_\alpha \uparrow$ | $F_\beta^{max} \uparrow$ | $F_\beta^{mean} \uparrow$ | $F_\beta^{adp} \uparrow$ | $E_\xi^{max} \uparrow$ | $E_\xi^{mean} \uparrow$ | $E_\xi^{adp} \uparrow$ | $\mathcal{M} \downarrow$ | $S_\alpha \uparrow$ | $F_\beta^{max} \uparrow$ | $F_\beta^{mean} \uparrow$ | $F_\beta^{adp} \uparrow$ | $E_\xi^{max} \uparrow$ | $E_\xi^{mean} \uparrow$ | $E_\xi^{adp} \uparrow$ | $\mathcal{M} \downarrow$ |
| R3Net[51] | CN | .8184 | .7498 | .6302 | .4165 | .9483 | .8294 | .6462 | .0171 | .8141 | .7456 | .7383 | .7379 | .8913 | .8681 | .8887 | .0399 |
| PoolNet[52] | CN | .8207 | .7545 | .6406 | .4611 | .9292 | .8193 | .6836 | .0210 | .8403 | .7706 | .6999 | .6166 | .9343 | .8650 | .8124 | .0358 |
| EGNet[53] | CN | .8601 | .7880 | .6967 | .5379 | .9570 | .8775 | .7566 | .0110 | .8721 | .8332 | .7500 | .6452 | .9731 | .9013 | .8226 | .0216 |
| GCPA[54] | CN | .8869 | .8347 | .7905 | .6723 | .9524 | .9167 | .8647 | .0102 | .9026 | .8687 | .8433 | .7861 | .9509 | .9341 | .9205 | .0168 |
| MINet[55] | CN | .9040 | .8344 | .8174 | .7705 | .9442 | .9346 | .9243 | .0093 | .9040 | .8761 | .8574 | .8251 | .9545 | .9454 | .9423 | .0144 |
| ITSD[56] | CN | .9050 | .8523 | .8221 | .7421 | .9556 | .9407 | .9103 | .0106 | .9050 | .8735 | .8502 | .8068 | .9601 | .9482 | .9335 | .0165 |
| GateNet[57] | CN | .9114 | .8566 | .8228 | .7109 | .9610 | .9385 | .8909 | .0095 | .9186 | .8871 | .8679 | .8229 | .9664 | .9538 | .9428 | .0137 |
| CSNet[58] | CN | .8364 | .8341 | .7656 | .6319 | .9535 | .8929 | .8339 | .0169 | .8910 | .8790 | .8285 | .7615 | .9628 | .9171 | .9068 | .0186 |
| SAMNet[59] | CN | .8622 | .7813 | .7214 | .6114 | .9421 | .8700 | .8284 | .0132 | .8761 | .8137 | .7531 | .6843 | .9478 | .8818 | .8656 | .0217 |
| HVPNet[60] | CN | .8734 | .8036 | .7377 | .6202 | .9482 | .8721 | .8270 | .0110 | .8610 | .7938 | .7396 | .6726 | .9320 | .8717 | .8471 | .0225 |
| SUCA[61] | CN | .8988 | .8229 | .7949 | .7260 | .9520 | .9277 | .9082 | .0097 | .8989 | .8484 | .8237 | .7748 | .9584 | .9400 | .9194 | .0145 |
| PA-KRN[62] | CN | .9192 | .8639 | .8358 | .7993 | .9616 | .9536 | .9416 | .0104 | .9239 | .8890 | .8727 | .8548 | .9680 | .9620 | .9579 | .0139 |
| VST[35] | TN | .9208 | .8716 | .8263 | .7089 | .9743 | .9442 | .8941 | .0067 | .9365 | .9095 | .8817 | .8262 | .9810 | .9621 | .9466 | .0094 |
| DPORTNet[63] | CN | .8960 | .8363 | .7937 | .7545 | .9423 | .9116 | .9150 | .0150 | .8827 | .8309 | .8184 | .7970 | .9214 | .9139 | .9083 | .0220 |
| DNTD[64] | CN | .8957 | .8189 | .7962 | .7288 | .9378 | .9225 | .9047 | .0113 | .8698 | .8231 | .8020 | .7645 | .9286 | .9086 | .9081 | .0217 |
| ICON[65] | TN | .9185 | .8622 | .8371 | .8065 | .9687 | .9619 | .9497 | .0073 | .9256 | .8939 | .8671 | .8444 | .9704 | .9637 | .9554 | .0116 |
| LVNet[15] | CO | .8630 | .7794 | .7328 | .6284 | .9254 | .8801 | .8445 | .0146 | .8815 | .8263 | .7995 | .7506 | .9456 | .9259 | .9195 | .0207 |
| DAFNet[14] | CO | .9166 | .8614 | .7845 | .6427 | .9861 | .9291 | .8446 | .0060 | .9191 | .8928 | .8511 | .7876 | .9771 | .9539 | .9360 | .0113 |
| SARNet[21] | CO | .9240 | .8719 | .8541 | .8304 | .9620 | .9555 | .9536 | .0099 | .9134 | .8850 | .8619 | .8512 | .9557 | .9477 | .9464 | .0187 |
| MJRBM[22] | CO | .9197 | .8656 | .8239 | .7066 | .9646 | .9350 | .8897 | .0099 | .9204 | .8842 | .8566 | .8022 | .9623 | .9415 | .9328 | .0163 |
| EMFINet[23] | CO | .9290 | .8720 | .8486 | .7984 | .9711 | .9604 | .9501 | .0084 | .9366 | .9002 | .8856 | .8617 | .9737 | .9671 | .9663 | .0109 |
| ERPNet[18] | CO | .9210 | .8632 | .8304 | .7554 | .9603 | .9401 | .9228 | .0089 | .9254 | .8974 | .8745 | .8356 | .9710 | .9566 | .9520 | .0135 |
| ACCoNet[16] | CO | .9290 | .8837 | .8552 | .7969 | .9727 | .9653 | .9450 | .0074 | .9437 | .9149 | .8971 | .8806 | .9796 | .9754 | .9721 | .0088 |
| CorrNet[17] | CO | .9289 | .8778 | .8620 | .8311 | .9696 | .9646 | .9593 | .0083 | .9380 | .9129 | .9002 | .8875 | .9790 | .9746 | .9721 | .0098 |
| MCCNet[20] | CO | .9327 | .8904 | .8604 | .8137 | .9755 | .9685 | .9538 | .0066 | .9437 | .9155 | .9054 | .8957 | .9800 | .9758 | .9735 | .0087 |
| HFANet[66] | TO | .9380 | .8876 | .8681 | .8365 | .9740 | .9679 | .9644 | .0070 | .9399 | .9112 | .8981 | .8819 | .9770 | .9712 | .9722 | .0092 |
| **Ours-VGG** | CO | .9241 | .8721 | .8616 | .8382 | .9723 | .9636 | .9622 | .0080 | .9252 | .9023 | .8932 | .8806 | .9744 | .9651 | .9655 | .0130 |
| **Ours-Res** | CO | .9271 | .8723 | .8621 | .8481 | .9692 | .9651 | .9644 | .0071 | .9307 | .9042 | .8934 | .8826 | .9774 | .9714 | .9709 | .0098 |
| **Ours-SwinT** | TO | .9259 | .8774 | .8649 | .8528 | .9794 | .9752 | .9713 | .0055 | .9410 | .9203 | .9093 | .9038 | .9829 | .9779 | .9805 | .0080 |
| **Ours-PVT** | TO | .9376 | .8923 | .8781 | .8641 | .9828 | .9766 | .9750 | .0064 | .9469 | .9254 | .9128 | .9035 | .9860 | .9815 | .9786 | .0079 |

