

Partitioned Saliency Ranking with Dense Pyramid Transformers

Chengxiao Sun*
School of Software Engineering,
Huazhong University of Science and
Technology
Wuhan, China
sunchengxiao@hust.edu.cn

Yan Xu*
School of Software Engineering,
Huazhong University of Science and
Technology
Wuhan, China
yan_xu@hust.edu.cn

Jialun Pei
School of Computer Science and
Engineering, The Chinese University
of Hong Kong
Hong Kong, China
jialunpei@cuhk.edu.hk

Haopeng Fang
School of Software Engineering,
Huazhong University of Science and
Technology
Wuhan, China
haopengfang@hust.edu.cn

He Tang†
School of Software Engineering,
Huazhong University of Science and
Technology
Wuhan, China
hetang@hust.edu.cn

ABSTRACT

In recent years, saliency ranking has emerged as a challenging task focusing on assessing the degree of saliency at instance-level. Being subjective, even humans struggle to identify the precise order of all salient instances. Previous approaches undertake the saliency ranking by directly sorting the rank scores of salient instances, which have not explicitly resolved the inherent ambiguities. To overcome this limitation, we propose the ranking by partition paradigm, which segments unordered salient instances into partitions and then ranks them based on the correlations among these partitions. The ranking by partition paradigm alleviates ranking ambiguities in a general sense, as it consistently improves the performance of other saliency ranking models. Additionally, we introduce the Dense Pyramid Transformer (DPT) to enable global cross-scale interactions, which significantly enhances feature interactions with reduced computational burden. Extensive experiments demonstrate that our approach outperforms all existing methods. The code for our method is available at <https://github.com/ssecv/PSR>.

CCS CONCEPTS

• **Computing methodologies** → **Interest point and salient region detections.**

KEYWORDS

saliency ranking, instance segmentation, partition, self-attention

1 INTRODUCTION

Saliency detection is a crucial research area in computer vision. Prior studies have focused on pixel-level salient object detection [6, 12, 13, 24, 37, 43, 46–49] and salient instance segmentation [18, 23, 30, 39, 45]. Though these works have achieved promising results, they do not take into account the relative saliency ranks among objects, which is more aligned with the human visual system [35]. In order to tackle this problem, Islam *et al.* [17] introduce a new task called saliency ranking (SR), which not only segments the salient instances from the image but also predicts the relative saliency

*Both authors contributed equally to this research.

†Corresponding author: He Tang (E-mail: hetang@hust.edu.cn)

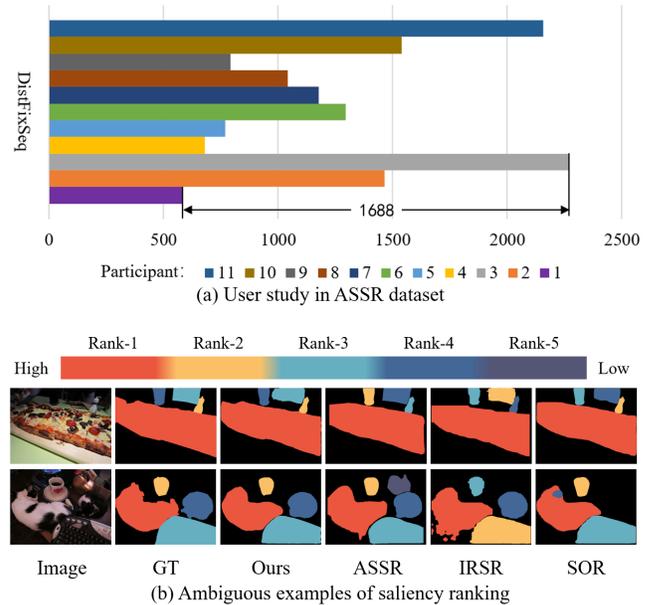


Figure 1: Challenges in saliency ranking. (a) The user study conducted on the ASSR dataset revealed a discrepancy among participants regarding the methodology for determining the degree of saliency. (b) Previous works have not been successful in predicting saliency ranks and masks in ambiguous and complex situations.

ranks of these instances. Subsequent works [10, 22, 36, 38] advance the pixel-level prediction [17] to instance-level prediction. This is promising to be applied in other vision tasks, like person re-identification [3, 34], human gaze communication [8] and video conversion [50].

For generating the most consistent ground truth (GT) of saliency ranking, Siris *et al.* [36] compare nine methods via a user study with 11 participants across 2,500 images. Among the methods, DistFixSeq gains the highest number of image picks from participants and being adopted to generate the ground truth of saliency ranking. However,

as shown in Fig. 1(a), there is a discrepancy in the number of image picks between participant-1 and participant-3, which reaches 1688 images with a proportion of 67.5%. This shows that GT of saliency ranking can still be ambiguous even when using the most consistent generating method. This ambiguity may be due to the fact that saliency ranking is difficult even for humans to determine the exact ranks among salient objects without dispute [22].

The intrinsic ambiguity of saliency ranking makes it difficult for modern methods to accurately predict the ranks of salient instances. As shown in the top row of Fig. 1(b), while methods such as ASSR [36], IRSR [22], and SOR [10] can predict reasonable masks of salient instances, their predicted ranks differ among each other, especially for the ones with lower ranks.

Considering the paradigm of previous works, [10, 22, 36, 38] directly predict ranks of salient instances by sorting the rank scores, *i.e.*, ranking by sorting refer to Fig. 2(a). But this paradigm is not designed to explicitly address the inherent ambiguities in saliency ranking, resulting in incorrect assessment for salient instances with inferior ranks, as shown in top row of Fig. 1(b).

On the other hand, ASSR [36], IRSR [22] and SOR [10] are based on sparse interaction of proposal features [15]. They obtain object proposals at the first stage, then interact these proposals to build relative relationship at the second stage. However, these models select the salient object proposals and discard the background and the objects with lower degrees of saliency. The missing features that not participate in interactions are also useful for saliency ranking. As shown in the bottom row of Fig. 1(b), results of the sparse interaction-based methods ASSR [36], IRSR [22] and SOR [10] lead to false positives and false negatives.

