# Position-Aware Relation Learning for RGB-Thermal Salient Object Detection

Heng Zhou, *Graduate Student Member, IEEE,* Chunna Tian*, Zhenxi Zhang, *Graduate Student Member, IEEE,* Chengyang Li, *Graduate Student Member, IEEE,* Yuxuan Ding, Yongqiang Xie*, and Zhongbo Li,

*Abstract*—RGB-Thermal salient object detection (SOD) combines two spectra to segment visually conspicuous regions in images. Most existing methods use boundary maps to learn the sharp boundary. These methods ignore the interactions between isolated boundary pixels and other confident pixels, leading to sub-optimal performance. To address this problem, we propose a position-aware relation learning network (PRLNet) for RGB-T SOD based on swin transformer. PRLNet explores the distance and direction relationships between pixels to strengthen intra-class compactness and inter-class separation, generating salient object masks with clear boundaries and homogeneous regions. Specifically, we develop a novel signed distance map auxiliary module (SDMAM) to improve encoder feature representation, which takes into account the distance relation of different pixels in boundary neighborhoods. Then, we design a feature refinement approach with directional field (FRDF), which rectifies features of boundary neighborhood by exploiting the features inside salient objects. FRDF utilizes the directional information between object pixels to effectively enhance the intra-class compactness of salient regions. In addition, we constitute a pure transformer encoder-decoder network to enhance multispectral feature representation for RGB-T SOD. Finally, we conduct quantitative and qualitative experiments on three public benchmark datasets. The results demonstrate that our proposed method outperforms the state-of-the-art methods.

*Index Terms*—Salient object detection, RGB-Thermal images, swin transformer, position-aware relation learning.

## I. INTRODUCTION

SALIENT object detection (SOD) is to segment the main conspicuous objects in the image at the pixel level by simulating the human visual system. In applications of image quality assessment [1], [2], image editing [3], [4], person re-identification [5] and robotics [6], [7], SOD extracts informative objects in images to help scene analysis and understanding. Traditional SOD methods mainly use low-level features and certain priors, such as color contrast and background priors, to detect targets [8].

In recent years, CNN-based SOD methods [9]–[12] have shown advantages over traditional hand-crafted feature-based methods in terms of model accuracy and generalization. The application of SOD is also extended from visible light images to multispectral ones [13]. Thermal sensors rely on the thermal radiation of the object to generate images, which are not easily affected by environmental conditions, such as weather, illumination, *etc*. [14]. For example, the quality of thermal images is noticeably better than RGB images in low illumination. RGB-T image pairs have both the radiometric intensity of
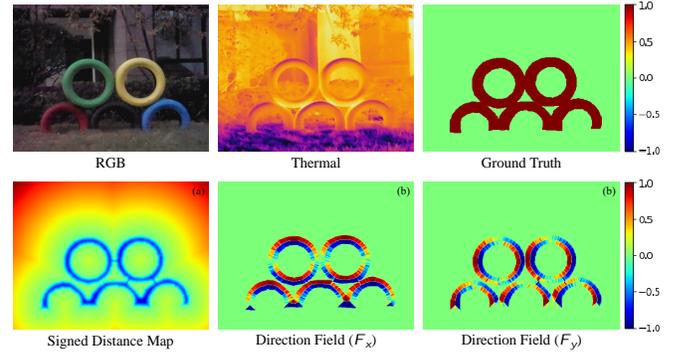
Fig. 1. RGB-T SOD with position-aware relation learning. **(a)** The signed distance map (SDM) calculates the distance from a pixel to the nearest boundary, where the sign indicates that the pixel is inside (+) or outside (-) the salient object. The zero level set is the boundary of the salient object. **(b)** The directional field $(\mathcal{F}_x, \mathcal{F}_y)$ of a salient object points from the pixel to its nearest boundary pixel.

infrared and the detail information of visible light. Compared with single RGB images, RGB-T multispectral fusion can generate discriminative and robust saliency features [15], [16]. Therefore, the RGB-T SOD method achieves a more robust generalization performance in real-world scenes.

To obtain accurate salient object masks, many CNN-based models [12], [17], [18] focus on generating clear contours by learning edge maps. The consistency between pixels mainly includes intra-class compactness and inter-class separation. However, these methods ignore the relation learning between boundary pixels and region pixels (target and background regions), resulting in unsatisfactory results. In this paper, we propose a position-aware relation learning network (PRLNet) to model the distance relations and direction relations between pixels, which enhances the intra-class compactness and inter-class separability of features.

The relative distance information between pixels can effectively alleviate the misprediction of salient pixels [19]. Inspired by the level set method [20], [21], the signed distance map (SDM) models the distance relation between region pixels (target and background regions) and boundary pixels. As shown in Fig. 1, SDM provides object boundary interaction information which is the distance between the foreground-background region and the boundary. SDM prediction task can be served as an auxiliary task to guide the feature extraction of the encoder. Different from multi-task learning of SDM in decoder [22]–[24], we propose the SDM auxiliary module (SDMAM) to enhance the boundary-awareness of the encoder. SDMAM assists the encoder to learn the relative distance

between the region pixels and the boundary, and increases the inter-class separation of the pixels.

Not only the distance relationship, but also the directional relationship between salient pixels is crucial in position-aware relationship learning. Fig. 1 shows the visualization of the horizontal and vertical directions of the direction field [25]. As illustrated in Fig. 1, the direction field can simply yet efficiently represent the directional information between intra-class pixels and boundary pixels. To strengthen the compactness of intra-class pixels, we propose a feature refinement approach with direction field (FRDF) to rectify the output feature maps of the decoder. FRDF exploits the direction relationship between saliency pixels to effectively reinforce the compactness of intra-class features. Meanwhile, we design a novel direction-aware loss function to improve the smoothing loss [26], [27], which guides the model to generate homogeneous regions and sharp boundaries.

In addition, we take full advantage of transformer [28] in modeling long-range contexts to generate cross-spectral robust RGB-T features. Swin transformer [29] adopts a hierarchical architecture to effectively solve the visual multi-scale problem and reduce the computational complexity, achieving state-of-the-art (SOTA) performance in semantic segmentation and instance segmentation [29]–[31]. In our paper, we use swin transformer as the backbone network for RGB-T feature extraction. At the same time, we design a novel patch separating layer to upsample the encoder features, which build a swin transformer decoder. At the same time, the reverse swin transformer is designed to decode the RGB-T patches. Finally, we propose a position-aware relation learning network (PRLNet) based on pure transformer for RGB-T SOD.

In summary, the main contributions of this paper are as follows.

- We propose the novel PRLNet with pure swin transformer to generate salient object masks with clear boundaries and homogeneous regions by learning the distance and direction relationships between different pixels.
- Specifically, the SDM auxiliary module is suggested to learn the distance relation of each pixel to the boundary, enhancing boundary-based inter-class (foreground-background) separation.
- In order to strengthen the intra-class compactness, we design a feature refinement approach with direction field (FRDF) and direction-aware smoothness loss.
  The features close to the boundary are refined by utilizing the internal features of salient objects.
- The qualitative and quantitative experimental results on three public benchmark datasets demonstrate that our proposed model outperforms the state-of-the-art models.

The rest of this paper is organized as follows. Section II overviews the existing methods mainly on RGB and RGB-T SOD and swin transformer. In Section III, we introduce our proposed position-aware relation learning network for RGB-T SOD. Extensive experiments and visualization results on the three benchmark datasets are given in Section IV. Finally, we conclude our work in Section V.

## II. RELATED WORKS

In this section, we review the previous SOD methods for RGB and RGB-T images. Meanwhile, related works about swin transformer are also included in this section.

### A. RGB Salient Object Detection

Recently, most CNN-based SOD methods adopt a fully convolutional network (FCN) structure [32], [33]. To improve the accuracy of prediction results, multi-level feature fusion [34]–[36] and multi-task learning [17], [37] have been widely studied. *Deng et al.* [35] use the low-level and high-level features of FCN to learn residuals between intermediate saliency predictions and ground truth for refining saliency maps. *Wu et al.* [38] propose a cascaded partial decoder (CPD) that discards large-resolution features in shallow layers for acceleration, and fuses features in deep layers to obtain accurate saliency maps. *Liu et al.* [39] present pool-based modules to progressively refine features at multiple scales producing detailed results. The boundary prediction task [12] captures accurate boundary information of salient objects. *Qin et al.* [17] design a hybrid loss for predicting the boundaries of salient objects. However, boundary supervision lacks the consideration of the interaction between boundary pixels and target pixels. Inspired by the level set method [20], [21], we develop a novel signed distance map auxiliary module (SDMAM) to improve encoder features. SDMAM takes into account the distance relation of pixels in boundary neighborhoods. The distance relationship between foreground-background region pixels and boundary pixels can effectively enhance the inter-class separability of features.