CN: CNN-based NSI-SOD method, TN: Transformer-based NSI-SOD method, CO: CNN-based ORSI-SOD method, TO: Transformer-based ORSI-SOD method.

ORSI-4199 dataset, we compare our GeleNet with 19 of the above 26 methods, whose saliency maps on the ORSI-4199 dataset are available, and additional five NSI-SOD methods (*i.e.*, PiCANet [75], BASNet [68], CPD [48], RAS [76], ENFNet [77]) provided by the public benchmark[5] [22]. Here, for a comprehensive comparison, in addition to GeleNet with the backbone of PVT-v2-b2 (*i.e.*, Ours-PVT), we also provide three variants of our GeleNet with backbones of VGG, ResNet, and Swin Transformer, named Ours-VGG, Ours-Res, and Ours-SwinT, respectively.

*1) Quantitative Comparison on the EORSSD and ORSSD Datasets.* We report the quantitative comparison results of our method and other 26 compared methods on the EORSSD and ORSSD datasets in Tab. I. We observe that Ours-PVT outperforms all compared methods on both datasets, except for $S_\alpha$, $E_\xi^{max}$ and $\mathcal{M}$ on the EORSSD dataset. Concretely, on the EORSSD dataset, Ours-PVT greatly surpasses the second-best method by 1.00%, 2.76%, and 1.06% in terms of $F_\beta^{mean}$, $F_\beta^{adp}$, and $E_\xi^{adp}$, respectively. In $E_\xi^{max}$ and $\mathcal{M}$, Ours-PVT is marginally lower than the best method by 0.33% and 0.0004,

[5]https://github.com/wchao1213/ORSI-SOD

respectively. On the ORSSD dataset, Ours-PVT is better than the second-best method in terms of $S_\alpha$ (0.9469 v.s. 0.9437), $F_\beta^{max}$ (0.9254 v.s. 0.9155), $E_\xi^{max}$ (0.9860 v.s. 0.9810), and $\mathcal{M}$ (0.0079 v.s. 0.0087). Notably, Ours-PVT is the only one whose $F_\beta^{adp}$ exceeds 0.9, *i.e.*, 0.9035. In addition, we plot the PR curve and F-measure curve of Ours-PVT and the compared methods on the EORSSD and ORSSD datasets in Fig. 3 (a-b). We can find that under different thresholds, Ours-PVT maintains its superiority and consistently achieves excellent performance.

Moreover, Ours-SwinT achieves competitive performance on the EORSSD dataset, and outperforms 26 compared methods in $F_\beta^{adp}$, $E_\xi^{mean}$, $E_\xi^{adp}$, and $\mathcal{M}$. Ours-SwinT ranks first out of seven metrics and second out of one metric compared to 26 compared methods on the ORSSD dataset. Since our modules are designed specifically for the global features of transformer, the performance of our two CNN-based variants, *i.e.*, Ours-VGG and Ours-Res, is inferior to that of Ours-SwinT and Ours-PVT, and is comparable to that of ERPNet, EMFINet, and CorrNet.