In this paper, we tackle the above limitations by alleviating ambiguity and enhancing feature interaction when ranking the salient instances. It is assumed that humans are easier to identify the *entire* N most attractive objects than to determine their exact relative order *one-by-one*. Based on this assumption, we propose a new paradigm to alleviate the ambiguity problem in saliency ranking, namely, ranking by partition. As shown in Fig. 2(b), we produce N saliency partitions, where N is equal to the maximum rank of a dataset.

Each partition consists of a set of unordered salient instances that are prioritized as equal to or higher than the corresponding rank. Specifically, the partition- n produces at most n unordered salient instances.

To implement this idea, we design partition heads to predict probabilities of N partitions for each instance, and propose Partition to Rank (P2R) to infer the exact rank of each salient instance via the correlations among these saliency partitions. This ranking by partition paradigm is helpful for alleviating saliency ambiguity.

Furthermore, we sufficiently enhance the feature interaction in a global and cross-scale manner. In order to avoid the computational complexity of quadratic to the product of the number of scales and spatial resolution, we design the Dense Pyramid Transformer (DPT), dividing the interaction process into three routes: row attention, column attention and cross-scale attention. The proposed DPT outperforms Twin Transformer [14] and Deformable Transformer [51] with aligned speed.

To sum up, the contributions of this work are as follows:

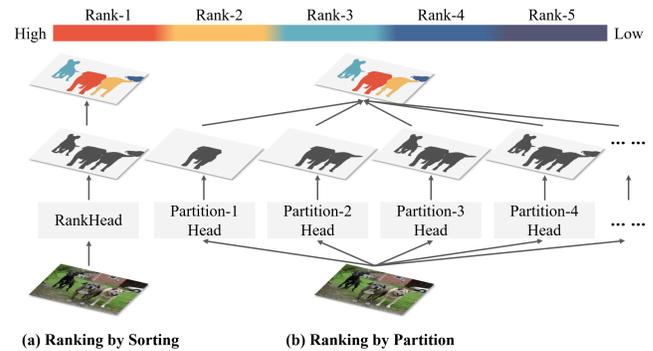


Figure 2: (a) Ranking by sorting and (b) ranking by partition. Ranking by sorting is directly predicting rank of salient instance by sorting the rank scores. Ranking by partition involves identifying saliency partitions of unordered salient instances and ranking the salient instances by the correlations among these partitions.

- To alleviate the ambiguity in saliency ranking, we propose a new paradigm, ranking by partition other than ranking by sorting.
- We propose dense pyramid transformer (DPT) that achieves global cross-scale interaction with reduced computational burden.
- We conduct extensive experiments to analyze our approach and verify its superior performance over the state-of-the-art methods.

2 RELATED WORK

2.1 Saliency Ranking

Saliency Ranking task is motivated by real-world scenarios where human prioritize certain objects over others. Islam *et al.* [17] first introduced saliency ranking in the computer vision community. They treated it as a pixel-wise regression problem. Siris *et al.* [36] introduced a new dataset ASSR according to human attention shift, and proposed a method that used both bottom-up and top-down attention mechanisms to predict the ranks and masks of salient instances. This was the first instance-level saliency ranking method, the subsequent works followed this objective. Liu *et al.* [22] built another dataset IRSR that depended on the duration of gaze instead of the sequential order. In addition, they designed a graph convolution based network to predict relative saliency ranking and proposed a ranking loss that incorporated rank order of GT. Fang *et al.* [10] proposed the first transformer-based position-preserved attention module and used an end-to-end multi-task model that simultaneously performed instance segmentation and saliency ranking. Tian *et al.* [38] proposed a bi-directional proposal interaction method with a selective object saliency module. In the domain of video, [42] proposed a technique to rank the saliency of objects based on the relative fixations of human observers.

By modeling the interaction among proposals and context, the above methods predict mask and relative ranks of salient instances. However, they follow the ranking by sorting paradigm that not

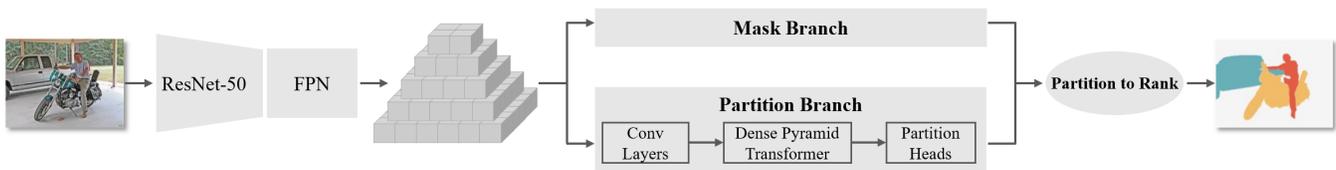


Figure 3: The overall architecture of the proposed partitioned saliency ranking.

adequately address the label ambiguity challenge. In contrast, in this paper, we propose a new ranking by partition paradigm with dense pyramid transformers to alleviate saliency ambiguity and promote cross-scale global feature interaction.

2.2 Regular Instance Segmentation

Regular Instance Segmentation (RIS) is a fundamental vision task that predicts all masks and semantic classes of objects in a scene. Mask R-CNN [15] was a typical instance segmentation method which improved on [33] by using FCN [26] to predict class and mask for each proposal generated by ROIAlign layer. Ma *et al.* [27] introduced an implicit feature refinement module for high-quality instance segmentation. SOLO series [40, 41] segmented objects by location without box detection. MaskFormer series [4, 5] predicted a set of binary masks with a single global class label prediction, implementing semantic and instance segmentation in a unified framework. Recently, bipartite matching [2] based methods [7, 11] achieved promising results without NMS in postprocessing. Besides, video instance segmentation also achieved significant improvements with tracklet query, tracklet proposal [32] and generative model [32]. RIS and SR are similar on the segmentation task, meanwhile, the difference between them is that the predicted semantic class of RIS is independent among each object rather than the relative rank of saliency in SR. Nevertheless, modern RIS models like Mask R-CNN [15] and QueryInst [11] are adopted by [22, 36, 38] as base models for saliency ranking. As the simplification, we select SOLOv2 [41] as our base model.