### B. RGB-T Salient Object Detection

Compared to RGB images, RGB-T images offer more information of salient objects [40]. In recent years, synergistic SOD between thermal and visible images has been widely studied [41]–[43]. The dual encoders extract RGB-T features respectively, and the decoder outputs the salient prediction results [44]. The RGB-T SOD methods take full advantage of the complementary capabilities between multimodal sensors to generate cross-modal robust fusion features [45]–[47]. *Tu et al.* [48] suggest a collaborative graph learning algorithm that uses superpixels as graph nodes to learn RGB-T node saliency. *Zhang et al.* [9] transform multi-spectral SOD into a CNN feature fusion problem, and propose to capture semantic information and visual details of RGB-T at different depths by fusing multi-level CNN features. *Tu et al.* [49] exploit the complementarity of different modalities of image content and multiple types of cues to extract multi-level multimodal features. *Zhou et al.* [50] propose an effective and consistent feature fusion network that combines features of different levels through a multi-level consistent fusion module to obtain complementary information. In this paper, in order to handle long-range dependencies between RGB-T, we develop a cross-spectra fusion transformer. Furthermore, we propose a feature refinement approach with direction field (FRDF) to enhance the intra-class compactness of salient objects. FRDF exploits feature far from the boundary to refine the features of pixels close to the boundary.
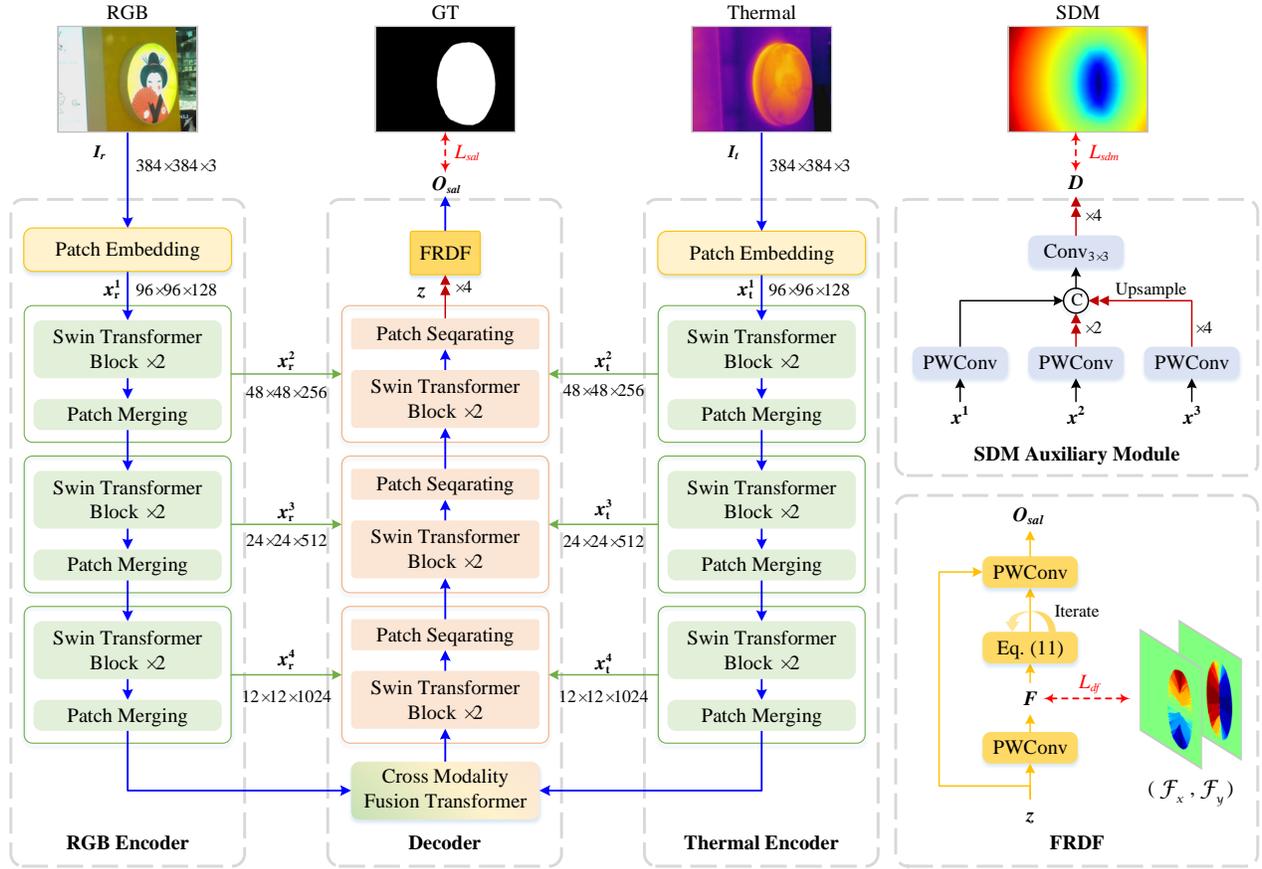
Fig. 2. The framework of our proposed PRLNet. Our network consists of four main parts, namely dual-stream encoders, RGB-T decoder, SDM auxiliary module and a feature refinement approach with directional fields (FRDF). First, multiscale features of RGB-T are extracted by a dual-stream swin transformer encoder. (Sec. III-A). Then, we construct SDMAM for encoders to learn the distance relationship between regional pixels and boundary pixels (Sec. III-B). Next, the reverse swin transformer decoder aggregates the complementarity between different levels of RGB-T features (Sec. III-C), where © denotes concatenation. In addition, FRDF is further designed by exploring direction information between salient pixels to strengthen the intra-class compactness of the salient features (Sec. III-D). Finally, we propose a novel position-aware relation learning loss to generate object masks with clear boundaries and homogeneous regions (Sec. III-E).

## C. Swin Transformer

Compared with CNN, transformer has an advantage in modeling long-range dependencies [28], [51]. ViT [52] and DETR [53] apply transformer to computer vision tasks and achieve promising performance. However, the computational complexity of the transformer is proportional to the square of the image size. To handle high-resolution images, swin transformer [29] introduces the hierarchical structure commonly used in CNN and achieves SOTA results on dense prediction tasks. Swin transformer gradually becomes a powerful general backbone network for SOD [54]. *Liu et al.* [55] propose a cross-modal fusion network based on the swin transformer for RGB-T SOD, bridging the gap between two modalities through an attention mechanism. *Zhu et al.* [56] encode multi-scale features via the swin transformer in a coarse-to-fine manner to learn salient region feature representations. In this paper, swin transformer block is used as the backbone for both the encoder stage and decoder stage. Specifically, we employ dual-swin transformer encoders to extract multi-scale features from RGB and thermal images, respectively. Referring to the patch merging layer, we design a patch separating layer to decode RGB-T hierarchical features and generate robust

results with multispectral complementarity.

## III. PRLNET

In this section, we elaborate on our proposed PRLNet for RGB-T SOD with swin transformer. The overall architecture is illustrated in Fig. 2, which consists of four main parts: Dual-stream encoders for both RGB-T images, a decoder for pixel-by-pixel prediction, a SDM auxiliary module (SDMAM) and a feature refinement approach with directional fields (FRDF). They are simultaneously optimized during the training process.

As shown in Fig. 2, PRLNet takes the RGB-T image pair as input, and segments the precise mask of the salient objects. We first use the dual-stream swin transformer encoder to generate multi-scale features of RGB and thermal images (Sec. III-A). Then, to improve the boundary perception of the encoder, we introduce SDMAM to learn the distance relationship between regional pixels and boundary pixels. SDMAM enhances the separability of foreground-background features (Sec. III-B). Next, we design a patch separating layer and construct an inverse swin transformer, which aggregates different levels of RGB-T features (Sec. III-C). To facilitate the robust cross-spectral features from the decoder, we further refine them with

the direction information between salient pixels to strengthen the intra-class compactness of the salient features (Sec. III-D). Finally, benefiting from the effective learning of position relations, we present a position-aware relation learning loss function to generate object masks with clear boundaries and homogeneous regions (Sec. III-E). The pipeline of PRLNet is illustrated in Algorithm 1.

### A. Dual-stream Swin Transformer Encoder

Swin Transformer introduces hierarchical feature mapping and shifted window attention, which has both the advantages of transformer and CNN structure [29]. We employ two Swin Transformers to extract efficient features for RGB-T image pairs, respectively. Concretely, the images are first divided into $4 \times 4$ patches and then input to the patch embedding layer, which is a $4 \times 4$ convolution with stride 4. Next, as shown in Fig. 2, RGB-T salient features are extracted by three swin transformer layers (ST), consisting of swin transformer block (STB) and patch merging layer (PM). That is,

$$
\begin{aligned}
\mathbf{R} &= \{\mathbf{x}_r^i\}_{i=1}^4 = \mathrm{ST}_r\left(\boldsymbol{I}_r\right), \\
\mathbf{T} &= \{\mathbf{x}_t^i\}_{i=1}^4 = \mathrm{ST}_t\left(\boldsymbol{I}_t\right).
\end{aligned} \tag{1}
$$

The dual encoder outputs the hierarchical representation $\mathbf{R}$ and $\mathbf{T}$, where $\mathbf{x}_r^i$ and $\mathbf{x}_t^i$ denote the $i$-th layer features of the RGB and thermal encoder, respectively. $\boldsymbol{I}_r$ and $\boldsymbol{I}_t$ indicate the RGB and thermal images, which are the input of th encoder. The bold symbols indicate the matrix. The $\mathrm{ST}_r$ and $\mathrm{ST}_t$ functions are composed of STB and PM, and represent the standard swin transformer backbones in RGB and thermal branches, respectively.

Different from ViT block [52], STB replaces the multihead attention mechanism (MSA) of ViT with window-based MSA (W-MSA) and shifted window-based MSA (SW-MSA). More formally, STB is defended as

$$
\begin{aligned}
\hat{z}^l &= \text{W-MSA}\left(\mathrm{LN}\left(z^{l-1}\right)\right) + z^{l-1}, \\
z^l &= \mathrm{MLP}\left(\mathrm{LN}\left(\hat{z}^l\right)\right) + \hat{z}^l, \\
\hat{z}^{l+1} &= \text{SW-MSA}\left(\mathrm{LN}\left(z^l\right)\right) + z^l, \\
z^{l+1} &= \mathrm{MLP}\left(\mathrm{LN}\left(\hat{z}^{l+1}\right)\right) + \hat{z}^{l+1}.
\end{aligned} \tag{2}
$$

where $\hat{z}^l$ and $z^l$ denote the output feature of the (S)W-MSA module and the MLP module for block $l$, respectively. Fig. 3 demonstrates the detailed architecture of STB. STB uses a layernorm (LN) layer before each MSA module and each multilayer perceptron (MLP), followed by residual connections. As shown in Fig. 4 (a), PM reduces the resolution of the features and increases the number of channels of the features.

### B. SDM Auxiliary Module

The signed distance map (SDM) [20], [21] models the distance relationship between pixels in the foreground-background region and the boundary, and further distinguishes foreground and background with positive and negative signs.
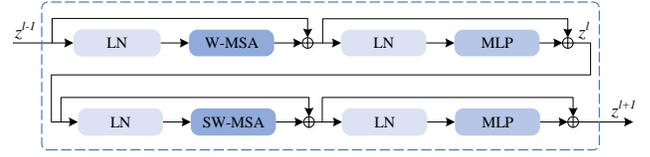


Fig. 3. The architecture of the swin transformer block (STB). W-MSA calculates the pairwise attention of each token in the window. SW-MSA shifts the window of W-MSA by half the window length.

According to the ground truth $\boldsymbol{G}$, SDM transformation $\mathcal{D}(p)$ for each pixel $p \in \boldsymbol{G}$ is given by:

$$
\mathcal{D}(p) = \begin{cases}
-\inf\limits_{\forall b \in \partial \mathcal{S}} d(p, b), & p \in \mathcal{S}_{\mathrm{sal}} \\
0, & p \in \partial \mathcal{S} \\
+\inf\limits_{\forall b \in \partial \mathcal{S}} d(p, b), & p \in \mathcal{S}_{\mathrm{bg}}
\end{cases} \tag{3}
$$

where $\inf$ denotes the infimum, $b$ is the boundary pixel. In Eq. (3), $\partial \mathcal{S}$ is the zero level set and also represents the pixel set of the target boundary. $\mathcal{S}_{\mathrm{sal}}$ and $\mathcal{S}_{\mathrm{bg}}$ indicate the salient object pixel set and background pixel set, respectively. In our work, $d(\cdot)$ indicates the Euclidean distance. As shown in Fig. 1, SDM not only perceives the boundary of an object, but also predicts whether the pixel is located inside or outside the object. For each pixel $p \in \boldsymbol{G}$, the sign of $\mathcal{D}(p)$ indicates whether it is located outside (i.e., $\mathcal{D}(p) > 0$) or inside (i.e., $\mathcal{D}(p) < 0$) the object. $\mathcal{D}(p) = 0$ denotes the boundary of the object. $|\mathcal{D}(p)|$ represents the distance from pixel $p$ to the boundary.