*2) Quantitative Comparison on the ORSI-4199 Dataset.* Due to slight differences in the comparison methods, we report
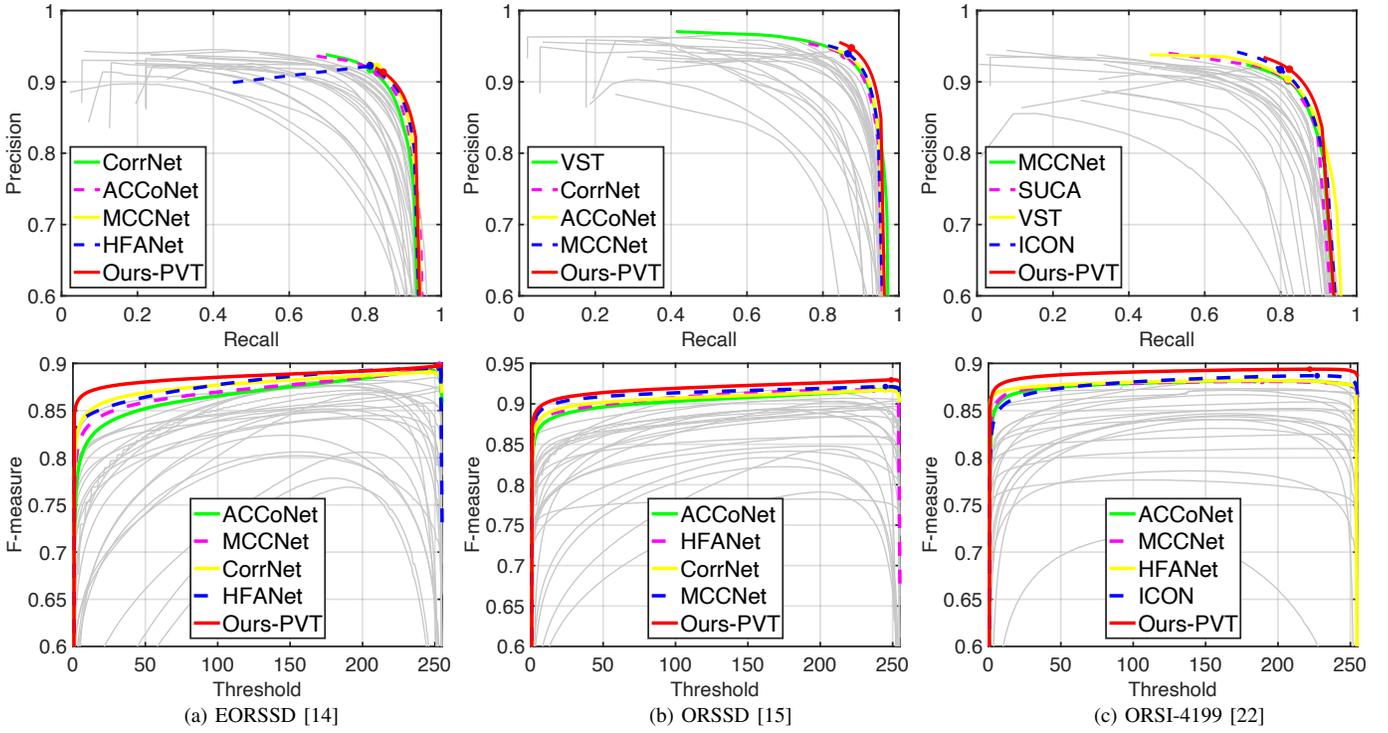
Fig. 3. Quantitative comparison on PR curve (the first row) and F-measure curve (the second row) in three datasets. We show the top five methods in different colors and the other compared methods in gray.
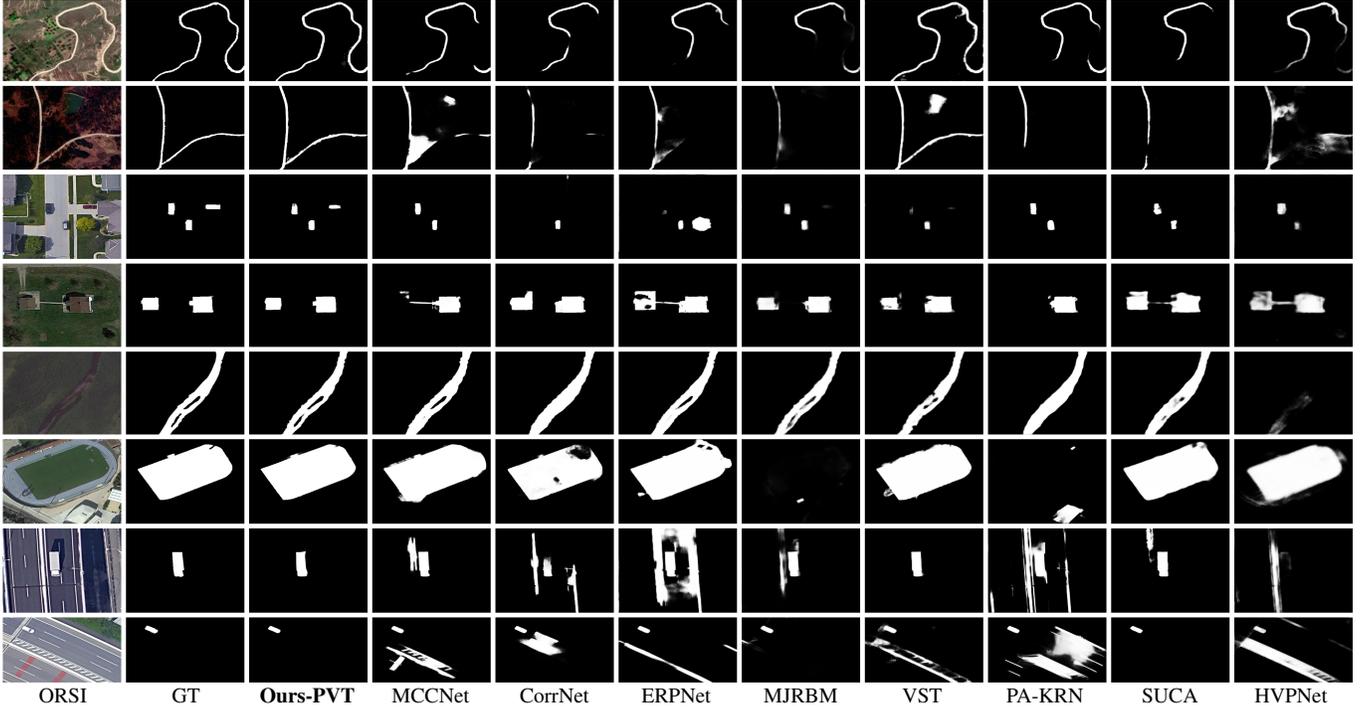


Fig. 4. Visual comparisons with eight representative state-of-the-art methods on three datasets.

the quantitative comparison results of Ours-PVT and other 24 compared methods on the ORSI-4199 dataset separately in Tab. II. The ORSI-4199 dataset is the biggest and the most challenging dataset for ORSI-SOD. The performance of Ours-PVT on this dataset is impressive, outperforming the second-best method by 0.23%~1.63% in terms of S-measure, F-measure, and E-measure. And the MAE score of Ours-PVT is only 0.0264, which is one of only three methods with the MAE score below 0.03. The advantage of Ours-PVT is easier to spot on the PR curve and F-measure curve, especially the latter one, as plotted in Fig. 3 (c). The above excellent performance on the challenging ORSI-4199 dataset strongly demonstrates the effectiveness of Ours-PVT. But to be honest, there is still a lot of room for improvement on the ORSI-4199 dataset.

TABLE II
QUANTITATIVE COMPARISONS WITH STATE-OF-THE-ART NSI-SOD AND
ORSI-SOD METHODS ON THE ORSI-4199 DATASET. WE MARK THE TOP
TWO RESULTS IN RED AND BLUE, RESPECTIVELY.