2.3 Salient Instance Segmentation

Salient Instance Segmentation (SIS) aims to segment individual salient objects rather than pixel-level SOD. Li *et al.* [18] first introduced the concept of salient instance segmentation and proposed a corresponding dataset containing 1000 images with instance and pixel-level as well as contours annotations. They also proposed a multi-scale network MSRNet to segment salient instances with the help of salient object contours. Subsequent works [9, 23, 45] segmented salient instance in a single-stage manner without using contours annotations. Weakly supervised SIS methods [31, 39] predicted mask of each salient object without using the instance-level annotations. Recently, Pei *et al.* proposed a method [30] which leveraged the long-range dependencies of transformers and built a new dataset SIS10K.

Considering the maximum number of salient instances, SIS can be regarded as a special case of SR. Since SIS also segment masks of each salient instance as SR, but do not distinguish the saliency

degree and ranking order. If the binary (salient/non-salient) classification head of a SIS model is extended to the maximum number of saliency ranking, it can also be served as a satisfactory saliency ranking model.

3 PROBLEM FORMULATION

3.1 Ranking by sorting

Previous saliency ranking methods [10, 22, 36, 38] segment the image $I \in \mathbb{R}^{3 \times H \times W}$ into a set of masks, then sort the rank scores of salient instances for ranking, *i.e.*, ranking by sorting. This process is depicted in Fig. 2(a). We assume that N is equal to the maximum rank of the dataset and $f(\cdot)$ is the forward process of the network. Let $M = \{m_i | m_i \in \{0, 1\}^{H \times W}\}_{i=1}^N$, denote the masks of instances, and $R = \{r_i | r_i \in \{1, 2, \dots, N\}\}_{i=1}^N$, denote the ranks of instances. The prediction including mask and saliency rank of each instance is formulated as

$$\langle m_i, r_i \rangle = f(I). \quad (1)$$

However, as shown in Fig. 1(b), ranking by sorting based methods [10, 22, 36] fail to predict masks and ranks of salient instances with lower ranks.

3.2 Ranking by partition

Generally speaking, it is less complicated for individuals to select N most salient instances concurrently than to identify the precise sequential ordering of salient instances. Inspired by this perspective, we propose a novel paradigm, *i.e.*, ranking by partition which reduces the ambiguity of ranking each salient instance.

As delineated in Fig. 2(b), each partition consists of a set of unordered salient instances that are prioritized as equal to or higher than the corresponding rank. Let $P = \{p_n\}_{n=1}^N$ denotes the probabilities of each instance belonging to the N partitions. The outputs of the network comprise predictions in the form of mask and partition probabilities for each instance, which can be formulated as

$$\langle M, p_1, p_2, \dots, p_N \rangle = f(I). \quad (2)$$

We propose the Partition to Rank (P2R) to infer the rank of each instance by parsing the correlations among these partitions. The parsing result is formulated as

$$\langle m_i, r_i = n \rangle = P2R(M, p_1, p_2, \dots, p_n). \quad (3)$$

4 METHODOLOGY

4.1 Overall architecture

The overall architecture is depicted in Fig. 3. The network first utilizes a convolutional backbone, ResNet-50 [16], along with a

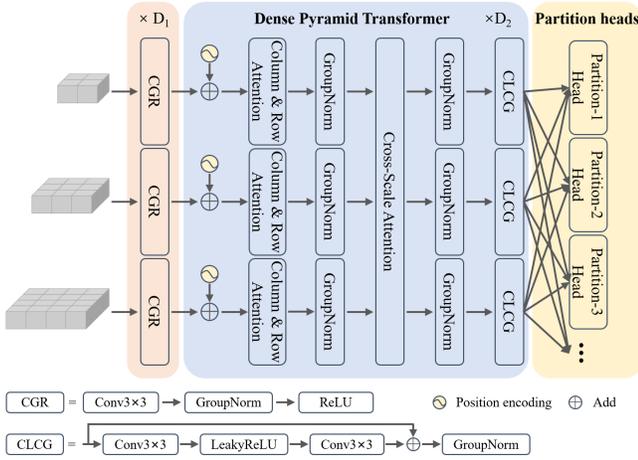


Figure 4: Partition branch. It includes convolutional layers, dense pyramid transformer layers and partition heads. For graphical perspicuity, the illustration omits two scales of feature maps.

feature pyramid network (FPN) [19]. Given an input image $I \in \mathbb{R}^{3 \times H \times W}$, we extract multi-scale features $C = \{c_i\}_{i=2}^6$ from the FPN.

The multi-scale features C are partitioned into grids, with the size of the grids varying across scales. The grids can be formulated as $G = \{g_i | g_i \in \mathbb{R}^{E \times s_i \times s_i}\}_{i=2}^6$, where E denotes the number of channels and $\{s_i\}_{i=2}^6$ signifies the side length of grids. The grids G serve as the input to the mask branch and partition branch.

Following SOLOv2 [41], we adopt dynamic instance segmentation as our mask branch. Leveraging the grids G , dynamic convolution kernels are computed. We then integrate the features from FPN to obtain a new global feature map. The instance masks are obtained by performing dynamic convolutions on the global feature map.

To predict saliency partitions, we design a partition branch consisting of three components as delineated in Fig. 4. The first component comprises convolutional layers which perform preliminary processing on G . The second part, DPT, enables comprehensive interaction among all grid cells across scales. Finally, there are N partition heads to predict partition probabilities for each instance utilizing the outputs of the DPT.

4.2 Partition branch

4.2.1 Convolutional Layers. Previous works [10, 22, 38] provide the evidence that feature interaction is crucial for saliency ranking. We attempt to establish global cross-scale feature interaction with transformers.