In order to precisely perceive the boundaries of salient objects, we present a SDM auxiliary module (SDMAM) to learn the distance relation between region pixels and boundary pixels. Benefiting from SDM, SDMAM can effectively strengthen the inter-class separability of foreground-background region features. The upper right part of Fig. 2 shows the structure of SDMAM in detail. The shallow high-resolution features contain rich texture information [57]. SDMAM integrates RGB-T shallow features to predict the distance relationship between pixels. Formally,

$$
\boldsymbol{D} = \mathrm{SDMAM}\left(\mathbf{x}^1, \mathbf{x}^2, \mathbf{x}^3\right). \tag{4}
$$

where $\mathbf{x}^i = \mathrm{concat}\left(\mathbf{x}_r^i, \mathbf{x}_t^i\right)$, $i = 1, 2, 3$. $\boldsymbol{D} \in \mathbb{R}^{h \times w \times 1}$ represents the prediction result of SDMAM. The dimensions of $\mathbf{x}^1$, $\mathbf{x}^2$ and $\mathbf{x}^3$ are $\mathbb{R}^{\frac{h}{4} \times \frac{w}{4} \times c}$, $\mathbb{R}^{\frac{h}{8} \times \frac{w}{8} \times 2c}$ and $\mathbb{R}^{\frac{h}{16} \times \frac{w}{16} \times 4c}$, respectively. In this paper, $h = w = 384$, $c = 128$.

Specifically, the multi-scale features $\{\mathbf{x}^1, \mathbf{x}^2, \mathbf{x}^3\}$ are further fused by pointwise convolution (PWConv) [58] with ReLU and upsampling operations,

$$
\mathbf{y}^i = \left[\mathrm{PWConv}\left(\mathbf{x}^i\right)\right]^{\times (2^{i-1})}, \tag{5}
$$

where $i = 1, 2, 3$. $[\cdot]^{\times(n)}$ denotes upsampling the features by $n$ times. In Eq. (5), the different scale high-resolution features $\mathbf{y}^i \in \mathbb{R}^{\frac{h}{4} \times \frac{w}{4} \times 32}$. Finally, the multi-scale multi-spectral features $\mathbf{y}$ are fused by $3 \times 3$ convolution and upsampled to the resolution of the input image.

$$
\begin{aligned}
\mathbf{y} &= \mathrm{concat}\left(\mathbf{y}^1, \mathbf{y}^2, \mathbf{y}^3\right), \\
\boldsymbol{D} &= \tanh\left[\mathrm{Conv}_{3\times3}(\mathbf{y})\right]^{\times(4)},
\end{aligned} \tag{6}
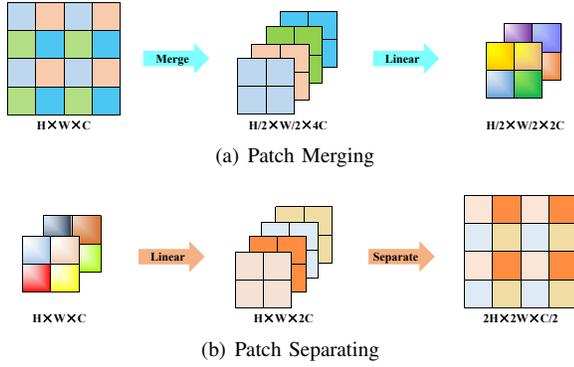$$

(a) Patch Merging



(b) Patch Separating

Fig. 4. **(a)** Patch merging layer (PM) merges the neighboring patches into a new patch, thus reducing the resolution. **(b)** Our proposed patch separating layer (PS) upsamples features by expanding each patch into multiple sub-patches.

where the output of SDMAM is $\boldsymbol{D} \in \mathbb{R}^{h \times w \times 1}$, $\mathrm{Conv}_{3 \times 3}$ indicate $3 \times 3$ convolution with stride 1.

### C. Reverse Swin Transformer Decoder

Our decoder is designed to decode patches as saliency maps. Hence, we propose a novel patch upsampling method with multi-level patch fusion and a patch-based SOD decoder.

*1) Cross Spectrum Fusion Transformer:* Concretely, the RGB-T encoder feature map $\mathbf{x}^4 = \mathrm{concat}\left(\mathbf{x}_r^4, \mathbf{x}_t^4\right)$ is flattened into an input sequence. A set of queries $\mathbf{Q}$, keys $\mathbf{K}$ and values $\mathbf{V}$ is computed by embedding the input sequence into three weight matrices.

In Eq. (7), we compute cross-spectral attention $\mathbf{z}^4$ as in [28].

$$\mathbf{z}^4 = \mathrm{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathrm{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}, \quad (7)$$

where $\mathbf{z}^4 \in \mathbb{R}^{\frac{h}{32} \times \frac{w}{32} \times 8c}$, $\sqrt{d}$ is an adjustment factor that prevents the $\mathrm{softmax}$ function from having too large an input value resulting in too small a partial derivative.

*2) Reverse Swin Transformer Decoder:* In swin transformer, the patch merging (PM) integrates patches of different windows to reduce the spatial resolution of feature maps. Inspired by PM, we design patch separating (PS) to upsample patches by separating each patch for multiple sub-patches. As shown in Fig. 4 (b), Based on PS, we propose a reverse swin transformer layer (RST) for the decoder.

The reverse swin transformer decoder is illustrated in the middle of Fig. 2. RST layer includes STB and PS. For RGB-T features of encoders, RST generates more patches and progressively decodes the patches into high-resolution saliency maps, as in Eq. (8).

$$\mathbf{z}^i = \mathrm{RST}\left(\mathbf{z}^{i+1}, \mathbf{x}^{i+1}\right), \quad (8)$$

where $i = 1, 2, 3$. The dimensions of $\mathbf{z}^1$, $\mathbf{z}^2$ and $\mathbf{z}^3$ are $\mathbb{R}^{\frac{h}{4} \times \frac{w}{4} \times c}$, $\mathbb{R}^{\frac{h}{8} \times \frac{w}{8} \times 2c}$ and $\mathbb{R}^{\frac{h}{16} \times \frac{w}{16} \times 4c}$, respectively. The salient swin transformer decoder output $\mathbf{z} \in \mathbb{R}^{h \times w \times 64}$ is obtained by upsampling $\mathbf{z}^1$ by a factor of 4.

### D. Feature Refinement Approach with Direction Field

The direction field (DF) [59] offers the directional relationship between salient pixels. The direction vector of each pixel points from the boundary to the center. The mathematical definition of the direction field function $\mathcal{F}$ is shown in Eq. (9). The direction of $\mathcal{F}(p)$ is from $b$ pointing to $p$, and $b$ is the nearest pixel to $p$ on the boundary. For the pixel $p \in \boldsymbol{G}$,

$$\mathcal{F}(p) = \begin{cases} \inf_{\forall b \in \partial \mathcal{S}} \vec{bp}, & p \in \mathcal{S}_{\mathrm{sal}} \\ (0, 0). & p \in \mathcal{S}_{\mathrm{bg}} \end{cases} \quad (9)$$

where $\mathcal{S}_{\mathrm{sal}}$ and $\mathcal{S}_{\mathrm{bg}}$ denote the salient object pixel set and background pixel set, respectively.

The refinement of the initial predicted features provides an effective way to improve salient object masks. Based on this idea, we design a feature refinement approach with directional field (FRDF). FRDF uses features inside the object to improve the feature representation near the boundary with the help of directional information between pixels. FRDF progressively enforces the intra-class compactness of salient region features through several iterative updates. As shown in the bottom right of Fig. 2, we first use the decoder feature $\mathbf{z}$ to predict the directional field feature $\boldsymbol{F} \in \mathbb{R}^{h \times w \times 2}$ in Eq. (10).

$$\boldsymbol{F} = \mathrm{PWConv}\left(\mathbf{z}\right). \quad (10)$$

Then, the initial predicted saliency feature map is refined step by step iteratively according to Eq. (11).

$$\mathbf{z}_k(p) = \mathbf{z}_k\left(p_x + \boldsymbol{F}_x(p), p_y + \boldsymbol{F}_y(p)\right), \quad (11)$$

where $\mathbf{z}_k$ denotes the salient feature map after the $k$-th iteration. The number of iterations is setting as $K = 5$, which is further ablated with experiments in Sec. IV-D3. $p_x$ and $p_y$ indicate the $x$ and $y$ coordinates of pixel $p$, respectively. The output of FRDF is the refined feature $\mathbf{z}^* \in \mathbb{R}^{h \times w \times 2c}$.

Finally, the $\mathrm{PWConv}$ layer combines initial feature $\mathbf{z}$ with the rectified feature $\mathbf{z}^*$ to generate the salient mask $\boldsymbol{O}_{sal} \in \mathbb{R}^{h \times w \times 1}$. Both SDM and DF prediction are supervised, which will be discussed in Sec. III-E. Based on the ground truth $\boldsymbol{G}$, we obtain the true supervised signal for the SDMAM and FRDF.

### E. Loss Function

According to the ground truth $\boldsymbol{G}$ of the image, the true SDM and the true direction field of the salient object can be calculated by mathematical models, *i.e.*, Eq. (3) and Eq. (9).

$$\begin{aligned} \boldsymbol{D}_{gt} &= \mathcal{D}(\boldsymbol{G}) \\ \boldsymbol{F}_{gt} &= \mathcal{F}(\boldsymbol{G}) \end{aligned} \quad (12)$$

In Eq. (12), $\boldsymbol{D}_{gt}$ is the true SDM and $\boldsymbol{F}_{gt}$ is the true DF. They guide SDMAM and FRDF to enhance intra-class compactness and inter-class separability, which are weakly supervision for RGB-T SOD.

*1) SDM Loss:* SDM loss is

$$\mathcal{L}_{sdm} = \sum_{p\in\Omega} \|\boldsymbol{D} - \boldsymbol{D}_{gt}\|^2, \qquad (13)$$

where $\Omega$ denotes all pixels, $\boldsymbol{D}$ is the predicted result of SDMAM. $\mathcal{L}_{sdm}$ drives PRLNet to learn the distance relationship between foreground-background regions and boundaries, effectively enhancing the inter-class differences of salient features.