| Methods | Type | ORSI-4199 [22] | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | $S_\alpha \uparrow$ | $F_\beta^{\max} \uparrow$ | $F_\beta^{\mathrm{mean}} \uparrow$ | $F_\beta^{\mathrm{adp}} \uparrow$ | $E_\xi^{\max} \uparrow$ | $E_\xi^{\mathrm{mean}} \uparrow$ | $E_\xi^{\mathrm{adp}} \uparrow$ | $\mathcal{M} \downarrow$ |
| R3Net[18] [51] | CN | .8142 | .7847 | .7790 | .7776 | .8880 | .8722 | .8645 | .0401 |
| PiCANet[18] [75] | CN | .7114 | .6461 | .5684 | .5933 | .7946 | .6927 | .7511 | .0974 |
| PoolNet[19] [52] | CN | .8271 | .8010 | .7779 | .7382 | .8964 | .8676 | .8531 | .0541 |
| EGNet[19] [53] | CN | .8464 | .8267 | .8041 | .7650 | .9161 | .8947 | .8620 | .0440 |
| BASNet[19] [68] | CN | .8341 | .8157 | .8042 | .7810 | .9069 | .8881 | .8882 | .0454 |
| CPD[19] [48] | CN | .8476 | .8305 | .8169 | .7960 | .9168 | .9025 | .8883 | .0409 |
| RAS[20] [76] | CN | .7753 | .7343 | .7141 | .7017 | .8481 | .8133 | .8308 | .0671 |
| CSNet[20] [58] | CN | .8241 | .8124 | .7674 | .7162 | .9096 | .8586 | .8447 | .0524 |
| SAMNet[21] [59] | CN | .8409 | .8249 | .8029 | .7744 | .9186 | .8938 | .8781 | .0432 |
| HVPNet[21] [60] | CN | .8471 | .8295 | .8041 | .7652 | .9201 | .8956 | .8687 | .0419 |
| ENFNet[21] [77] | CN | .7766 | .7285 | .7177 | .7271 | .8370 | .8107 | .8235 | .0608 |
| SUCA[21] [61] | CN | .8794 | .8692 | .8590 | .8415 | .9438 | .9356 | .9186 | .0304 |
| PA-KRN[21] [62] | CN | .8491 | .8415 | .8324 | .8200 | .9280 | .9168 | .9063 | .0382 |
| VST[21] [35] | TN | .8790 | .8717 | .8524 | .7947 | .9481 | .9348 | .8997 | **.0281** |
| DPORTNet[22] [63] | CN | .8094 | .7789 | .7701 | .7554 | .8759 | .8687 | .8628 | .0569 |
| DNTD[22] [64] | CN | .8444 | .8310 | .8208 | .8065 | .9158 | .9050 | .8963 | .0425 |
| ICON[23] [65] | TN | .8752 | .8763 | .8664 | .8531 | .9521 | .9438 | .9239 | .0282 |
| MJRBM[22] [22] | CO | .8593 | .8493 | .8309 | .7995 | .9311 | .9102 | .8891 | .0374 |
| EMFINet[22] [23] | CO | .8675 | .8584 | .8479 | .8186 | .9340 | .9257 | .9136 | .0330 |
| ERPNet[22] [18] | CO | .8670 | .8553 | .8374 | .8024 | .9290 | .9149 | .9024 | .0357 |
| ACCoNet[22] [16] | CO | .8775 | .8686 | .8620 | .8581 | .9412 | .9342 | .9167 | .0314 |
| CorrNet[22] [17] | CO | .8623 | .8560 | .8513 | .8534 | .9330 | .9206 | .9142 | .0366 |
| MCCNet[22] [20] | CO | .8746 | .8690 | .8630 | .8592 | .9413 | .9348 | .9182 | .0316 |
| HFANet[22] [66] | TO | .8767 | .8700 | .8624 | .8323 | .9431 | .9336 | .9191 | .0314 |
| **Ours-VGG** | CO | .8540 | .8444 | .8374 | .8345 | .9283 | .9098 | .9086 | .0391 |
| **Ours-Res** | CO | .8670 | .8601 | .8549 | .8516 | .9383 | .9284 | .9178 | .0329 |
| **Ours-SwinT** | TO | .8828 | .8806 | .8734 | .8681 | .9537 | .9482 | .9261 | .0264 |
| **Ours-PVT** | TO | .8862 | .8842 | .8785 | .8755 | .9544 | .9478 | .9265 | .0264 |

TABLE III
ABLATION RESULTS OF EVALUATING THE INDIVIDUAL CONTRIBUTION OF
EACH MODULE IN GELENET. THE BEST ONE IN EACH COLUMN IS **BOLD**.

| No. | Baseline | D-SWSAM | KTM | SWSAM | EORSSD [14] | | |
|---|---|---|---|---|---|---|---|
| | | | | | $S_\alpha \uparrow$ | $F_\beta^{\max} \uparrow$ | $E_\xi^{\max} \uparrow$ |
| 1 | ✓ | | | | 0.9249 | 0.8717 | 0.9712 |
| 2 | ✓ | ✓ | | | 0.9305 | 0.8827 | 0.9764 |
| 3 | ✓ | | ✓ | | 0.9301 | 0.8812 | 0.9778 |
| 4 | ✓ | | | ✓ | 0.9309 | 0.8836 | 0.9775 |
| 5 | ✓ | | ✓ | ✓ | 0.9350 | 0.8872 | 0.9796 |
| 6 | ✓ | ✓ | | ✓ | 0.9328 | 0.8871 | 0.9786 |
| 7 | ✓ | ✓ | ✓ | | 0.9339 | 0.8863 | 0.9791 |
| 8 | ✓ | ✓* | ✓ | | 0.9355 | 0.8879 | 0.9798 |
| 9 | ✓ | | ✓ | ✓* | 0.9366 | 0.8911 | 0.9802 |
| 10 | ✓ | ✓ | ✓ | ✓ | **0.9376** | **0.8923** | **0.9828** |

✓*: using this module to enhance the lowest- and highest-level features.

convolution unit of D-SWSAM. The second scene contains multiple salient objects, as in the third and fourth cases of Fig. 4. Most methods only locate some of these objects and their saliency maps are relatively rough, but our method finely segments all salient objects. This is due to the precise location capability of SWSAM. The third scene contains objects with fine structure, as in the fifth and sixth cases of Fig. 4. Our method successfully delineates the same fine structure of salient objects as GTs, such as the islands in the river and the shape of the playground. The last scene is low contrast, where the color of foreground and background is similar, as in the last two cases of Fig. 4. Due to the global modeling capability of PVT and the local enhancement of proposed modules, our method accurately distinguishes white vehicles in both cases without the interference of white zebra crossings. While other methods are confused by the white zebra crossing and wrongly highlight them.

*C. Ablation Studies*

We conduct comprehensive ablation studies on the EORSSD dataset to evaluate the effectiveness of each module of our GeleNet and each component of our three modules. Accordingly, we analyze 1) the individual contribution of three modules, 2) the effectiveness of each component in D-SWSAM, 3) the rationality of the way of modeling knowledge in KTM, and 4) the effectiveness of each component in SWSAM. We conduct these ablation studies on the GeleNet with the backbone of PVT-v2-b2, and adopt the same parameter settings and dataset partitioning as in Sec. IV-A for all variants.

*1) Individual Contribution of Three Modules.* To investigate the individual contribution of the proposed three modules, *i.e.*, D-SWSAM, KTM, and SWSAM, we design various combinations of the above three modules for a total of seven variants: 1) Baseline, in which we remove all proposed modules and adopt element-wise summation to fuse $f_t^2$ and $f_t^3$, 2) Baseline+D-SWSAM, 3) Baseline+KTM, 4) Baseline+SWSAM, 5) Baseline+KTM+SWSAM, 6) Baseline+D-SWSAM+SWSAM, and 7) Baseline+D-SWSAM+KTM. We report the quantitative results in Tab. III.

Ours-SwinT consistently outperforms all compared methods in all eight metrics on the ORSI-4199 dataset, and achieves the best performance in $E_\xi^{\mathrm{mean}}$ and $\mathcal{M}$, even compared to Ours-PVT. Similar to the performance on the EORSSD and ORSSD datasets, the performance of Ours-VGG and Ours-Res is relatively average on the ORSI-4199 dataset, which further confirms that our modules is specifically designed for the global features of transformer.