According to [28], self-attention and convolutions are complementary. Hence, we adopt convolutional layers before global cross-scale feature interaction for harmonization. The input to the convolutional layers is the grids G . The process of harmonizing is described as CGR in Fig. 4. The output of CGR is formulated as

$$\hat{G} = \text{ReLU}(\text{GN}(\text{Conv}(G))), \quad (4)$$

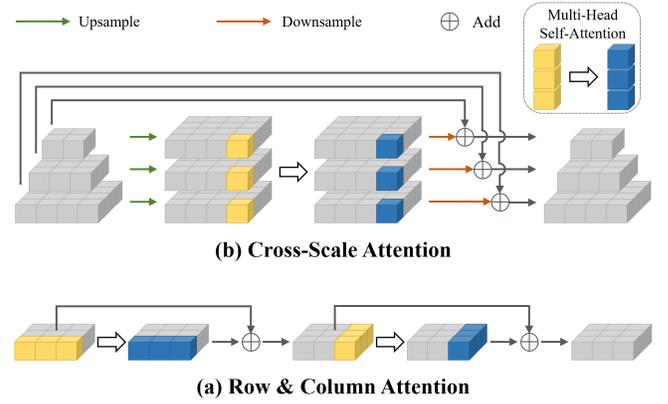


Figure 5: Dense pyramid transformer consists of (a) Row & Column Attention and (b) Cross-Scale Attention. For graphical perspicuity, the illustration omits some scales of feature maps.

where \hat{G} denote the harmonized feature.

4.2.2 Dense Pyramid Transformer. To achieve global cross-scale feature interaction, one feasible approach is to flatten all features across scales and concatenate them before enabling interaction through multi-head self-attention (MHSA), refer to all-scale transformer in this paper. However, if transformer layers are designed this way, the computational complexity per layer would be $O((S \times H \times W)^2)$, where S is the number of scales.

We find there are two elegant transformers which achieve better interaction with less computation. One is twin transformer [14] which decomposes 2D attention into column and row attentions. Another is deformable transformer [51] which can interact with cross-scale features efficiently. Inspired by them, we propose a novel transformer architecture, dense partition transformer (DPT), which divides the interaction process into three components: row attention, column attention and cross-scale attention. The proposed DPT first facilitates comprehensive interaction among features within the same scale through row and column attentions. It then establishes cross-scale attention to promote interaction between features at different scales. This reduces the computational complexity to $O(S \times H \times W^2 + S \times H^2 \times W + S^2 \times H \times W)$, while preserving the comprehensiveness of interaction.

The detailed structure of DPT is shown in Fig. 4. We first add position encoding to \hat{G} and acquire positioned feature F . Subsequently, we utilize the row and column attention, as shown in Fig. 5(a), to enable interaction within each scale of F . In each scale, MHSA is first applied within each row of features, which are then combined with the original features through residual connections. This procedure is then repeated within each column of features. In summary, the process of row and column attentions is described as

$$F_{:,y,z}^R = \text{MHSA}(F_{:,y,z}) + F_{:,y,z}, \quad (5)$$

$$F_{x,,:z}^{RC} = \text{MHSA}(F_{x,,:z}^R) + F_{x,,:z}^R, \quad (6)$$

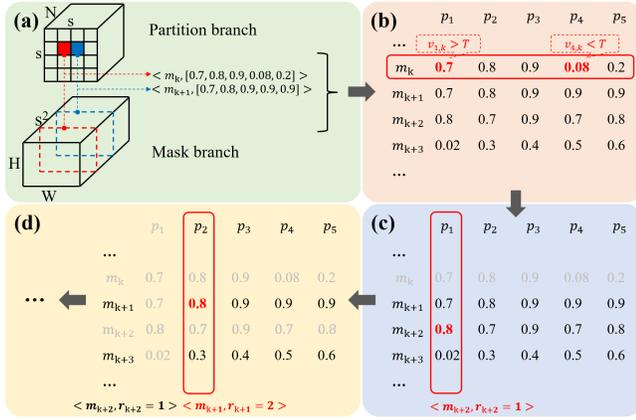


Figure 6: The illustration of our partition to rank (P2R), the ranks are ranging from 1 to 5. For graphical perspicuity, the illustration omits some scales of feature maps. (a) The associated output of partition branch and mask branch; (b) ambiguity alleviation of saliency ranking; (c) rank-1 instance selection; (d) instance of other ranks selection.

where F^R and F^{RC} are the outputs of row attention and row & column attention respectively. We normalize F^{RC} into F' by group normalization [44]. In cross-scale attention, we need a pre-processing step to guarantee that the grid sizes of each level are uniform, which enables cross-scale feature interaction in the subsequent steps. As can be seen in Fig. 5(b), we upsample the grids to the size of the largest one. After applying MHSA to features across different scales at equivalent locations, we restore the grid of each scale to their original shape and add a residual connection. The output of cross-scale attention can be expressed as

$$F_{x,y,:}^{RCC} = \text{Downsample}(\text{MHSA}(\text{Upsample}(F'_{x,y,:}))) + F'_{x,y,:} \quad (7)$$

where F^{RCC} is the output of cross-scale attention. Then we normalize F^{RCC} into F'' by group normalization. The terminal component of DPT adopts a convolutional neural network architecture denoted as CLCG. CLCG comprises two 3×3 convolutional layers connected by a LeakyReLU and followed by group normalization. Prior to the group normalization layer, residual connection is implemented. The output of DPT \hat{F} can be computed as

$$\hat{F} = \text{GN}(\text{Conv}(\text{LeakyReLU}(\text{Conv}(F'')))) + F'' \quad (8)$$

The addition of a convolution operation has been proposed as a useful complement to the attention mechanism, as it can better capture local information and enhance feature representation.

4.2.3 Partition Heads. As can be seen in Fig. 2, the partition heads of our proposed approach are designed as N binary saliency classification heads, where N represents the maximum number of ranks. Each partition head is a convolutional layer of $E \times 3 \times 3 \times 1$, where E denotes the number of input channels, it connects to all scales of grid cells from DPT. The output of partition heads is $P = \{p_i | p_i = [v_{i,1}, \dots, v_{i,\sum_{j=2}^6 s_j \times s_j}]^T\}_{i=1}^N$, if we distinguish the salient instances into N ranks, where each logit $v \in [0, 1]$ and s_j is the side length of each grid. In this way, each partition predicts a

set of unordered salient instances that ranks equal or higher than the partition index.

4.3 Partition to Rank

Fig. 6 delineates the process by which the outputs from the mask branch and partition branch are associated to engender the final prediction of saliency ranking.