*2) Direction Field Loss:* DF loss is

$$\mathcal{L}_{df} = \sum_{p\in\Omega} \left( \|\boldsymbol{F} - \hat{\boldsymbol{F}}\|_2 + \left\| \cos^{-1}\langle \boldsymbol{F}, \boldsymbol{F}_{gt}\rangle \right\|^2 \right), \qquad (14)$$

where $\boldsymbol{F}$ and $\boldsymbol{F}_{gt}$ indicate the predicted DF and the corresponding ground truth, respectively. $\mathcal{L}_{df}$ guides the model to learn the directional relationship between pixels, which rectifies features of boundary neighborhood by exploiting the features inside salient objects.

*3) Direction-aware Smoothness Loss:* We develop a novel direction-aware smoothness loss ($\mathcal{L}_{DS}$) that enhances the compactness of regions and the boundary clearness. We calculate the first order derivative of the saliency map in the smooth term [26], [49]. $\mathcal{L}_{DS}$ is defined as follows,

$$\mathcal{L}_{DS}(\boldsymbol{O}, \boldsymbol{G}) = \sum_{p\in\Omega}\sum_{\partial_{x,y}} w(p)\,\psi\left( |\partial\boldsymbol{O}|\, e^{-\alpha|\partial\boldsymbol{G}|} \right), \qquad (15)$$

$$w(p) = \begin{cases} \|\mathcal{F}(p)\|^{-1}, & p \in \mathcal{S} \\ 1, & p \in \mathcal{S}_{\text{bg}} \end{cases} \qquad (16)$$

where $\psi(m) = \sqrt{m^2 + 0.001^2}$, $\boldsymbol{O}$ and $\boldsymbol{G}$ represent the predicted salient result and ground truth, respectively. $\partial_{x,y}$ denotes the partial derivatives in $x$ and $y$ directions. Same as [27], we set $\alpha = 10$ to balance the contribution of the edges. In Eq. (16), $w(p)$ indicates the weight on pixel $p$. Therefore, the saliency loss is

$$\mathcal{L}_{sal} = \mathcal{L}_{DS}(\boldsymbol{O}_{sal}, \boldsymbol{G}). \qquad (17)$$

Finally, our position-aware relation learning (PRL) loss is

$$\mathcal{L}_{prl} = \mathcal{L}_{sal} + \lambda_1\mathcal{L}_{sdm} + \lambda_2\mathcal{L}_{df}, \qquad (18)$$

where $\lambda_1$ and $\lambda_2$ are the hyper-parameters balancing the contributions of the two losses, which are set via ablative analysis in Sec. IV-D3. The proposed PRLNet is optimized through Eq. (18) jointly. Our proposed position-aware relation learning loss can effectively guide the network to pay more attention to the pixels around the object boundary, thereby helping the network to predict salient masks with sharp boundaries and homogeneous regions.

## IV. EXPERIMENTS

In this section, we first introduce the three RGB-T datasets, implementation details, and evaluation metrics. We then give the details of our experiments. In particular, we evaluate our method on three widely used datasets to compare with SOTA methods. Moreover, ablation studies are also conducted to further validate the validity of our network.

---

**Algorithm 1:** The pipeline of PRLNet

---

1 **Input**: RGB-T images $\{\boldsymbol{I}_r, \boldsymbol{I}_t\}$, ground truth $\boldsymbol{G}$
2 **Output**: Salient object mask $\boldsymbol{O}_{sal}$, signed distance map $\boldsymbol{D}$, direction field $\boldsymbol{F}$
  1: Init true SDM $\boldsymbol{D}_{gt} \leftarrow \mathcal{D}(\boldsymbol{G})$ using Eq. (3)
  2: Init true direction field $\boldsymbol{F}_{gt} \leftarrow \mathcal{F}(\boldsymbol{G})$ using Eq. (9)
  3: Init the number of FRDF iterations $K = 5$
  4: **while** $epoch < N$ **do**
  5:    Extract RGB-T features $\mathbf{R}$ and $\mathbf{T}$ using Eq. (1)
  6:    Generate signed distance map $\boldsymbol{D}$ using Eq. (4)
  7:    Generate decoder output feature $\mathbf{z}$ using Eq. (8)
  8:    Generate direction field $\boldsymbol{F}$ using Eq. (10)
  9:    **for** $k \leq K$ **do**
 10:      Iterative refinement of decoder feature $\mathbf{z} \leftarrow \mathbf{z} + \boldsymbol{F}$ using Eq. (11)
 11:    **end for**
 12:    Generate salient mask result $\boldsymbol{O}_{sal}$ using initial features $\mathbf{z}$ and refined features $\mathbf{z}^*$
 13:    Calculate the SDM loss $\mathcal{L}_{sdm} \leftarrow \mathcal{L}(\boldsymbol{D}, \boldsymbol{D}_{gt})$ using Eq. (13)
 14:    Calculate the DF loss $\mathcal{L}_{df} \leftarrow \mathcal{L}(\boldsymbol{F}, \boldsymbol{F}_{gt})$ using Eq. (14)
 15:    Calculate the direction-aware smoothness loss $\mathcal{L}_{sal} \leftarrow \mathcal{L}(\boldsymbol{O}_{sal}, \boldsymbol{G})$ using Eq. (17)
 16:    Calculate the overall loss function of PRLNet $\mathcal{L}_{prl} \leftarrow \mathcal{L}_{sal} + \lambda_1\mathcal{L}_{sdm} + \lambda_2\mathcal{L}_{df}$,
 17:    $\boldsymbol{O}_{sal}, \boldsymbol{D}, \boldsymbol{F} \leftarrow \arg\min \mathcal{L}_{prl}$
 18: **end while**
 19: **return** $\boldsymbol{O}_{sal}$

---

### A. Experimental Setup

*1) Datasets:* There are three available benchmark datasets for RGB-T SOD tasks, including VT821 [41], VT1000 [48] and VT5000 [44], which have 821, 1000, and 5000 aligned image pairs, respectively. Compared with VT821, VT1000 dataset has more images and scenes, and the quality of the thermal images is better. VT5000 provides a large-scale dataset for TGB-T SOD. In addition, VT5000 does not require manual RGB-T image pair alignment, which reduces the errors caused by manual alignment. VT5000 contains a variety of complex scenes with diverse objects and covers 13 challenges of RGB-T SOD [44], the details are shown in TABLE I. As reported in TABLE I, VT5000 simulates image saliency detection under real-world conditions mainly in terms of target diversity (BSO, SSO, MSO, CB, CIB, OF, SA and TC), scene complexity (IC, LI and BW) and spectral effectiveness (RGB and T).

*2) Implementation Details:* To extract multispectral features, we initialize our backbone networks through the parameters of the pre-trained Swin-B model [29]. The whole network is then trained on a large-scale dataset with the proposed position-aware relation learning loss in an end-to-end manner.

For a fair comparison, we use the same setting as in [35], [44], [48], [55] where half of VT5000 dataset is applied as the training set. VT821, VT1000 and the other half of VT5000 are treated as the test set. In addition, each image is random flipping, cropping and rotation ($-15° \sim 15°$), and then resized

TABLE I
DETAILS OF 13 CHALLENGES IN VT5000 DATASET.

| Challenge | Describe |
|---|---|
| BSO | Big salient object: the proportion of pixels of salient objects to the image is more than 0.26. |
| SSO | Small salient object: the percentage of the number of salient pixels is less than 0.05. |
| MSO | Multiple salient objects. |
| CB | Center bias: the salient object is out of the center of the image. |
| CIB | Cross image boundary: a part of the salient object is outside the image. |
| OF | Out of focus: out of focus causes the whole image to be blurred. |
| SA | Similar appearance: the salient object is similar to the color and texture of the background. |
| TC | Thermal crossover: the salient object is similar to its surrounding temperature. |
| IC | Image clutter: the scene is cluttered. |
| LI | Low illumination: the scene is cloudy or at night. |
| BW | Bad weather: the scene is rainy or foggy. |
| RGB | Objects are not clear in RGB image. |
| T | Objects are not clear in the thermal image. |

to $384 \times 384$. We train our models by using adaptive optimizer Adam. The initial learning rate of the network is set to $10^{-5}$ and is decayed by 0.1 every 100 epochs. The total epoch number is set to 300. The mini-batch size is set as 6. Our framework is implemented by PyTorch. The experiment is conducted on a computer with 3.0 GHz CPU, 128 GB RAM, and four NVIDIA GeForce RTX 3090 GPUs.

### B. Evaluation Metrics

In order to facilitate the comparison of the performance of different RGB-T methods, we use the evaluation metrics commonly used in SOD model: P-R curves [60], S-measure ($S_\alpha \uparrow$) [61], F-measure ($F_\beta \uparrow$) [62], E-measure ($E_m \uparrow$) [63] and MAE ($\mathcal{M} \downarrow$) [64]. $\uparrow$ and $\downarrow$ indicate that the higher the better and the lower the better, respectively. The P-R curves and $F_\beta$ evaluate the quality of the prediction results in terms of Precision and Recall. $S_\alpha$ and $E_m$ mainly measure the structural similarity between predicted saliency mask and GT. $\mathcal{M}$ counts the error of the prediction result pixel by pixel. We use the above metrics to evaluate the model accurately and comprehensively. The formal definition is as follows.

**P-R curves.** We first demonstrate the performance of our model through standard P-R curves [60]. Different thresholds ($[0, 255]$) are applied to the prediction to generate binarized result that produces pairs of PRECISION-RECALL values. A set of thresholds provides the P-R curve of the model. Formally, the P and R are defined based on the binarized salient object mask and the corresponding ground truth in Eq. (19).

$$P = \frac{TP}{TP + FP}, \quad R = \frac{TP}{TP + FN}, \tag{19}$$

where TP, FP and FN denote true positive, false positive and false negative, respectively.

**S-measure.** The structure measure ($S_\alpha$) can effectively evaluate the spatial structure compactness between prediction and ground truth [61].

$$S_\alpha = \alpha S_o + (1 - \alpha)S_r, \tag{20}$$

where $\alpha$ is set as 0.5 empirically [61]. In Eq. (20), $S_\alpha$ integrates object-aware structural similarity $S_o$ and region-aware structural similarity $S_r$.

**F-measure.** $F_\beta$ takes into account precision and recall [62], and calculates the weighted harmonic mean of P and R:

$$F_\beta = \frac{(1 + \beta^2) \times P \times R}{\beta^2 \times P + R}, \tag{21}$$

where we set $\beta^2 = 0.3$ to weigh precision more than recall.