In addition, Ours-PVT and two other transformer-based method (*i.e.*, VST and ICON) perform almost the best among their respective types of methods, *i.e.*, ORSI-SOD method and NSI-SOD method, on three datasets. This means that transformer-based methods can continue to drive the development of ORSI-SOD. The performance of the specialized ORSI-SOD method is generally better than that of NSI-SOD method on three datasets, which motivates us to develop better specialized ORSI-SOD solutions.

*3) Visual Comparison.* We show the visual comparison of Ours-PVT and eight representative state-of-the-art methods in Fig. 4. There are eight cases in Fig. 4 belonging to four typical and challenging ORSI scenes from three datasets. The first scene is objects with various orientations, which is unique to ORSIs, as in the first two cases of Fig. 4. We observe that only our method accurately highlights salient objects without including background. In contrast, another transformer-based method, *i.e.*, VST, incorrectly highlights some background regions, and all CNN-based methods fail to fully highlight objects. This is attributed to the directional

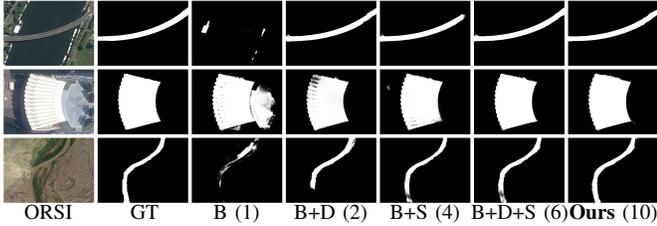| ORSI | GT | B (1) | B+D (2) | B+S (4) | B+D+S (6) | **Ours** (10) |

Fig. 5. Visual comparisons of different variants. B, D, and S are Baseline, D-SWSAM, and SWSAM, respectively. The numbers in parentheses are the ordinal numbers of these variants in Tab. III.

TABLE IV
ABLATION RESULTS OF EVALUATING THE EFFECTIVENESS OF EACH COMPONENT OF THE PROPOSED THREE MODULES. THE BEST ONE IN EACH COLUMN IS **BOLD**. D-SW. MEANS D-SWSAM, AND SWSA. MEANS SWSAM.

| | Models | EORSSD [14] | | |
|---|---|---|---|---|
| | | $S_\alpha \uparrow$ | $F_\beta^{\max} \uparrow$ | $E_\xi^{\max} \uparrow$ |
| | GeleNet (**Ours**) | **0.9376** | **0.8923** | **0.9828** |
| D-SW. | *w/o DirConv* | 0.9366 | 0.8911 | 0.9802 |
| | *w/o SWSAM* | 0.9346 | 0.8886 | 0.9800 |
| KTM | *w/ sum* | 0.9353 | 0.8886 | 0.9808 |
| | *w/ product* | 0.9334 | 0.8876 | 0.9813 |
| SWSA. | *w/o shuffle* | 0.9320 | 0.8895 | 0.9798 |
| | *w/o weights* | 0.9319 | 0.8872 | 0.9788 |

From the first four rows in Tab. III, we can find that each module can individually improve "Baseline" by around 0.5% in $S_\alpha$, around 1.0% in $F_\beta^{\max}$, and around 0.6% in $E_\xi^{\max}$, which directly proves that the proposed three modules are effective. The fifth to seventh rows of Tab. III present the performance of pairwise cooperation of modules. We can conclude that the cooperation of different modules can further improve the robustness of our method, resulting in better performance. Therefore, with all three modules working together, our full model significantly outperforms "Baseline" by 1.27% in $S_\alpha$, 2.06% in $F_\beta^{\max}$, and 1.16% in $E_\xi^{\max}$.

We provide two variants to prove the necessity of enhancing the lowest-level and highest-level features with different modules: 8) using D-SWSAM to enhance both lowest-level and highest-level features, and 9) using SWSAM to enhance both lowest-level and highest-level features. As shown in Tab. III, we observe that the performance of the above two variants is not as good as our method with different enhancements. This means that enhancing the lowest-level and highest-level features with the same module is suboptimal, and our different enhancements to the lowest-level and highest-level features are necessary.

Furthermore, we show the saliency maps for the first, second, fourth, sixth variants, and our full model in Fig. 5 to visually illustrate the role of modules. Without the help of any modules, "Baseline" performs badly, and its saliency maps have the problems of wrongly highlighting, introducing background, and incomplete highlighting. With the addition of D-SWSAM which can perceive the orientation information and perform local enhancement, the saliency maps generated

TABLE V
COMPARING THE PROPOSED SWSA WITH TWO CLASSIC ATTENTION OPERATIONS, *i.e.*, THE TRADITIONAL SPATIAL ATTENTION [41] AND SGE [44]. THE BEST ONE IN EACH COLUMN IS **BOLD**.

| Models | EORSSD [14] | | |
|---|---|---|---|
| | $S_\alpha \uparrow$ | $F_\beta^{\max} \uparrow$ | $E_\xi^{\max} \uparrow$ |
| *w/ SWSAM* (**Ours**) | **0.9376** | **0.8923** | **0.9828** |
| *w/ SA* | 0.9293 | 0.8850 | 0.9784 |
| *w/ SGE* | 0.9324 | 0.8843 | 0.9791 |

by "B+D" successfully highlight the salient objects with the correct direction (*i.e.*, the first and last cases) and suppress the background (*i.e.*, the second case). Since SWSAM is responsible for location enhancement in the highest-level features, the salient objects in the saliency maps generated by "B+S" are highlighted correctly and completely. Naturally, the combination of D-SWSAM and SWSAM, *i.e.*, "B+D+S", inherits all the advantages of the two modules. With the additional help of KTM, the saliency maps generated by our full model are visually indistinguishable from GTs. The above analysis proves that the proposed three modules are effective and play their respective functions.

*2) Effectiveness of Each Component in D-SWSAM.* D-SWSAM consists of a directional convolution unit and SWSAM. Here, we provide two variants of D-SWSAM to investigate the effectiveness of the above components: 1) removing directional convolution unit (*i.e.*, *w/o DirConv* which is the same as No.9 in Tab. III), and 2) removing SWSAM (*i.e.*, *w/o SWSAM*). As shown in the second and third rows of Tab. IV, removing either component reduces detection accuracy, which demonstrates both components are necessary for D-SWSAM. Notably, the performance of *w/o SWSAM* degrades more than that of *w/o DirConv*, indicating that SWSAM is more important in D-SWSAM.

*3) Rationality of the Way of Modeling Knowledge in KTM.* Due to the product and sum of $f_t^2$ and $f_t^3$ are complementary, we model the contextual correlation knowledge between the product and sum of $f_t^2$ and $f_t^3$ in KTM. To investigate the rationality of this way of modeling knowledge, we design two alternative modeling strategies: 1) removing product then modeling knowledge only from sum (*i.e.*, *w/ sum*), and 2) removing sum then modeling knowledge only from product (*i.e.*, *w/ product*). As shown in the fourth and fifth rows of Tab. IV, *w/ sum* and *w/ product* perform worse. As detailed in Sec. III-C, due to the complementarity between the product and sum of $f_t^2$ and $f_t^3$, the contextual correlation knowledge modeled from both is more conducive to inferring salient objects. Modeling knowledge from only one of them is suboptimal.