- Association.** We associate the predicted instance masks and their corresponding partition probabilities on their reference grid cells. This association yields tuples that compose each row of the partition matrix, as illustrated in Fig. 6(a).
- Ambiguity alleviation.** Presented in Fig. 6(b), we discard one instance when 1) the probability of partition- i is below the threshold T as $v_{i,k} < T$; and 2) there exists a probability of partition- j exceed T as $v_{j,k} \geq T$; and 3) $j < i$. The threshold T is set to 0.3 in this paper for best performance.
- Rank-1 selection.** We consider the instance with the largest probability of partition-1 as rank-1, e.g., m_{k+2} with a 0.8 probability of partition-1, as shown in Fig. 6(c). Non-Maximum Suppression (NMS) is then performed on remaining instances, suppressing the instances with an Intersection over Union (IoU) exceeding 0.5.
- Other ranks selection.** In Fig. 6(d), we discard the rank-1 instance from the partition matrix and use the same approach as rank-1 selection to process the remaining instances. This process is repeated until all N salient instances are selected or all instances in the partition matrix are discarded.

Finally, we apply a post-process for the predictions of saliency ranking. A threshold of 0.5 is adopted to convert the predicted soft masks into binary masks.

4.4 Loss Function

In order to obtain the ground truth of the partitions, for each instance mask, we convert the corresponding rank label into a boolean vector of length N . The n -th boolean value in the vector signifies whether the salient instance is prioritized equal to or higher than rank- n .

The training loss function is defined as follows

$$L = \sum_{n=1}^N \lambda_{partition} L_{partition} + \lambda_{mask} L_{mask} \quad (9)$$

We train the partition branch by $L_{partition}$ which is computed by focal loss [20]. L_{mask} is computed by Dice loss for segmentation. In addition, $\lambda_{partition}$ and λ_{mask} are the coefficients for partition loss and mask loss.

5 EXPERIMENTS

5.1 Datasets and Evaluation Metrics

5.1.1 Datasets. We conduct experiments on two publicly accessible datasets ASSR [36] and IRSR [22]. The ASSR dataset employs mouse tracking to collect gaze data, ranks saliency by the first five unique fixated objects. The dataset has 7,464 training, 1,436 validation, and 2,418 test images. The IRSR dataset discards over or under segmented images and those containing over eight or under two

Table 1: Performance comparison with 16 state-of-the-art instance-level segmentation methods on two datasets. Smaller MAE, larger SA – SOR and SOR indicates better performance. The best result is in bold and the second is underlined.

Method	Task	ASSR Dataset			IRSR Dataset			#Para.(M)	FPS
		MAE↓	SA-SOR↑	SOR↑	MAE↓	SA-SOR↑	SOR↑		
Mask2Former [4]	RIS	0.085	0.694	0.862	0.097	0.493	0.833	44.0	30.1
SparseInst [7]		0.103	0.669	0.858	0.112	0.464	0.848	31.6	50.4
SOTR [14]		0.092	0.682	0.867	0.111	0.470	0.823	63.1	11.3
SOLO [40]		0.117	0.643	0.827	0.125	0.465	0.820	36.1	14.7
SOLOv2 [41]		0.096	0.676	0.848	0.114	0.479	0.831	45.0	22.0
Mask R-CNN [15]		0.128	0.624	0.731	0.131	0.437	0.813	43.8	22.8
Cascade R-CNN [1]		0.109	0.661	0.809	0.121	0.452	0.825	76.8	16.2
QueryInst [11]		0.088	0.698	0.832	0.105	0.519	0.833	172.2	15.7
RDPNet [45]	S/C IS	0.119	0.641	0.826	0.134	0.491	0.809	44.3	20.7
SCG [23]		0.107	0.659	0.842	0.123	0.518	0.816	47.9	7.2
OQTR [30]		0.094	0.677	0.865	0.116	0.535	0.831	43.1	25.4
OSFormer [29]		0.118	0.626	0.831	0.112	0.438	0.817	46.6	26.3
ASSR [36]	SR	0.104	0.661	0.787	0.134	0.377	0.710	43.9	22.1
IRSR [22]		0.105	0.705	0.813	<u>0.088</u>	<u>0.564</u>	<u>0.806</u>	102.4	18.5
SOR [10]		<u>0.083</u>	0.717	0.836	0.091	0.543	0.797	100.0	36.1
OCOR [38]		0.085	<u>0.723</u>	<u>0.877</u>	-	-	-	-	-
PSR (proposed)		0.075	0.738	0.892	0.080	0.651	0.878	50.9	20.9
OCOR-Swin-L [38]	SR	<u>0.078</u>	<u>0.738</u>	<u>0.904</u>	<u>0.079</u>	<u>0.578</u>	<u>0.834</u>	-	-
PSR-Swin-L		0.071	0.746	0.915	0.075	0.664	0.890	225.3	7.7

salient instances. This dataset includes 6,059 training and 2,929 test images.

5.1.2 Evaluation Metrics. We utilize the identical evaluation methods as [22, 36] for an equitable comparison, namely Salient Object Ranking (SOR), Segmentation-Aware SOR (SA-SOR), and Mean Absolute Error (MAE). SOR computes the Spearman’s rank-order correlation between the predictions and GT. On the other hand, SA-SOR filters the objects at the instance-level using the mask Intersection over Union (IoU) and employs the Pearson correlation coefficient to measure the linear correlation between the prediction and GT. The MAE metric compares the average per-pixel difference between the predictions and GT, accounting for both segmentation and ordering quality as a comprehensive metric.

5.2 Implementation Details

Following preceding works [10, 17, 22, 36], we adopt ResNet-50 [16] weights pretrained on the MS-COCO [21] 2017 training split and resize input images to 640×480 . Besides, our model is trained on the ASSR dataset training set. The partition branch comprises 3 convolutional layers and 3 transformer layers. The weights of the partition loss and mask loss are 1 and 3, respectively. Furthermore, we implement the stochastic gradient descent (SGD) optimizer with a learning rate of $2.5e-5$ and employ a warm-up strategy in the initial 1,000 iterations. We train the model for 60 epochs with a batch size of 4 on a NVIDIA RTX 3090 GPU. We perform multi-step decay with a decay factor of $1e-4$ at the 42nd and 54th epochs, respectively.

