# Rethinking Lightweight Salient Object Detection via Network Depth-Width Tradeoff

Jia Li, *Senior Member, IEEE*, Shengye Qiao, Zhirui Zhao, Chenxi Xie and Xiaowu Chen, *Senior Member, IEEE*, Changqun Xia

*Abstract*—Existing salient object detection methods often adopt deeper and wider networks for better performance, resulting in heavy computational burden and slow inference speed. This inspires us to rethink saliency detection to achieve a favorable balance between efficiency and accuracy. To this end, we design a lightweight framework while maintaining satisfying competitive accuracy. Specifically, we propose a novel trilateral decoder framework by decoupling the U-shape structure into three complementary branches, which are devised to confront the dilution of semantic context, loss of spatial structure and absence of boundary detail, respectively. Along with the fusion of three branches, the coarse segmentation results are gradually refined in structure details and boundary quality. Without adding additional learnable parameters, we further propose Scale-Adaptive Pooling Module to obtain multi-scale receptive filed. In particular, on the premise of inheriting this framework, we rethink the relationship among accuracy, parameters and speed via network depth-width tradeoff. With these insightful considerations, we comprehensively design shallower and narrower models to explore the maximum potential of lightweight SOD. Our models are purposed for different application environments: 1) a tiny version CTD-S (1.7M, 125FPS) for resource constrained devices, 2) a fast version CTD-M (12.6M, 158FPS) for speed-demanding scenarios, 3) a standard version CTD-L (26.5M, 84FPS) for high-performance platforms. Extensive experiments validate the superiority of our method, which achieves better efficiency-accuracy balance across five benchmarks.

*Index Terms*—Salient object detection, lightweight framework, trilateral decoder

## I. INTRODUCTION

IN recent years, salient object detection (SOD) [2], [3] has made great progress with the development of Deep Neural Networks (DNNs). As an efficient preprocessing technique, SOD plays an important role in many downstream computer vision tasks, such as image retrieval [4], visual tracking [5] and semantic segmentation [6].

Jia Li is with the State Key Laboratory of Virtual Reality Technology and Systems, School of Computer Science and Engineering, Beihang University, Beijing, 100191, China, and also with Peng Cheng Laboratory, Shenzhen, 518000, China.

Shengye Qiao, Zhirui Zhao and Chenxi Xie are with the State Key Laboratory of Virtual Reality Technology and Systems, School of Computer Science and Engineering, Beihang University, Beijing, 100191, China.

Xiaowu Chen is with the State Key Laboratory of Virtual Reality Technology and Systems, School of Computer Science and Engineering, Beihang University, Beijing, 100191, China, and also with Peng Cheng Laboratory, Shenzhen, 518000, China.

Changqun Xia is with Peng Cheng Laboratory, Shenzhen, 518000, China. Corresponding author: Changqun Xia(E-mail: xiachq@pcl.ac.cn)

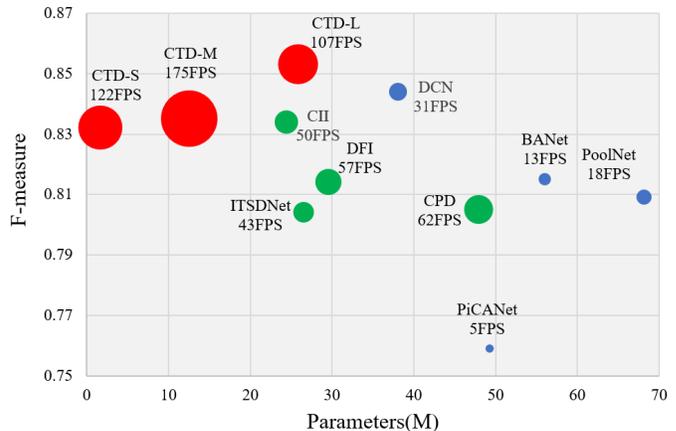A preliminary version of this work has appeared in ACM MM 2021 [1]



Fig. 1. Comparisons of our proposed CTDNet with other ResNet-based SOD models in accuracy, parameters and speed. We calculate $mF_\beta$ on DUTS-TE dataset as an example. The green circles represent real-time SOD methods, blue circles represent non-real-time SOD methods, and red circles represent our method. The size of circle indicates speed and larger circle indicates faster speed.