*4) Effectiveness of Each Component in SWSAM.* SWSAM plays an important role in our GeleNet. We use it twice in our GeleNet on the lowest-level and highest-level features. Here, we provide two variants of SWSAM to investigate the effectiveness of its components: 1) removing channel shuffle (*i.e.*, *w/o shuffle*), and 2) removing learnable parameter $w^n$ in Eq. 3 (*i.e.*, *w/o weights*). Notably, these two variants are

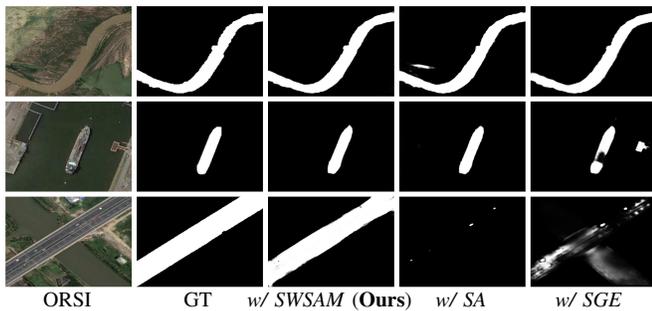| ORSI | GT | *w/ SWSAM* (**Ours**) | *w/ SA* | *w/ SGE* |

Fig. 6. Visual comparisons of different attention mechanisms.

applied in SWSAM of D-SWSAM and SWSAM of the highest level. As shown in the last two rows of Tab. IV, *w/o shuffle* and *w/o weights* are almost the worst of all variants in Tab. IV, which illustrates the importance of both operations. Actually, channel shuffle in SWSAM serves two different purposes. *w/o shuffle* lets SWSAM in D-SWSAM to generate four spatial attention maps directly from four sub-features with single direction instead of sub-features with uniform four directions, and weakens the channel communication of the global highest-level features. *w/o weights* does not take into account the differences between different spatial attention maps and simply fuses spatial attention maps. Therefore, the performance of both variants is degraded.

In addition, we compare the proposed SWSAM with two classic attention mechanisms, *i.e.*, the traditional spatial attention [41] and SGE, to further investigate the effectiveness of our SWSAM. We provide two variants: 1) replacing SWSAM with the traditional spatial attention (*i.e.*, *w/ SA*), and 2) replacing SWSAM with SGE (*i.e.*, *w/ SGE*). As reported in Tab. V, the effectiveness of these two attention mechanisms is lower than that of our SWSAM, *i.e.*, *w/ SWSAM*, for ORSI-SOD. Moreover, in Fig. 6, we show the saliency maps generated by *w/ SWSAM*, *w/ SA*, and *w/ SGE* for the visual comparison. The first case is that some background regions are similar to salient objects. Traditional spatial attention generates the attention map in a global manner (*i.e.*, from all channels), which leads to the omission of valid information and is not conducive to generating an accurate attention map. Therefore, *w/ SA* incorrectly highlights background regions similar to salient objects. The second case is the scene with the irrelevant object. Since SGE extracts specific semantics from each sub-feature and does not adopt the same consistent attention map for enhancement, *w/ SGE* mistakenly highlights the irrelevant object in the scene. The last case is the elevated highway with cars. Since the comprehensive valid information in traditional spatial attention is omitted, *w/ SA* only highlights cars on the elevated highway instead of the elevated highway. *w/ SGE* takes into account the semantics of different sub-features, so it highlights more regions than *w/ SA*, but meanwhile introduces other background regions. Differently, our SWSAM aggregates multiple attention maps generated from different sub-features in an adaptive way, resulting in a comprehensive attention map for consistent enhancement. Therefore, our *w/ SWSAM* can effectively handle the above cases.

## V. CONCLUSION

In this paper, we propose the first transformer-driven ORSI-SOD solution, namely GeleNet. GeleNet mainly follows the global-to-local paradigm, while also considering cross-level contextual interactions. GeleNet employs PVT to extract global features, SWSAM and D-SWSAM to achieve local enhancement, and KTM to activate cross-level contextual interactions. Specifically, SWSAM is an improved spatial attention module, which is responsible for location enhancement in the highest-level features. To adapt to various object orientations in ORSIs, directional convolutions are used in D-SWSAM to explicitly perceive orientation information of the lowest-level features, followed by SWSAM to achieve detail enhancement. KTM is built on the self-attention mechanism, and models the complementary information between the product and the sum of two middle-level features to generate discriminative features. The cooperation of components makes GeleNet a successful salient object detector for ORSIs. Extensive comparisons and ablation studies demonstrate the superiority of GeleNet and the effectiveness of the three proposed modules. Moreover, the proposed D-SWSAM and SWSAM can be used as plug-and-play modules for related tasks [1], [2], [77]–[79].

## REFERENCES

[1] 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, Oct. 2019.

[2] A. Borji, M.-M. Cheng, H. Jiang, and J. Li, "Salient object detection: A benchmark," *IEEE Trans. Image Process.*, vol. 24, no. 12, pp. 5706–5722, Dec. 2015.

[3] G. Li, Z. Liu, and H. Ling, "ICNet: Information conversion network for RGB-D based salient object detection," *IEEE Trans. Image Process.*, vol. 29, pp. 4873–4884, Mar. 2020.

[4] K. Gu, S. Wang, H. Yang, W. Lin, G. Zhai, X. Yang, and W. Zhang, "Saliency-guided quality assessment of screen content images," *IEEE Trans. Multimedia*, vol. 18, no. 6, pp. 1098–1110, Jun. 2016.

[5] S. Yang, Q. Jiang, W. Lin, and Y. Wang, "SGDNet: An end-to-end saliency-guided deep neural network for no-reference image quality assessment," in *Proc. ACM MM*, Oct. 2019, pp. 1383–1391.

[6] G. Li, Z. Liu, R. Shi, and W. Wei, "Constrained fixation point based segmentation via deep neural network," *Neurocomputing*, vol. 368, pp. 180–187, Nov. 2019.

[7] G. Li *et al.*, "Personal fixations-based object segmentation with object localization and boundary preservation," *IEEE Trans. Image Process.*, vol. 30, pp. 1461–1475, Jan. 2021.

[8] N. Liu, W. Zhao, L. Shao, and J. Han, "SCG: Saliency and contour guided salient instance segmentation," *IEEE Trans. Image Process.*, vol. 30, pp. 5862–5874, Jun. 2021.

[9] Q. En, L. Duan, and Z. Zhang, "Joint multisource saliency and exemplar mechanism for weakly supervised video object segmentation," *IEEE Trans. Image Process.*, vol. 30, pp. 8155–8169, Sept. 2021.