**E-measure.** The enhanced-alignment measure metric ($E_m$) considers both local pixel values and image-level averages. $E_m$ captures image-level statistics and local pixel matching information [63].

$$E_m = \frac{1}{HW} \sum_{i=1}^{H} \sum_{j=1}^{W} \phi_{ij}. \tag{22}$$

In Eq. (22), $H$ and $W$ are the height and width of the object map, respectively. $\phi$ is the enhanced alignment matrix [63].

**MAE.** The mean absolute error ($\mathcal{M}$) [64] measures the difference between saliency prediction $\boldsymbol{O} \in [0, 1]^{H \times W}$ and ground truth mask $\boldsymbol{G} \in \{0, 1\}^{H \times W}$,

$$\mathcal{M} = \frac{1}{HW} \sum_{i=1}^{H} \sum_{j=1}^{W} |\boldsymbol{O}_{ij} - \boldsymbol{G}_{ij}|. \tag{23}$$

### C. Comparison with State-of-the-Art Methods

To evaluate the validity of the proposed PRLNet, we conduct experiments compared with state-of-the-art methods on three datasets, which are shown in TABLE II and Fig. 6. Three RGB SOD methods include R3Net [35], CPD [38] and PoolNet [39]. Six RGB-T SOD methods include SGDL [48], ADF [44], FMCF [9], MIDD [49], ECFFNet [50] and Swin-Net [55].

*1) Qualitative Comparison:* The results visualized in Fig. 5 display a qualitative comparison of some challenging image pairs, such as SSO (column (a)-(c)), CB (column (d) and (e)), BSO (column (e), (f) and (q)-(t)), BW (column (g) and (r)), TC (column (h)-(j)), LI (column (f) and (k)), MSO (column (k)-(n)), SA (column (i)-(l)), CIB (column (o)-(r) and (u)), IC (column (s) and (v)), OF (column (p) and (r)), RGB images with low quality (column (a), (g), (k), (n) and (v)) and thermal images with low quality (column (a), (e), (i), (m) and (s)). As illustrated in Fig. 5, the results of our PRLNet are qualitatively superior to all SOTA methods. Our method takes full advantage of the discriminative feature representation capabilities of the swin transformer, while taking the position relations between pixels into account, *i.e.*, distance and direction relationships.

As shown in Fig. 5 (e), (h) and (o), the salient objects and background objects in certain spectral images have similar intensities, which can lead to confusion between foreground and background classes. Our proposed SDMAM effectively addresses this problem by explicitly constraining the foreground-background difference with signs and modeling the distance of different pixels from the boundary. SDMAM increases interclass separability. From the results in Fig. 5 (e), (h) and (o), it

TABLE II
QUANTITATIVE COMPARISON WITH SOTA METHOD ON THREE BENCHMARK DATASETS IN TERMS OF S-MEASURE ($S_\alpha \uparrow$), F-MEASURE ($F_\beta \uparrow$), E-MEASURE ($E_m \uparrow$) AND MAE ($\mathcal{M} \downarrow$). $\uparrow$ AND $\downarrow$ REPRESENT THE HIGHER THE BETTER AND THE LOWER THE BETTER, RESPECTIVELY. THE BEST RESULT IN EACH COLUMN IS IN RED, AND THE SECOND IS IN BLUE.

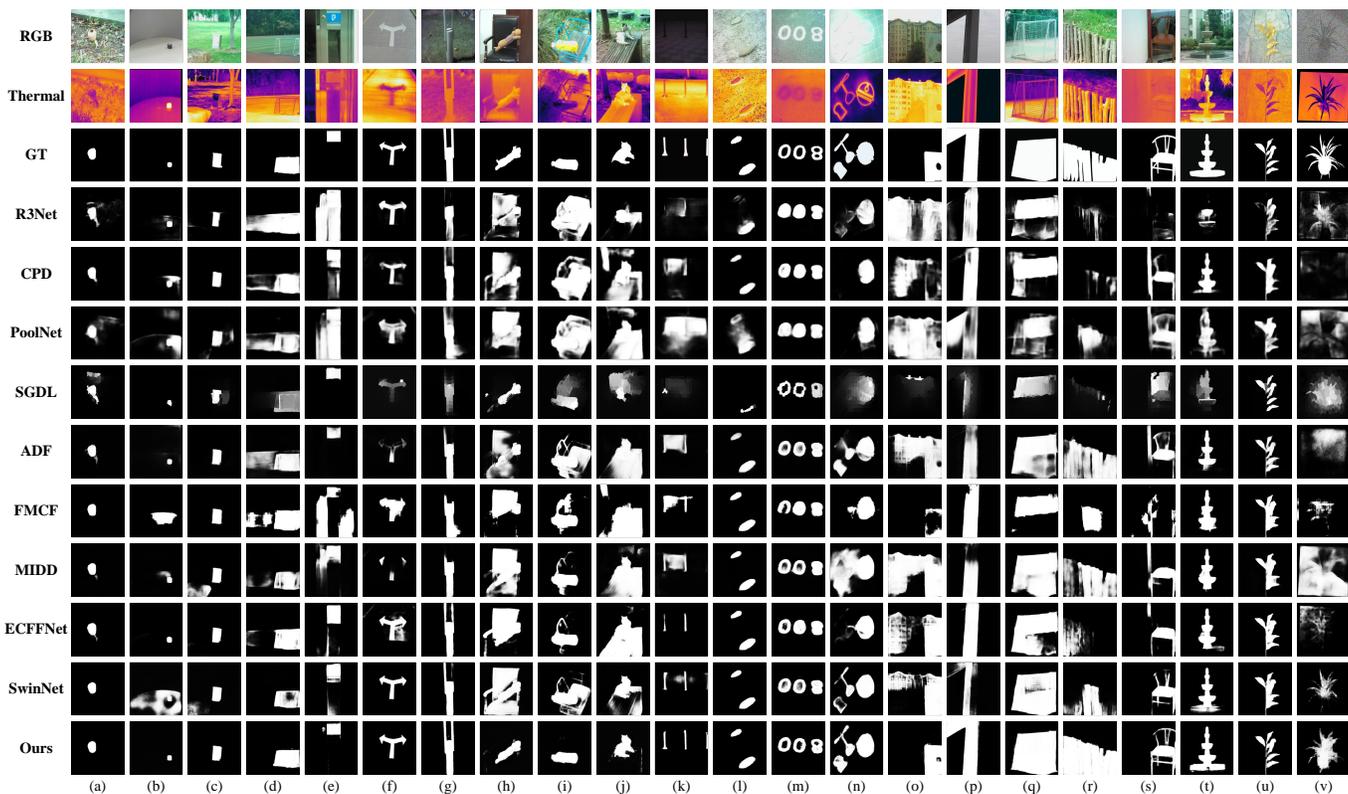| Methods | VT821 | | | | VT1000 | | | | VT5000 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $S_\alpha \uparrow$ | $F_\beta \uparrow$ | $E_m \uparrow$ | $\mathcal{M} \downarrow$ | $S_\alpha \uparrow$ | $F_\beta \uparrow$ | $E_m \uparrow$ | $\mathcal{M} \downarrow$ | $S_\alpha \uparrow$ | $F_\beta \uparrow$ | $E_m \uparrow$ | $\mathcal{M} \downarrow$ |
| R3Net | 0.785 | 0.809 | 0.660 | 0.073 | 0.842 | 0.859 | 0.761 | 0.055 | 0.757 | 0.790 | 0.615 | 0.083 |
| CPD | 0.818 | 0.718 | 0.843 | 0.079 | 0.907 | 0.863 | 0.923 | 0.031 | 0.855 | 0.787 | 0.894 | 0.046 |
| PoolNet | 0.751 | 0.739 | 0.578 | 0.109 | 0.834 | 0.813 | 0.714 | 0.067 | 0.769 | 0.755 | 0.588 | 0.089 |
| SGDL | 0.765 | 0.847 | 0.731 | 0.085 | 0.787 | 0.856 | 0.764 | 0.090 | 0.750 | 0.824 | 0.672 | 0.089 |
| ADF | 0.810 | 0.716 | 0.842 | 0.077 | 0.910 | 0.847 | 0.921 | 0.034 | 0.863 | 0.778 | 0.891 | 0.048 |
| FMCF | 0.760 | 0.796 | 0.640 | 0.080 | 0.873 | 0.899 | 0.823 | 0.037 | 0.814 | 0.864 | 0.734 | 0.055 |
| MIDD | 0.871 | 0.804 | 0.895 | 0.045 | 0.915 | 0.882 | 0.933 | 0.027 | 0.867 | 0.801 | 0.897 | 0.043 |
| ECFFNet | 0.877 | 0.810 | 0.902 | 0.034 | 0.923 | 0.876 | 0.930 | 0.021 | 0.874 | 0.806 | 0.906 | 0.038 |
| SwinNet | 0.904 | 0.847 | 0.926 | 0.030 | 0.938 | 0.896 | 0.947 | 0.018 | 0.912 | 0.865 | 0.942 | 0.026 |
| Our | **0.917** | **0.860** | **0.932** | **0.025** | **0.944** | **0.902** | **0.951** | **0.016** | **0.921** | **0.875** | **0.948** | **0.023** |



Fig. 5. Visual comparisons of different SOTA methods under various challenges, where each column indicates one input image. This figure shows that our proposed method (Ours) consistently generates saliency maps close to the Ground Truth (GT).

can be seen that for background objects similar to the target, such as *cabinets*, *chairs* and *buildings*, our method accurately excludes inter-class interference.

On the other hand, the foreground objects also contain many components with large differences, which leads to inconsistencies in the intra-class features, as shown in Fig. 5 (d), (n) and (q). Our proposed FRDF learns the directional relations of pixels in salient regions, enhancing the intra-class compactness of feature representations. From the results in Fig. 5 (d), (n) and (q), it can be seen that the salient object masks generated by our model are more homogenous compared to other methods. Overall, as shown in Fig. 5 (s), (t), (u) and (v), our PRLNet can generate masks with clear boundaries and

smooth regions for objects with fine structures, such as *stone fountains* and *leafy plants*. The extensive visualization results in Fig. 5 effectively prove that our method can handle a variety of complex scenarios with superior performance. Above all, the saliency masks generated by our PRLNet are consistently the closest to GT.

*2) Quantitative Comparison:* TABLE II and Fig. 6 provides a quantitative comparison of our model with other models on three datasets. First, it can be seen from TABLE II that our PRLNet achieves the highest results on VT821, VT1000 and VT5000. This benefits from the fact that our proposed position-aware relation learning can effectively enhance the intra-class compactness and inter-class separability of feature

TABLE III
PERFORMANCE COMPARISON (F-MEASURE, $F_\beta \uparrow$) WITH NINE METHODS ON 13 CHALLENGES OF THE VT5000 DATASET. BOLD FONT HIGHLIGHTS THE BEST RESULTS IN EACH COLUMN.