5.3 Comparisons with the State-of-the-arts

5.3.1 Quantitative Comparison. As shown in Tab. 1, our model is compared with 16 state-of-the-art (SOTA) methods, their original tasks include regular instance segmentation (RIS) [1, 4, 7, 11, 14, 15, 40, 41], salient/camouflage instance segmentation (S/C IS) [23, 29, 30, 45] and saliency ranking (SR) tasks [10, 22, 36, 38]. Unless otherwise specified, all experiments reported in Tab. 1 are conducted using the ResNet-50 backbone to ensure fair comparison. We use ‘-’ to indicate unavailable data due to the code not being open-sourced.

On the ASSR dataset, the proposed PSR surpasses SOTAs by 0.008, 0.015 and 0.015 in terms of MAE, SA-SOR and SOR metrics, respectively. For the IRSR dataset, PSR surpasses state-of-the-arts by 0.008, 0.087 and 0.072 in terms of MAE, SA-SOR and SOR metrics, respectively. Specially, the last two rows of Tab. 1 present a comparison between the proposed PSR and the latest OCOR [38] using the Swin-L [25] backbone. All the results indicate that PSR outperform previous methods in terms of all three metrics on both two datasets. Despite including multiple partition heads and the DPT, our FPS only drops by 1.1 compared to the base model, approaching the speed of ASSR and IRSR.

5.3.2 Qualitative Comparison. In Fig. 7, we present the visualization results for qualitative analysis between PSR and other instance-level saliency ranking methods include: ASSR [36], IRSR [22] and SOR [10]. Our method produces salient instance masks and ranks closer to GT when dealing with objects with ambiguous rank levels in complex scenes. For example, in the first column, though all models predict correct masks of salient instances, ranks of the



Figure 7: Visual comparison between the proposed PSR and other instance-level saliency ranking methods. Our PSR improves both mask and ranking precision compared to ASSR [36], IRSR [22] and SOR [10].

right two people are hard to determine, only PSR produces the correct ranking order. In the second column, the size of sheep is variant, previous works do not consider global cross-scale feature interaction and thus gets the incorrect ranking order.

5.3.3 Confusion Matrix Comparison. The confusion matrices of the proposed PSR and other instance-level saliency ranking methods [10, 22, 36] are reported on Fig. 8. The principal diagonal line of each matrix reflects the number of true positives for each method. The proposed PSR obtains a summation of 4282, larger than that of 2314, 3134 and 3732 for ASSR [36], IRSR [22] and SOR [10] respectively. More importantly, with regard to the rank-5 and rank-4, the number of true positives of PSR reaches 248 and 403 respectively, achieving a relative improvement with 67.6% and 36.6% compared to SOR [10]. As the experimental phenomenon, we can see that 1) our PSR outperforms previous models in terms of the overall accuracy; 2) specially, PSR surpass previous methods by a large margin on difficult samples, such as rank-5 and rank-4.

5.4 Ablation Studies

5.4.1 Attention-based Interaction Modules. The twin transformer [14] lacks of cross-scale interaction and the deformable transformer [51] lacks global interaction. Hence, we design the DPT to enable comprehensive global cross-scale interaction. We conduct experiments to compare with twin transformer, deformable transformer and all-scale transformer to demonstrate the effectiveness of DPT.

As depicted in the top section of Tab. 2, we attempt to align the FPS by adjusting the number of layers in the twin transformer and deformable transformer. Among the three methods with similar

		Predicted label					
		bg	5	4	3	2	1
True label	bg	0	474	604	615	404	203
	5	904	66	109	154	176	135
	4	950	62	176	228	227	182
	3	968	67	179	304	340	290
	2	830	62	131	301	642	448
	1	558	38	93	186	413	1126

(a) ASSR

		Predicted label					
		bg	5	4	3	2	1
True label	bg	0	954	918	790	506	195
	5	1159	146	116	59	40	25
	4	1106	233	205	152	88	43
	3	946	251	284	348	240	81
	2	559	108	177	416	807	351
	1	215	18	32	118	407	1628

(b) IRSR

		Predicted label					
		bg	5	4	3	2	1
True label	bg	0	858	793	634	374	185
	5	1006	148	210	115	50	16
	4	975	130	295	287	104	36
	3	892	71	203	556	336	92
	2	645	38	88	253	1046	348
	1	312	21	20	58	320	1687

(c) SOR

		Predicted label					
		bg	5	4	3	2	1
True label	bg	0	865	607	377	223	95
	5	939	248	196	103	44	15
	4	753	250	403	276	102	43
	3	518	191	333	691	329	88
	2	239	101	187	387	1159	345
	1	76	19	49	123	370	1781

(d) PSR

Figure 8: Confusion matrix of the proposed PSR and other instance-level saliency ranking methods on ASSR dataset.

inference speed, the network employing DPT achieved the best performance.

Table 2: Comparing the performance of different transformer-based interaction modules.

Transformers	#Layers	MAE↓	SA-SOR↑	SOR↑	FPS↑
Twin [14]	4	0.080	0.719	0.883	20.7
Deformable [51]	4	0.080	0.720	0.877	21.2
Dense Pyramid	3	0.075	0.738	0.892	20.9
All-scale	3	0.079	0.726	0.885	13.1
Dense Pyramid	3	0.075	0.738	0.892	20.9

Table 3: Comparing the performance of layer variations of convolutions and transformers in partition branch.

#Conv Layers	#Transformer Layers	MAE↓	SA-SOR↑	SOR↑
-	6	0.080	0.725	0.877
3	3	0.075	0.738	0.892
6	-	0.083	0.720	0.873

As shown in the bottom part of Tab. 2, we align the number of transformer layers with 3. Though it reduces the computational burden, the proposed DPT outperform all-scale transformer in terms of MAE, SA-SOR, SOR, as well as FPS.

5.4.2 Layer Variations in Partition Branch. We exam the impact of different numbers of convolutional layers and transformer layers in partition branch. As shown in Tab. 3, the network achieves better performance when convolutional layers are followed by transformer layers. We ultimately adopt a setting of 3 convolutional layers and 3 transformer layers because it achieves best performance.