[10] G. Li, Y. Wang, Z. Liu, X. Zhang, and D. Zeng, "RGB-T semantic segmentation with location, activation, and sharpening," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 33, no. 3, pp. 1223–1235, Mar. 2023.

[11] H. Hadizadeh and I. V. Bajić, "Saliency-aware video compression," *IEEE Trans. Image Process.*, vol. 23, no. 1, pp. 19–33, Jan. 2014.

[12] W. Feng, R. Han, Q. Guo, J. Zhu, and S. Wang, "Dynamic saliency-aware regularization for correlation filter-based object tracking," *IEEE Trans. Image Process.*, vol. 28, no. 7, pp. 3232–3245, Jul. 2019.

[13] G. Li *et al.*, "Lightweight salient object detection in optical remote-sensing images via semantic matching and edge alignment," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, pp. 1–12, 2023.

[14] 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.

[15] 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, Nov. 2019.

[16] G. Li, Z. Liu, D. Zeng, W. Lin, and H. Ling, "Adjacent context coordination network for salient object detection in optical remote sensing images," *IEEE Trans. Cybern.*, vol. 53, no. 1, pp. 526–538, Jan. 2023.

[17] G. Li, Z. Liu, Z. Bai, W. Lin, and H. Ling, "Lightweight salient object detection in optical remote sensing images via feature correlation," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–12, 2022.

[18] 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.*, vol. 53, no. 1, pp. 539–552, Jan. 2023.

[19] Y. LeCun and others., "Backpropagation applied to handwritten zip code recognition," *Neural Comput.*, vol. 1, no. 4, pp. 541–551, Dec. 1989.

[20] G. Li, Z. Liu, W. Lin, and H. Ling, "Multi-content complementation network for salient object detection in optical remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–13, 2022.

[21] Z. Huang, H. Chen, B. Liu, and Z. Wang, "Semantic-guided attention refinement network for salient object detection in optical remote sensing images," *Remote Sens.*, vol. 13, no. 11, p. 2163, May 2021.

[22] Z. Tu, C. Wang, C. Li, M. Fan, H. Zhao, and B. Luo, "ORSI salient object detection via multiscale joint region and boundary model," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–13, 2022.

[23] X. Zhou, K. Shen, Z. Liu, C. Gong, J. Zhang, and C. Yan, "Edge-aware multiscale feature integration network for salient object detection in optical remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–15, 2022.

[24] C. Li *et al.*, "A parallel down-up fusion network for salient object detection in optical remote sensing images," *Neurocomputing*, vol. 415, pp. 411–120, Nov. 2020.

[25] 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. 1–11, 2022.

[26] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. ICLR*, May 2015, pp. 1–14.

[27] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE CVPR*, Jun. 2016, pp. 770–778.

[28] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. NeurIPS*, Dec. 2017, pp. 6000–6010.

[29] A. Dosovitskiy *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *Proc. ICLR*, 2021.

[30] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, Z. Jiang, F. E. H. Tay, J. Feng, and S. Yan, "Tokens-to-token ViT: Training vision transformers from scratch on imagenet," in *Proc. IEEE ICCV*, Oct. 2021, pp. 538–547.

[31] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin Transformer: Hierarchical vision transformer using shifted windows," in *Proc. IEEE ICCV*, Oct. 2021, pp. 9992–10 002.

[32] W. Wang *et al.*, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," in *Proc. IEEE ICCV*, Oct. 2021, pp. 548–558.

[33] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "PVT v2: Improved baselines with pyramid vision transformer," *Comput. Vis. Media*, vol. 8, pp. 415–424, Sept. 2022.

[34] Y. Mao, J. Zhang, Z. Wan, Y. Dai, A. Li, Y. Lv, X. Tian, D.-P. Fan, and N. Barnes, "Transformer transforms salient object detection and camouflaged object detection," *arXiv preprint arXiv:2104.10127*, 2021.

[35] N. Liu, N. Zhang, K. Wan, L. Shao, and J. Han, "Visual saliency transformer," in *Proc. IEEE ICCV*, Oct. 2021, pp. 4702–4712.

[36] Z. Liu, Y. Tan, Q. He, and Y. Xiao, "SwinNet: Swin Transformer drives edge-aware RGB-D and RGB-T salient object detection," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 32, no. 7, pp. 4486–4497, Jul. 2022.

[37] Z. Liu, Y. Wang, Z. Tu, Y. Xiao, and B. Tang, "TriTransNet: RGB-D salient object detection with a triplet transformer embedding network," in *Proc. ACM MM*, Oct. 2021, pp. 4481–4490.

[38] X. Fang, J. Zhu, X. Shao, and H. Wang, "GroupTransNet: Group transformer network for RGB-D salient object detection," *arXiv preprint arXiv:2203.10785*, 2022.

[39] Y. Su, J. Deng, R. Sun, G. Lin, and Q. Wu, "A unified transformer framework for group-based segmentation: Co-segmentation, co-saliency detection and video salient object detection," *arXiv preprint arXiv:2203.04708*, 2022.

[40] J. Hu, L. Shen, S. Albanie, G. Sun, and E. Wu, "Squeeze-and-excitation networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 8, pp. 2011–2023, Aug. 2020.

[41] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "CBAM: Convolutional block attention module," in *Proc. ECCV*, Sept. 2018, pp. 3–19.

[42] J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, and H. Lu, "Dual attention network for scene segmentation," in *Proc. IEEE CVPR*, Jun. 2019, pp. 3146–3154.

[43] Q. Wang, B. Wu, P. Zhu, P. Li, W. Zuo, and Q. Hu, "ECA-Net: Efficient channel attention for deep convolutional neural networks," in *Proc. IEEE CVPR*, Jun. 2020, pp. 11 531–11 539.

[44] X. Li, X. Hu, and J. Yang, "Spatial group-wise enhance: Improving semantic feature learning in convolutional networks," *arXiv preprint arXiv:1905.09646*, 2019.

[45] Q.-L. Zhang and Y.-B. Yang, "SA-Net: Shuffle attention for deep convolutional neural networks," in *Proc. IEEE ICASSP*, Jun. 2021, pp. 2235–2239.

[46] G. Li, Z. Liu, L. Ye, Y. Wang, and H. Ling, "Cross-modal weighting network for RGB-D salient object detection," in *Proc. ECCV*, Aug. 2020, pp. 665–681.

[47] Z. Bai, Z. Liu, G. Li, and Y. Wang, "Adaptive group-wise consistency network for co-saliency detection," *IEEE Trans. Multimedia*, vol. 25, pp. 764–776, Mar. 2023.

[48] Z. Wu, L. Su, and Q. Huang, "Cascaded partial decoder for fast and accurate salient object detection," in *Proc. IEEE CVPR*, Jun. 2019, pp. 3902–3911.