Many SOD methods have greatly benefited from very deep and wide models and achieves remarkable results. For example, BANet [7] and DCN [8] reach high accuracy compared to other models as shown in Fig. 1. However, the success comes at a price of heavy computation burden and slow running speed. These models increase the network depth and width by adjusting the number of layers and channels, which brings tremendous parameters and calculations. Taking these factors into consideration, a question arises: *Is shallow and narrow network able to achieve comparable performance of large counterparts?*

Some researches start to answer this question and try to compromise between efficiency and accuracy, such as GC-PANet [9] and ITSDNet [10]. These methods can lead to improvements in efficiency but still adopt strong classification models pretrained on ImageNet(e.g., ResNet-50 and ResNet-101) as backbone, which tend to bring about a large amount of parameters. For example, ITSDNet contains about 27M parameters and the ones of GCPANet are approaching 67M. Besides, since SOD does not care about the category of objects, the features extracted by overly deep networks may be overqualified. To this end, some methods have adopted efficient backbones when designing networks and shown great potential for establishing highly competitive models with fewer parameters and faster speed. For example, EDN
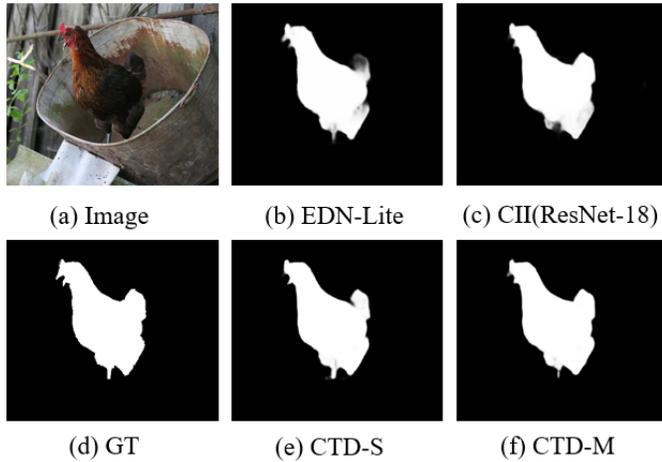
Fig. 2. The saliency maps of the same image from EDN-Lite, CII(ResNet-18), CTD-S and CTD-M

[11] using MobileNet-V2 [12] as the backbone only contains 1.8M parameters with fewer channels, and it has comparable performance due to inverted block and depth-wise separable convolutions. CII [13] using ResNet-18 [14] as the backbone is a representative shallow network with a depth of just 17 layers, which significantly improves the processing speed. Therefore, making full use of advantages of these efficient backbones to solve SOD problem is a way worth exploring and trying.

Besides backbone, the complexity of the framework itself is also an important factor affecting the efficiency of the model. It is worth noting that most state-of-the-art SOD methods belong to the U-shape structure [15], where the input image is gradually down-sampled to extract multi-scale features and then up-sampled to recover the resolution with skip connection, as shown in Fig. 3(a). Except for strong backbone models in the encoder, these methods design complex structures in the decoder to refine saliency maps, but inevitably introduce additional parameters. In addition, there are several problems for the U-shaped structure. First, spatial information is seriously lost during the process of frequent down-sampling and cannot be perfectly recovered by integrating the hierarchical features [16], which leads to incomplete structure details. Second, semantic and global context information may be gradually diluted in the top-down path, which leads to inaccurate object location. Third, boundary information is also ignored, which leads to poor boundary quality. To address these challenges, it is necessary to design a new and effective method to deal with the dilution of semantic context, loss of spatial structure and absence of boundary detail.

With the above considerations in mind, we propose a lightweight framework that performs competitively against other large counterparts. Instead of conventional backbone models, we rely on efficient backbones for feature extraction. We take advantage of these backbones more effectively to encode hierarchical features from low-level, mid-level to high-level, which can fully meet the requirements of SOD task at minimal computation cost. In addition to the encoder, we propose a novel Complementary Trilateral Decoder (CTD) that decouples the U-shape structure into three branches: Semantic