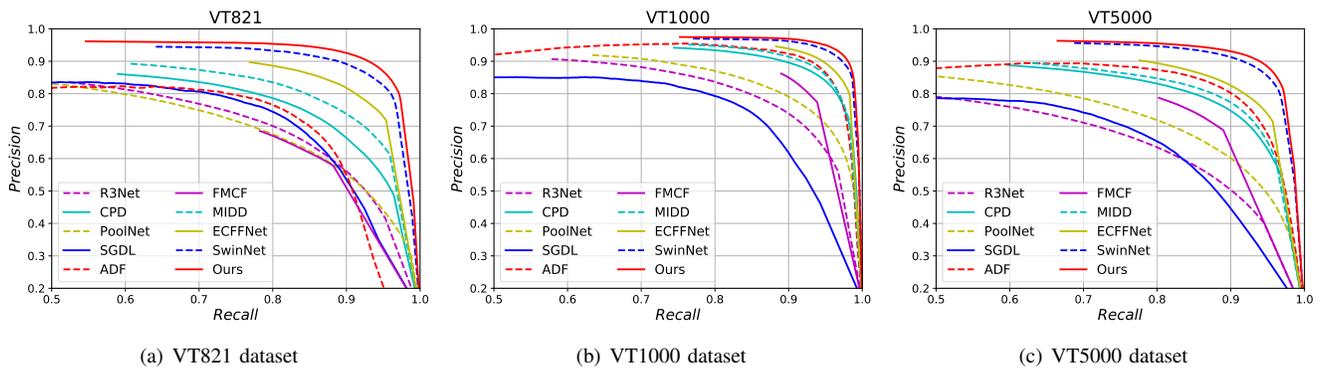| Challenge | BSO | SSO | MSO | CB | CIB | OF | SA | TC | IC | LI | BW | RGB | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| R3Net | 0.734 | 0.538 | 0.609 | 0.623 | 0.654 | 0.701 | 0.614 | 0.608 | 0.624 | 0.709 | 0.562 | 0.673 | 0.683 |
| CPD | 0.835 | 0.694 | 0.765 | 0.777 | 0.799 | 0.801 | 0.756 | 0.789 | 0.764 | 0.823 | 0.694 | 0.804 | 0.805 |
| PoolNet | 0.768 | 0.624 | 0.664 | 0.687 | 0.717 | 0.747 | 0.670 | 0.686 | 0.683 | 0.735 | 0.661 | 0.727 | 0.733 |
| SGDL | 0.722 | 0.715 | 0.660 | 0.656 | 0.654 | 0.707 | 0.598 | 0.621 | 0.631 | 0.697 | 0.583 | 0.705 | 0.710 |
| ADF | 0.858 | 0.737 | 0.806 | 0.821 | 0.837 | 0.806 | 0.791 | 0.792 | 0.803 | 0.845 | 0.771 | 0.840 | 0.842 |
| FMCF | 0.815 | 0.559 | 0.724 | 0.740 | 0.782 | 0.743 | 0.701 | 0.723 | 0.725 | 0.745 | 0.698 | 0.762 | 0.763 |
| MIDD | 0.848 | 0.696 | 0.781 | 0.803 | 0.818 | 0.799 | 0.755 | 0.778 | 0.768 | 0.797 | 0.756 | 0.817 | 0.817 |
| ECFFNet | 0.878 | 0.735 | 0.822 | 0.840 | 0.860 | 0.823 | 0.801 | 0.814 | 0.816 | 0.850 | 0.765 | 0.854 | 0.855 |
| SwinNet | 0.919 | 0.839 | 0.882 | 0.895 | 0.910 | 0.890 | 0.884 | 0.886 | 0.875 | 0.914 | 0.863 | 0.903 | 0.906 |
| Ours | **0.929** | **0.874** | **0.897** | **0.913** | **0.924** | **0.895** | **0.902** | **0.908** | **0.897** | **0.918** | **0.881** | **0.918** | **0.917** |



Fig. 6. P-R curves comparison of different methods on VT821, VT1000 and VT5000 datasets. The P-R curves show that our PRLNet (Ours) consistently outperforms the SOTA models on three datasets.

representations.

Specifically, our PRLNet achieves a marked superiority on VT821. As shown in the results of VT821 in TABLE II, our method improves on average by 0.101, 0.073, 0.152 and 0.043 over other nine methods for $S_\alpha$, $F_\beta$, $E_m$ and $\mathcal{M}$, respectively. Compared with other methods on VT1000, PRLNet has an average improvement of 0.063, 0.036, 0.094, and 0.026 on the four metrics, respectively. As reported in the results of VT5000 from TABLE II, the performance of our PRLNet has improved by an average of 0.092, 0.067, 0.155, and 0.034 on $S_\alpha$, $F_\beta$, $E_m$ and $\mathcal{M}$, respectively. Moreover, for salient masks, structural similarity ($S_\alpha$ and $E_m$) can better characterize the homogeneity of foreground-background regions and the sharpness of boundaries. From the above analysis, it can be seen that our PRLNet improves much higher in $S_\alpha$ and $E_m$ metrics than other two metrics. This indicates out that the salient mask of our method is more sophisticated and close to the ground truth. In addition, compared with the previous state-of-the-art method SwinNet [50] on three datasets, our PRLNet achieves an average gain of 1.02%, 1.12%, 0.57%, 13.11% w.r.t $S_\alpha$, $F_\beta$, $E_m$ and $\mathcal{M}$.

Meanwhile, the P-R curves in Fig. 6 also gives consistent results. As shown in Fig. 6, our curves noticeably lie above the others on VT821, VT1000 and VT5000 datasets. Our proposed method outperforms the state-of-the-art methods. Above all, both the P-R curves and quantization results on the three datasets demonstrate the validity and advantages of our PRLNet for RGB-T SOD.

*3) Quantitative Comparison on Challenge:* To further validate the performance of our PRLNet, we evaluate the performance of each model on all challenges of VT5000 dataset. TABLE I summarizes the challenges covered by the VT5000. Challenge-based quantitative comparison results are reported in TABLE III. The best performance of our PRLNet (Ours) is achieved on all 13 challenges. Compared with SwinNet, our method achieves an average improvement of 1.81% on all challenges. PRLNet achieves an average performance of 0.905 in handling diverse complex targets challenging, such as BSO, SSO, MSO, CB, CIB, OF, SA and TC.

Fig. 5 (b), (c), and (q) show the results on challenges on small objects, multi-object, and large object images, respectively. For example, *soccer goals* and *fences* shown in Fig. 5 (q) and (r) are two common BSO challenges. Compared with the SOTA methods, the BSO mask generated by our PRLNet maintains the global consistency of large objects with better intra-class compactness. TABLE III reports that our method achieves the highest performance of 0.929 on BSO. Some of the TC challenges are shown in Fig. 5 (h), (i) and (j), targets such as *ragdolls* on wooden chairs, *water bottles* in bicycle baskets, and *cats* on park seats have similar thermal radiation to their surroundings. The results from Fig. 5 suggest that the existing methods have difficulty in detecting *ragdolls*, *water bottles* and *cats* from the background with TC. In contrast, our PRLNet obviously suppresses the background objects with thermal crossover, and F-measure attains 0.908 on TC as shown in TABLE III.
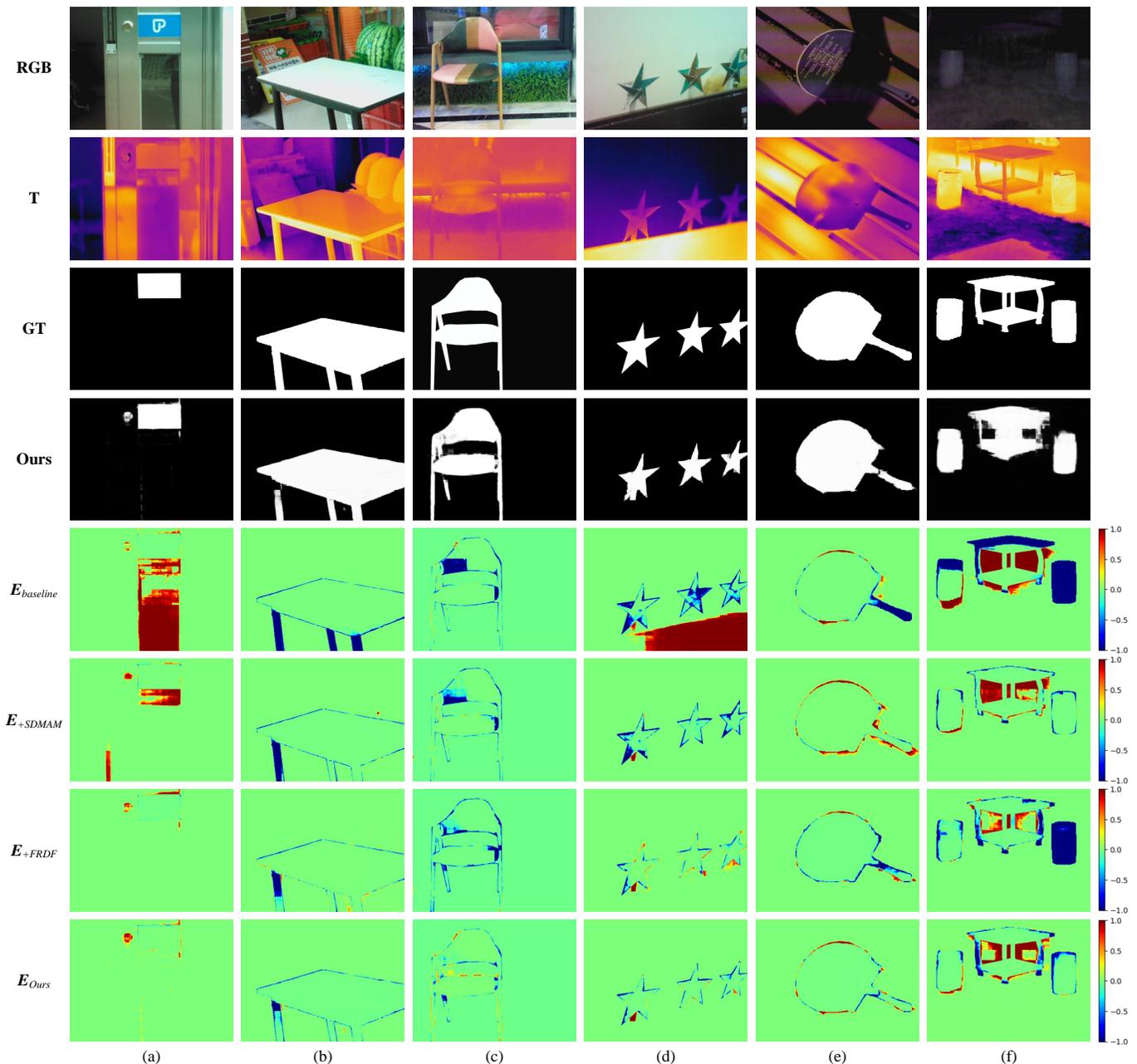
Fig. 7. Visualization results of the error map $E = O_{sal} - G$, where $E(p) > 0$ indicates a false positive pixel (FP), and $E(p) < 0$ indicates a false negative pixel (FN). $E_{baseline}$ represents the error map of the prediction results for the baseline model without SDMAM and FRDF. $E_{+SDMAM}$ and $E_{+FRDF}$ represent the error map after using the SDMAM and FRDF modules, respectively.