[49] J. Mei, R.-J. Li, W. Gao, and M.-M. Cheng, "CoANet: Connectivity attention network for road extraction from satellite imagery," *IEEE Trans. Image Process.*, vol. 30, pp. 8540–8552, 2021.

[50] X. Zhang, X. Zhou, M. Lin, and J. Sun, "ShuffleNet: An extremely efficient convolutional neural network for mobile devices," in *Proc. IEEE CVPR*, Jun. 2018, pp. 6848–6856.

[51] Z. Deng, X. Hu, L. Zhu, X. Xu, J. Qin, G. Han, and P.-A. Heng, "R$^3$Net: Recurrent residual refinement network for saliency detection," in *Proc. IJCAI*, Jul. 2018, pp. 684–690.

[52] J.-J. Liu, Q. Hou, M.-M. Cheng, J. Feng, and J. Jiang, "A simple pooling-based design for real-time salient object detection," in *Proc. IEEE CVPR*, Jun. 2019, pp. 3912–3921.

[53] J. 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. IEEE ICCV*, Oct. 2019, pp. 8779–8788.

[54] Z. Chen, Q. Xu, R. Cong, and Q. Huang, "Global context-aware progressive aggregation network for salient object detection," in *Proc. AAAI*, Feb. 2020, pp. 10 599–10 606.

[55] Y. Pang, X. Zhao, L. Zhang, and H. Lu, "Multi-scale interactive network for salient object detection," in *Proc. IEEE CVPR*, Jun. 2020, pp. 9410–9419.

[56] H. Zhou, X. Xie, J.-H. Lai, Z. Chen, and L. Yang, "Interactive two-stream decoder for accurate and fast saliency detection," in *Proc. IEEE CVPR*, Jun. 2020, pp. 9138–9147.

[57] X. Zhao, Y. Pang, L. Zhang, H. Lu, and L. Zhang, "Suppress and balance: A simple gated network for salient object detection," in *Proc. ECCV*, Aug. 2020, pp. 35–51.

[58] S.-H. Gao, Y.-Q. Tan, M.-M. Cheng, C. Lu, Y. Chen, and S. Yan, "Highly efficient salient object detection with 100k parameters," in *Proc. ECCV*, Aug. 2020, pp. 702–721.

[59] Y. Liu *et al.*, "SAMNet: Stereoscopically attentive multi-scale network for lightweight salient object detection," *IEEE Trans. Image Process.*, vol. 30, pp. 3804–3814, 2021.

[60] Y. Liu, Y.-C. Gu, X.-Y. Zhang, W. Wang, and M.-M. Cheng, "Lightweight salient object detection via hierarchical visual perception learning," *IEEE Trans. Cybern.*, vol. 51, no. 9, pp. 4439–4449, Sept. 2021.

[61] J. Li, Z. Pan, Q. Liu, and Z. Wang, "Stacked U-shape network with channel-wise attention for salient object detection," *IEEE Trans. Multimedia*, vol. 23, pp. 1397–1409, 2021.

[62] B. Xu, H. Liang, R. Liang, and P. Chen, "Locate globally, segment locally: A progressive architecture with knowledge review network for salient object detection," in *Proc. AAAI*, Feb. 2021, pp. 3004–3012.

[63] Y. Liu, D. Zhang, N. Liu, S. Xu, and J. Han, "Disentangled capsule routing for fast part-object relational saliency," *IEEE Trans. Image Process.*, vol. 31, pp. 6719–6732, 2022.

[64] C. Fang, H. Tian, D. Zhang, Q. Zhang, J. Han, and J. Han, "Densely nested top-down flows for salient object detection," *Sci. China Inf. Sci.*, vol. 65, no. 8, pp. 1–14, Aug. 2022.

[65] M. Zhuge, D.-P. Fan, N. Liu, D. Zhang, D. Xu, and L. Shao, "Salient object detection via integrity learning," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 3, pp. 3738–3752, Mar. 2023.

[66] Q. Wang, Y. Liu, Z. Xiong, and Y. Yuan, "Hybrid feature aligned network for salient object detection in optical remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–15, 2022.

[67] G. Li, Z. Liu, M. Chen, Z. Bai, W. Lin, and H. Ling, "Hierarchical alternate interaction network for RGB-D salient object detection," *IEEE Trans. Image Process.*, vol. 30, pp. 3528–3542, Mar. 2021.

[68] X. Qin, Z. Zhang, C. Huang, C. Gao, M. Dehghan, and M. Jagersand, "BASNet: Boundary-aware salient object detection," in *Proc. IEEE CVPR*, Jun. 2019, pp. 7479–7489.

[69] A. Paszke *et al.*, "PyTorch: An imperative style, high-performance deep learning library," in *Proc. NeurIPS*, Dec. 2019, pp. 8024–8035.

[70] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification," in *Proc. IEEE ICCV*, Dec. 2015, pp. 1026–1034.

[71] D. P. Kingma and J. L. Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, May 2015, pp. 1–15.

[72] D.-P. Fan, M.-M. Cheng, Y. Liu, T. Li, and A. Borji, "Structure-measure: A new way to evaluate foreground maps," in *Proc. IEEE ICCV*, Oct. 2017, pp. 4548–4557.

[73] R. Achanta, S. Hemami, F. Estrada, and S. Susstrunk, "Frequency-tuned salient region detection," in *Proc. IEEE CVPR*, Jun. 2009, pp. 1597–1604.

[74] D.-P. Fan, C. Gong, Y. Cao, B. Ren, M.-M. Cheng, and A. Borji, "Enhanced-alignment measure for binary foreground map evaluation," in *Proc. IJCAI*, Jul. 2018, pp. 698–704.

[75] N. Liu, J. Han, and M.-H. Yang, "PiCANet: Pixel-wise contextual attention learning for accurate saliency detection," *IEEE Trans. Image Process.*, vol. 29, pp. 6438–6451, Apr. 2020.

[76] S. Chen, X. Tan, B. Wang, H. Lu, X. Hu, and Y. Fu, "Reverse attention-based residual network for salient object detection," *IEEE Trans. Image Process.*, vol. 29, pp. 3763–3776, Jan. 2020.

[77] Z. Tu, Y. Ma, C. Li, J. Tang, and B. Luo, "Edge-guided non-local fully convolutional network for salient object detection," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 31, no. 2, pp. 582–593, Feb. 2021.

[78] Z. Tu, Z. Li, C. Li, and J. Tang, "Weakly alignment-free RGBT salient object detection with deep correlation network," *IEEE Trans. Image Process.*, vol. 31, pp. 3752–3764, May 2022.

[79] J. Tang, D. Fan, X. Wang, Z. Tu, and C. Li, "RGBT salient object detection: Benchmark and a novel cooperative ranking approach," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 30, no. 12, pp. 4421–4433, Dec. 2020.