5.4.3 Contribution of Different Components. To ascertain the efficacy of the Cross-Scale Attention (C-S Attn) and Row & Column Attention (R&C Attn) in DPT, as well as the proposed ranking by partition paradigm (Partition), we devise ablation experiments on the constituent components. The results of the ablation study are presented in Tab. 4, with SOLOv2 [41] serving as the base model, listed in the first column. Comparing to the baseline, by leveraging both Cross-Scale and Row & Column Attention, the performance improved by 0.010, 0.021 and 0.018 in terms of MAE, SA-SOR and SOR respectively. Employing only ranking by partition also improves the MAE, SA-SOR and SA-SOR by 0.016, 0.037 and 0.023, as this new paradigm mitigates the ambiguity challenges of ordering scenes. The overall results indicate that our method achieves the best performance when all components are included, which demonstrates the effectiveness and necessity of each component.

5.5 Generalization Experiment of Ranking by Partition

We propose a ranking by partition paradigm for the task of saliency ranking that can mitigate the effects of label ambiguity. The ranking by partition paradigm, in principle, can be applied to other models and enhance their performance on saliency ranking. We compare different models using ranking by partition versus the traditional ranking by sorting. As shown in Tab. 5, the ranking by partition

Table 4: Comparing the contribution of different components.

C-S Attn	R&C Attn	Partition	MAE↓	SA-SOR↑	SOR↑
			0.096	0.676	0.848
	√		0.092	0.688	0.853
√	√		0.086	0.697	0.866
		√	0.080	0.713	0.871
	√	√	0.077	0.724	0.884
√	√	√	0.075	0.738	0.892

Table 5: Generalization Experiment of Ranking by Partition

Method	Partition	MAE↓	SA-SOR↑	SOR↑
SOR [10]		0.083	0.637	0.836
	√	0.078 _{-0.005}	0.657 _{+0.020}	0.859 _{+0.023}
SOTR [14]		0.092	0.682	0.867
	√	0.081 _{-0.011}	0.711 _{+0.029}	0.893 _{+0.026}
M2F [4]		0.085	0.694	0.862
	√	0.077 _{-0.008}	0.721 _{+0.027}	0.880 _{+0.018}
SOLOv2 [41]		0.096	0.676	0.848
	√	0.080 _{-0.016}	0.713 _{+0.037}	0.871 _{+0.023}

paradigm improves the performance of saliency ranking not only on CNN-based models [41], but also on CNN-transformer-hybrid-based models [10, 14] and mask-classification-based models [4], demonstrating its generalizability.

6 CONCLUSIONS

In this paper, we propose a ranking by partition paradigm to alleviate saliency ambiguity in the saliency ranking task. The ranking by partition paradigm can be applied to other related models to boost saliency ranking. We also proposed a dense pyramid transformer (DPT) to facilitate cross-scale global feature interaction for saliency ranking. By leveraging the proposed ranking by partition and DPT, our partitioned saliency ranking (PSR) outperforms state-of-the-art methods.

7 ACKNOWLEDGMENTS

This work was supported by the National Natural Science Foundation of China Grant 61902139.