The challenges caused by weather or illumination, such as IC, LI and BW, degrade the performance of SOD models. As can be seen from TABLE III, our model still achieves the best performance of about 0.9 for the degraded scenario. In addition, for multispectral RGB-T images, PRLNet effectively learns robust cross-spectral fusion features and reduces the interference caused by spectral inconsistency. The thermal image in Fig. 5 (m) and the RGB image in Fig. 5 (v) have lower quality than the image in the other spectrum. PRLNet overcomes the effect of RGB-T spectral inconsistency and achieves a F-measure above 0.917.

Both the visualization results in Fig. 5 and the quantitative comparisons in TABLE III demonstrate that our method can effectively deal with a variety of salient objects. Above all, the challenge-based quantitative analysis and detailed visualization results consistently demonstrate that our method can effectively address various challenges and outperform state-of-the-art methods.

### D. Ablation Study

Our PRLNet mainly contains two key insights: SDM auxiliary module (SDMAM) and feature refinement approach with direction field (FRDF). Therefore, we conduct ablation

TABLE IV
ABLATION STUDY ON SDMAM AND FRDF ON VT5000 DATASET. BOLD
FONT HIGHLIGHTS THE BEST RESULTS IN EACH COLUMN.

| SDMAM | FRDF | $S_\alpha \uparrow$ | $F_\beta \uparrow$ | $E_m \uparrow$ | $\mathcal{M} \downarrow$ |
|:---:|:---:|:---:|:---:|:---:|:---:|
|  |  | 0.904 | 0.817 | 0.910 | 0.042 |
| ✓ |  | 0.916 | 0.868 | 0.913 | 0.033 |
|  | ✓ | 0.918 | 0.866 | 0.930 | 0.026 |
| ✓ | ✓ | **0.921** | **0.875** | **0.948** | **0.023** |



Fig. 8. Hyper-parameters analysis of $\lambda_1$ and $\lambda_2$ in the position-aware relation learning loss (PRL loss).



Fig. 9. Hyper-parameters analysis of $K$ in the feature refinement approach with directional field.

experiments to verify the validity of components and the involved hyper-parameters.

*1) Effectiveness of SDMAM:* TABLE IV reports the contributions of different components to the model, and Fig. 7 shows the corresponding visualization results. The first row of TABLE IV represents the baseline model, which does not use the SDMAM and FRDF modules. As can be seen from row 2 in TABLE IV, $S_\alpha$, $F_\beta$, $E_m$ and $\mathcal{M}$ attain 0.916, 0.868, 0.913 and 0.033, respectively. SDMAM improves the performance gain by 7.33% on average across the four metrics on the VT5000 dataset compared to the baseline model. The boundary discrimination of the features is enhanced by SDMAM to distinguish salient objects from the background.

To further prove the effectiveness and interpretability of our network, we visualize the error maps (*i.e.*, $\boldsymbol{E}_{+SDMAM}$ and $\boldsymbol{E}_{+FRDF}$) of the saliency maps generated by different components. As shown in Fig. 7 (row 6), SDMAM visibly reduces the error pixels and strengthens the separability of inter-class features. The results of $\boldsymbol{E}_{+SDMAM}$ in the Fig. 7 (a) and (d) illustrate that SDMAM notably suppresses the false alarm (*i.e.*, FP).

*2) Effectiveness of FRDF:* As can be seen from row 3 in TABLE IV, $S_\alpha$, $F_\beta$, $E_m$ and $\mathcal{M}$ attain 0.918, 0.866, 0.930 and 0.026, respectively. FRDF brings an average performance gain of 11.96% over the baseline model. This suggests that the directional information of object pixels is essential and indispensable for learning a fine feature. As shown in Fig. 7 (a), (b), and (e), the error map of the predicted result with FRDF effectively handles the missed detection (*i.e.*, FN) and generates object masks with clear contour and homogeneous regions. The visualization results $\boldsymbol{E}_{+FRDF}$ in Fig. 7 straightforwardly demonstrate the effectiveness of our proposed FRDF. In Section III, we argue that the distance relationship and direction relationship between pixels are crucial for salient object detection, which can be further proved by this experiment. Above all, we can conclude from TABLE IV and Fig. 7 that each component is integral and complementary, which together contribute to the final result.

*3) Hyper-Parameters Analysis:* The parameters $\lambda_1$ and $\lambda_2$ control the relative importance between SDM loss and DF loss in Eq. (18), which is the decisive hyperparameter for the detection results. The greater the value, the more importance lies in the proposed PRL loss. As shown in Fig. 8 (a), MAE score decreases as $\lambda_1$ grows from 0.01 to 1 and increases as $\lambda_1$ grows from 1 to 100, where the valley score reaches 0.023 when $\lambda_1 = 1$. As $\lambda_1$ becomes larger, the SDM loss dominates the PRL loss, leading to model performance degradation. The result in Fig. 8 (b) shows that MAE keeps decreasing when
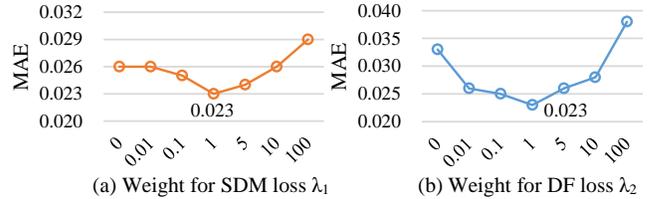
$\lambda_2$ increases to 1. As $\lambda_2$ grows further, the DF loss dominates model training, leading to an increase in the number of error pixels in the saliency mask. Hence, we fix $\lambda_1 = 1$, $\lambda_2 = 1$ in the following experimental settings.

The number of iterations, *i.e.*, $K$, is another important hyperparameter in our method. As $K$ controls the number of iterations of our proposed FRDF, we conduct experiments to verify the choice of $K = 5$ in Eq. (11). We vary $K$ from 0 to 8, as shown in Fig. 9. The MAE constantly decreases as $K$ is growing from 1 to 5. However, there is a slight increase after 5 due to over-refinement with too many iterations. From the above analysis, we choose $K = 5$ as the number of iterations for our feature refinement approach with direction field.

## V. CONCLUSION

In this paper, we have proposed novel a position-aware relation learning network (PRLNet) with pure transformer for RGB-T SOD. PRLNet explored the distance and direction relationships between pixels to strengthen intra-class compactness and inter-class separation. Specifically, we first constructed a dual-stream encoder and decoder framework based on swin transformer, where a patch separation layer was designed to decode the patches. Then, we proposed SDMAM to learn the distance relationship between foreground-background regions and boundaries, which enhanced the boundary perception capability of PRLNet. In addition, we designed FRDF to iteratively rectify the features of the bounding neighborhood using the internal features of the salient objects. FRDF strengthened the intra-class compactness of the salient regions. Extensive experiments and comparisons have shown that the proposed PRLNet consistently outperforms the state-of-the-art methods on three public RGB-T SOD benchmark datasets. Notably, visualization results not only demonstrated that the salient masks generated by our PRLNet have sharp boundaries and homogeneous regions, but also offered a new insight to investigate the relationship between pixels. In future work,

we will pay more attention to the following two directions: camouflage object detection (COD) and multispectral image fusion. Firstly, the proposed relation-aware learning can be applied to COD. In contrast to SOD, COD aims to identify objects embedded in the surrounding environment. Both SOD and COD need to effectively perceive the boundaries of objects and generate masks with clear boundaries and homogeneous regions. Secondly, we will study the efficient complementary fusion between RGB and thermal images. The quality of RGB and thermal images under different illumination conditions is modeled by uncertainty. The poor quality spectral images are used to enhance the good quality spectral images, which enhances the complementarity of multispectral image fusion.

## REFERENCES

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

[2] C. Chen, J. Wei, C. Peng, and H. Qin, "Depth-quality-aware salient object detection," *IEEE Transactions on Image Processing*, vol. 30, pp. 2350–2363, 2021.

[3] H. Chen, Y. Deng, Y. Li, T.-Y. Hung, and G. Lin, "Rgbd salient object detection via disentangled cross-modal fusion," *IEEE Transactions on Image Processing*, vol. 29, pp. 8407–8416, 2020.

[4] M. Feng, H. Lu, and E. Ding, "Attentive feedback network for boundary-aware salient object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1623–1632.

[5] S. Zhou, J. Wang, D. Meng, Y. Liang, Y. Gong, and N. Zheng, "Discriminative feature learning with foreground attention for person re-identification," *IEEE Transactions on Image Processing*, vol. 28, no. 9, pp. 4671–4684, 2019.

[6] C. Dawson, A. Jasour, A. Hofmann, and B. Williams, "Provably safe trajectory optimization in the presence of uncertain convex obstacles," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 6237–6244.

[7] S. Zhou, J. Wang, L. Wang, J. Zhang, F. Wang, D. Huang, and N. Zheng, "Hierarchical and interactive refinement network for edge-preserving salient object detection," *IEEE Transactions on Image Processing*, vol. 30, pp. 1–14, 2020.

[8] 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 Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 6, pp. 3239–3259, 2021.

[9] Q. Zhang, N. Huang, L. Yao, D. Zhang, C. Shan, and J. Han, "Rgb-t salient object detection via fusing multi-level cnn features," *IEEE Transactions on Image Processing*, vol. 29, pp. 3321–3335, 2019.

[10] Q. Liu, X. Li, Z. He, N. Fan, D. Yuan, W. Liu, and Y. Liang, "Multi-task driven feature models for thermal infrared tracking," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 11 604–11 611.

[11] N. Zhang, J. Han, N. Liu, and L. Shao, "Summarize and search: Learning consensus-aware dynamic convolution for co-saliency detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 4167–4176.

[12] W. Zhou, S. Dong, C. Xu, and Q. Yaguan, "Edge-aware guidance fusion network for rgb–thermal scene parsing," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.

[13] S. Song, H. Yu, Z. Miao, J. Fang, K. Zheng, C. Ma, and S. Wang, "Multi-spectral salient object detection by adversarial domain adaptation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 12 023–12 030.

[14] E. Bondi, R. Jain, P. Aggrawal, S. Anand, R. Hannaford, A. Kapoor, J. Piavis, S. Shah, L. Joppa, B. Dilkina *et al.*, "Birdsai: A dataset for detection and tracking in aerial thermal infrared videos," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020, pp. 1747–1756.

[15] Y. Hao, N. Wang, J. Li, and X. Gao, "Hsme: Hypersphere manifold embedding for visible thermal person re-identification," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 8385–8392.

[16] H. Zhou, C. Tian, Z. Zhang, Q. Huo, Y. Xie, and Z. Li, "Multi-spectral fusion transformer network for rgb-thermal urban scene semantic segmentation," *IEEE Geoscience and Remote Sensing Letters*, vol. 19, pp. 1–5, 2022.

[17] X. Qin, Z. Zhang, C. Huang, C. Gao, M. Dehghan, and M. Jagersand, "Basnet: Boundary-aware salient object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7479–7489.

[18] Z. Zhang, C. Tian, X. Gao, J. Li, Z. Jiao, C. Wang, and Z. Zhong, "Collaborative boundary-aware context encoding networks for error map prediction," *Pattern Recognition*, vol. 125, p. 108515, 2022.

[19] X. Zhang, B. Ma, H. Chang, S. Shan, and X. Chen, "Location sensitive network for human instance segmentation," *IEEE Transactions on Image Processing*, vol. 30, pp. 7649–7662, 2021.

[20] W. Zhang, X. Wang, W. You, J. Chen, P. Dai, and P. Zhang, "Resls: Region and edge synergetic level set framework for image segmentation," *IEEE Transactions on Image Processing*, vol. 29, pp. 57–71, 2019.

[21] Q. Cai, Y. Qian, S. Zhou, J. Li, Y.-H. Yang, F. Wu, and D. Zhang, "Avlsm: Adaptive variational level set model for image segmentation in the presence of severe intensity inhomogeneity and high noise," *IEEE Transactions on Image Processing*, vol. 31, pp. 43–57, 2021.

[22] A. A. Farag, H. E. Abd El Munim, J. H. Graham, and A. A. Farag, "A novel approach for lung nodules segmentation in chest ct using level sets," *IEEE Transactions on Image Processing*, vol. 22, no. 12, pp. 5202–5213, 2013.

[23] D. Chai, S. Newsam, and J. Huang, "Aerial image semantic segmentation using dcnn predicted distance maps," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 161, pp. 309–322, 2020.

[24] L. Lin, Z. Wang, J. Wu, Y. Huang, J. Lyu, P. Cheng, J. Wu, and X. Tang, "Bsda-net: A boundary shape and distance aware joint learning framework for segmenting and classifying octa images," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021, pp. 65–75.

[25] F. Cheng, C. Chen, Y. Wang, H. Shi, Y. Cao, D. Tu, C. Zhang, and Y. Xu, "Learning directional feature maps for cardiac mri segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2020, pp. 108–117.

[26] C. Godard, O. Mac Aodha, and G. J. Brostow, "Unsupervised monocular depth estimation with left-right consistency," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 270–279.

[27] Y. Wang, Y. Yang, Z. Yang, L. Zhao, P. Wang, and W. Xu, "Occlusion aware unsupervised learning of optical flow," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4884–4893.

[28] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in Neural Information Processing Systems*, vol. 30, 2017.

[29] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10 012–10 022.

[30] X. He, Y. Zhou, J. Zhao, D. Zhang, R. Yao, and Y. Xue, "Swin transformer embedding unet for remote sensing image semantic segmentation," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–15, 2022.

[31] X. Yu, D. Shi, X. Wei, Y. Ren, T. Ye, and W. Tan, "Soit: Segmenting objects with instance-aware transformers," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 3, 2022, pp. 3188–3196.

[32] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.

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

[34] L. Zhang, X. Fang, H. Bo, T. Wang, and H. Lu, "Deep multi-level networks with multi-task learning for saliency detection," *Neurocomputing*, vol. 312, pp. 229–238, 2018.

[35] Z. Deng, X. Hu, L. Zhu, X. Xu, J. Qin, G. Han, and P.-A. Heng, "R3net: Recurrent residual refinement network for saliency detection," in *Proceedings of the 27th International Joint Conference on Artificial Intelligence*. AAAI Press Menlo Park, CA, USA, 2018, pp. 684–690.

[36] Y.-H. Wu, Y. Liu, L. Zhang, M.-M. Cheng, and B. Ren, "Edn: Salient object detection via extremely-downsampled network," *IEEE Transactions on Image Processing*, vol. 31, pp. 3125–3136, 2022.

[37] X. Li, L. Zhao, L. Wei, M.-H. Yang, F. Wu, Y. Zhuang, H. Ling, and J. Wang, "Deepsaliency: Multi-task deep neural network model for salient object detection," *IEEE Transactions on Image Processing*, vol. 25, no. 8, pp. 3919–3930, 2016.

[38] Z. Wu, L. Su, and Q. Huang, "Cascaded partial decoder for fast and accurate salient object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3907–3916.

[39] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3917–3926.

[40] Z. Wei, X. Yang, N. Wang, and X. Gao, "Syncretic modality collaborative learning for visible infrared person re-identification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 225–234.

[41] G. Wang, C. Li, Y. Ma, A. Zheng, J. Tang, and B. Luo, "Rgb-t saliency detection benchmark: Dataset, baselines, analysis and a novel approach," in *Chinese Conference on Image and Graphics Technologies*. Springer, 2018, pp. 359–369.

[42] Z. Tu, T. Xia, C. Li, Y. Lu, and J. Tang, "M3s-nir: Multi-modal multi-scale noise-insensitive ranking for rgb-t saliency detection," in *2019 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)*. IEEE, 2019, pp. 141–146.

[43] W. Gao, G. Liao, S. Ma, G. Li, Y. Liang, and W. Lin, "Unified information fusion network for multi-modal rgb-d and rgb-t salient object detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 4, pp. 2091–2106, 2021.

[44] Z. Tu, Y. Ma, Z. Li, C. Li, J. Xu, and Y. Liu, "Rgbt salient object detection: A large-scale dataset and benchmark," *IEEE Transactions on Multimedia*, 2022.

[45] J. Wang, K. Song, Y. Bao, L. Huang, and Y. Yan, "Cgfnet: Cross-guided fusion network for rgb-t salient object detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 5, pp. 2949–2961, 2021.

[46] F. Huo, X. Zhu, L. Zhang, Q. Liu, and Y. Shu, "Efficient context-guided stacked refinement network for rgb-t salient object detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 5, pp. 3111–3124, 2021.

[47] Y. Liang, G. Qin, M. Sun, J. Qin, J. Yan, and Z. Zhang, "Multi-modal interactive attention and dual progressive decoding network for rgb-d/t salient object detection," *Neurocomputing*, vol. 490, pp. 132–145, 2022.

[48] Z. Tu, T. Xia, C. Li, X. Wang, Y. Ma, and J. Tang, "Rgb-t image saliency detection via collaborative graph learning," *IEEE Transactions on Multimedia*, vol. 22, no. 1, pp. 160–173, 2019.

[49] Z. Tu, Z. Li, C. Li, Y. Lang, and J. Tang, "Multi-interactive dual-decoder for rgb-thermal salient object detection," *IEEE Transactions on Image Processing*, vol. 30, pp. 5678–5691, 2021.

[56] H. Zhu, X. Sun, Y. Li, K. Ma, S. K. Zhou, and Y. Zheng, "Dftr: Depth-supervised hierarchical feature fusion transformer for salient object detection," *arXiv preprint arXiv:2203.06429*, 2022.

[50] W. Zhou, Q. Guo, J. Lei, L. Yu, and J.-N. Hwang, "Ecffnet: Effective and consistent feature fusion network for rgb-t salient object detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 3, pp. 1224–1235, 2021.

[51] N. Zhang, J. Han, and N. Liu, "Learning implicit class knowledge for rgb-d co-salient object detection with transformers," *IEEE Transactions on Image Processing*, 2022.

[52] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations*, 2021. [Online]. Available: https://openreview.net/forum?id=YicbFdNTTy

[53] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European Conference on Computer Vision*. Springer, 2020, pp. 213–229.

[54] C. Zeng and S. Kwong, "Dual swin-transformer based mutual interactive network for rgb-d salient object detection," *arXiv preprint arXiv:2206.03105*, 2022.

[55] Z. Liu, Y. Tan, Q. He, and Y. Xiao, "Swinnet: Swin transformer drives edge-aware rgb-d and rgb-t salient object detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 7, pp. 4486–4497, 2022.

[57] Y. Liu, J. Han, Q. Zhang, and C. Shan, "Deep salient object detection with contextual information guidance," *IEEE Transactions on Image Processing*, vol. 29, pp. 360–374, 2019.

[58] M. Lin, Q. Chen, and S. Yan, "Network in network," *arXiv preprint arXiv:1312.4400*, 2013.

[59] Z. Zhang, C. Tian, H. X. Bai, Z. Jiao, and X. Tian, "Discriminative error prediction network for semi-supervised colon gland segmentation," *Medical Image Analysis*, vol. 79, p. 102458, 2022.

[60] Q. Zhang, R. Cong, C. Li, M.-M. Cheng, Y. Fang, X. Cao, Y. Zhao, and S. Kwong, "Dense attention fluid network for salient object detection in optical remote sensing images," *IEEE Transactions on Image Processing*, vol. 30, pp. 1305–1317, 2020.

[61] D.-P. Fan, M.-M. Cheng, Y. Liu, T. Li, and A. Borji, "Structure-measure: A new way to evaluate foreground maps," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 4548–4557.

[62] R. Achanta, S. Hemami, F. Estrada, and S. Susstrunk, "Frequency-tuned salient region detection," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 1597–1604.

[63] D.-P. Fan, C. Gong, Y. Cao, B. Ren, M.-M. Cheng, and A. Borji, "Enhanced-alignment measure for binary foreground map evaluation," in *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence*, 2018, pp. 698–704.

[64] F. Perazzi, P. Krähenbühl, Y. Pritch, and A. Hornung, "Saliency filters: Contrast based filtering for salient region detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 733–740.