REFERENCES

- [1] Zhaowei Cai and Nuno Vasconcelos. 2018. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. IEEE, Salt Lake City, UT, USA, 6154–6162.
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*. Springer, 213–229.
- [3] Cuiqun Chen, Mang Ye, Meibin Qi, Jingjing Wu, Yimin Liu, and Jianguo Jiang. 2022. Saliency and granularity: Discovering temporal coherence for video-based person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology* 32, 9 (2022), 6100–6112.
- [4] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. 2022. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, New Orleans, LA, USA, 1290–1299.
- [5] Bowen Cheng, Alex Schwing, and Alexander Kirillov. 2021. Per-pixel classification is not all you need for semantic segmentation. *Advances in Neural Information Processing Systems* 34 (2021), 17864–17875.
- [6] Ming-Ming Cheng, Niloy J Mitra, XiaoLei Huang, Philip HS Torr, and Shi-Min Hu. 2014. Global contrast based salient region detection. *IEEE transactions on pattern analysis and machine intelligence* 37, 3 (2014), 569–582.
- [7] Tianheng Cheng, Xinggang Wang, Shaoyu Chen, Wenqiang Zhang, Qian Zhang, Chang Huang, Zhaoxiang Zhang, and Wenyu Liu. 2022. Sparse instance activation for real-time instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4433–4442.
- [8] Lifeng Fan, Wenguan Wang, Siyuan Huang, Xinyu Tang, and Song-Chun Zhu. 2019. Understanding human gaze communication by spatio-temporal graph reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5724–5733.
- [9] Ruochen Fan, Ming-Ming Cheng, Qibin Hou, Tai-Jiang Mu, Jingdong Wang, and Shi-Min Hu. 2019. S4net: Single stage salient-instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 6103–6112.
- [10] Hao Fang, Daoxin Zhang, Yi Zhang, Minghao Chen, Jiawei Li, Yao Hu, Deng Cai, and Xiaofei He. 2021. Salient object ranking with position-preserved attention. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 16331–16341.
- [11] Yuxin Fang, Shusheng Yang, Xinggang Wang, Yu Li, Chen Fang, Ying Shan, Bin Feng, and Wenyu Liu. 2021. Instances as queries. In *Proceedings of the IEEE/CVF international conference on computer vision*. 6910–6919.
- [12] Mengyang Feng, Huchuan Lu, and Errui Ding. 2019. Attentive feedback network for boundary-aware salient object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 1623–1632.
- [13] Shang-Hua Gao, Yong-Qiang Tan, Ming-Ming Cheng, Chengze Lu, Yunpeng Chen, and Shuicheng Yan. 2020. Highly efficient salient object detection with 100k parameters. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI*. Springer, 702–721.
- [14] Ruohao Guo, Dantong Niu, Liao Qu, and Zhenbo Li. 2021. Sotr: Segmenting objects with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 7157–7166.
- [15] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*. 2961–2969.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [17] Md Amirul Islam, Mahmoud Kalash, and Neil DB Bruce. 2018. Revisiting salient object detection: Simultaneous detection, ranking, and subitizing of multiple salient objects. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7142–7150.
- [18] Guanbin Li, Yuan Xie, Liang Lin, and Yizhou Yu. 2017. Instance-level salient object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2386–2395.
- [19] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2117–2125.
- [20] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. 2017. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*. 2980–2988.
- [21] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*. Springer, 740–755.
- [22] Nian Liu, Long Li, Wangbo Zhao, Junwei Han, and Ling Shao. 2021. Instance-level relative saliency ranking with graph reasoning. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44, 11 (2021), 8321–8337.
- [23] Nian Liu, Wangbo Zhao, Ling Shao, and Junwei Han. 2021. SCG: Saliency and contour guided salient instance segmentation. *IEEE Transactions on Image Processing* 30 (2021), 5862–5874.
- [24] Tie Liu, Zejian Yuan, Jian Sun, Jingdong Wang, Nanning Zheng, Xiaou Tang, and Heung-Yeung Shum. 2010. Learning to detect a salient object. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33, 2 (2010), 353–367.
- [25] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*. 10012–10022.
- [26] Jonathan Long, Evan Shelhamer, and Trevor Darrell. 2015. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3431–3440.
- [27] Lufan Ma, Tiancai Wang, Bin Dong, Jiangpeng Yan, Xiu Li, and Xiangyu Zhang. 2021. Implicit feature refinement for instance segmentation. In *Proceedings of the 29th ACM International Conference on Multimedia*. 3088–3096.
- [28] Namuk Park and Songkuk Kim. 2022. How do vision transformers work? *arXiv preprint arXiv:2202.06709* (2022).
- [29] Jialun Pei, Tianyang Cheng, Deng-Ping Fan, He Tang, Chuanbo Chen, and Luc Van Gool. 2022. Osformer: One-stage camouflaged instance segmentation with transformers. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVIII*. Springer, 19–37.
- [30] Jialun Pei, Tianyang Cheng, He Tang, and Chuanbo Chen. 2022. Transformer-based efficient salient instance segmentation networks with orientative query. *IEEE Transactions on Multimedia* (2022).
- [31] Jialun Pei, He Tang, Wanru Wang, Tianyang Cheng, and Chuanbo Chen. 2022. Salient instance segmentation with region and box-level annotations. *Neurocomputing* 507 (2022), 332–344.
- [32] Zheyun Qin, Xiankai Lu, Xiushan Nie, Xiantong Zhen, and Yilong Yin. 2021. Learning hierarchical embedding for video instance segmentation. In *Proceedings of the 29th ACM International Conference on Multimedia*. 1884–1892.
- [33] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems* 28 (2015).
- [34] Brandon Richard Webster, Brian Hu, Keith Fieldhouse, and Anthony Hoogs. 2022. Doppelganger Saliency: Towards More Ethical Person Re-Identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2847–2857.
- [35] Anil K Seth and Tim Bayne. 2022. Theories of consciousness. *Nature Reviews Neuroscience* 23, 7 (2022), 439–452.
- [36] Avishek Siris, Jianbo Jiao, Gary KL Tam, Xianghua Xie, and Rynson WH Lau. 2020. Inferring attention shift ranks of objects for image saliency. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12133–12143.
- [37] Avishek Siris, Jianbo Jiao, Gary KL Tam, Xianghua Xie, and Rynson WH Lau. 2021. Scene context-aware salient object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4156–4166.
- [38] Xin Tian, Ke Xu, Xin Yang, Lin Du, Baocai Yin, and Rynson WH Lau. 2022. Bi-directional object-context prioritization learning for saliency ranking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5882–5891.
- [39] Xin Tian, Ke Xu, Xin Yang, Baocai Yin, and Rynson WH Lau. 2020. Weakly-supervised salient instance detection. *arXiv preprint arXiv:2009.13898* (2020).
- [40] Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. 2020. Solo: Segmenting objects by locations. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16*. Springer, 649–665.
- [41] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. 2020. Solov2: Dynamic and fast instance segmentation. *Advances in Neural Information Processing Systems* 33 (2020), 17721–17732.
- [42] Zheng Wang, Xinyu Yan, Yahong Han, and Meijun Sun. 2019. Ranking video salient object detection. In *Proceedings of the 27th ACM International Conference on Multimedia*. 873–881.
- [43] Jun Wei, Shuhui Wang, Zhe Wu, Chi Su, Qingming Huang, and Qi Tian. 2020. Label decoupling framework for salient object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 13025–13034.
- [44] Yuxin Wu and Kaiming He. 2018. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*. 3–19.
- [45] Yu-Huan Wu, Yun Liu, Le Zhang, Wang Gao, and Ming-Ming Cheng. 2021. Regularized densely-connected pyramid network for salient instance segmentation. *IEEE Transactions on Image Processing* 30 (2021), 3897–3907.
- [46] Chuan Yang, Lihe Zhang, Huchuan Lu, Xiang Ruan, and Ming-Hsuan Yang. 2013. Saliency detection via graph-based manifold ranking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3166–3173.
- [47] Miao Zhang, Tingwei Liu, Yongri Piao, Shunyu Yao, and Huchuan Lu. 2021. Auto-msfnet: Search multi-scale fusion network for salient object detection. In *Proceedings of the 29th ACM international conference on multimedia*. 667–676.
- [48] Jia-Xing Zhao, Jiang-Jiang Liu, Deng-Ping Fan, Yang Cao, Jufeng Yang, and Ming-Ming Cheng. 2019. EGNNet: Edge guidance network for salient object detection.

- In *Proceedings of the IEEE/CVF international conference on computer vision*. 8779–8788.
- [49] Zhirui Zhao, Changqun Xia, Chenxi Xie, and Jia Li. 2021. Complementary trilateral decoder for fast and accurate salient object detection. In *Proceedings of the 29th acm international conference on multimedia*. 4967–4975.
- [50] Tun Zhu, Daoxin Zhang, Yao Hu, Tianran Wang, Xiaolong Jiang, Jianke Zhu, and Jiawei Li. 2021. Horizontal-to-vertical video conversion. *IEEE Transactions on Multimedia* 24 (2021), 3036–3048.
- [51] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. 2020. Deformable DETR: Deformable Transformers for End-to-End Object Detection. In *International Conference on Learning Representations*.