Path, Spatial Path and Boundary Path. As illustrated in Fig. 3(b), these three branches are derived from different levels of the shared encoder and complementary to each other. Along with the fusion of three branches, the coarse segmentation results are refined in structure details and boundary quality. Specifically, Semantic Path is designed to capture rich semantic context and global context of high-level features, which can form initial coarse saliency maps with accurate locations of salient objects. Spatial Path is introduced to preserve more spatial details of mid-level features and then is combined with Semantic Path by the proposed Cross Aggregation Module (CAM), which can produce relative fine saliency maps with precise structures of salient objects. Boundary Path is utilized to extract salient boundary contour by fusing low-level local features and high-level location features with an explicit edge supervision. Finally, the output of CAM and Boundary Path are merged by the proposed Boundary Refinement Module (BRM) to further refine boundary, which can generate final finer saliency maps with clear boundaries of salient objects. Besides, although semantic category features are unnecessary for SOD task, the multi-scale representation ability of high-level features is indispensable. Therefore, we take the advantages of pooling operation and propose a novel structure named Scale-Adaptive Pooling (SAP) to capture multi-scale context information from multiple receptive fields without additional learnable parameters.

More importantly, on the basis of inheriting this framework, we rethink the relationship among accuracy, parameters and speed according to the different depth and width characteristics of efficient backbones, thus discover more potential of lightweight SOD. To facilitate the practical application in different environments, we provide three versions: 1) a tiny version CTD-S (1.7M, 105FPS) for resource-constrained devices, 2) a fast version CTD-M (12.6M, 160FPS) for speed-demanding scenarios, 3) a standard version CTD-L (25.9M, 100FPS) for high-performance platforms. Extensive experiments on five popular SOD datasets demonstrate the generality and superiority of our method, which achieves better balance between efficiency and accuracy.

In general, the main contributions of our work are summarized as follows:

1) We rethink SOD task from the perspective of network depth and width, and then design lightweight and efficient saliency detection models with shallow and narrow networks, which can perform surprisingly well against other large counterparts with less parameters and higher efficiency.

2) We propose a novel Complementary Trilateral Decoder (CTD) framework that decouples the U-shape structure into three branches: Semantic Path, Spatial Path and Boundary Path. The Cross Aggregation Module (CAM) and Boundary Refinement Module (BRM) are constructed to gradually merge these three complementary branches according to "coarse-fine-finer" strategy, which significantly improves the region accuracy and boundary quality.

3) We take the advantages of pooling operation to enhance the multi-scale representation ability of high-level features and propose a novel Scale-Adaptive Pooling (SAP) structure to obtain multi-scale receptive filed without increasing parameters.
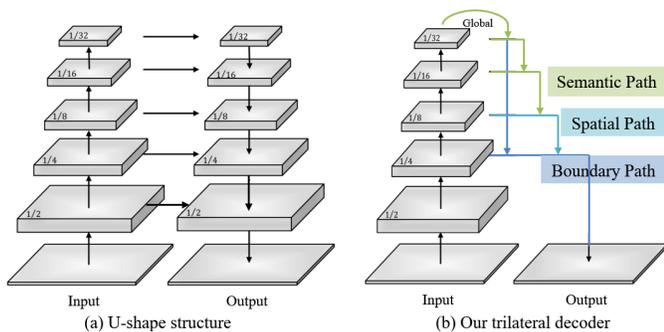
Fig. 3. The U-shape structure gradually recovers the spatial information by leveraging lateral connections and topdown path in the decoder part, while our trilateral decoder separately treats semantic context, spatial detail and boundary information with three branches.

4) We provide three versions for different application environments and conduct extensive experiments on 5 benchmarks to demonstrate the generality and superiority of our method, which achieves a favorable trade-off between efficiency and accuracy.

Compared with the conference version of CTDNet, we optimize the parameters and accuracy from the perspective of width and depth under the premise of being lightweight, and propose three versions of the model based on complementary trilateral decoder framework: CTD-S, CTD-M and CTD-L. Parameters and speeds of different magnitudes in these three models can be more comprehensively adapted to various devices and scenarios. In addition to the original three branches and the corresponding fusion module, we introduce the SAP structure in the decoder. Meanwhile, the decoding process of the three models is also significantly different from the conference version, which will be explained in detail in the methodology section.

## II. RELATED WORK

### A. Accurate SOD

Traditional SOD methods mostly rely on heuristic priors (e.g., color, texture and contrast) to generate saliency maps. However, these hand-crafted features can hardly capture high-level information, which are not robust enough for complex scenarios [17]. In recent years, many FCN-based methods have achieved remarkable progress thanks to its powerful representation capability. As one of the most representative networks, the U-shape structure has been widely followed for accurate saliency detection. PiCANet [18] proposed a pixel-wise contextual attention network to learn informative context locations for each pixel. BMPM [19] designed a bi-directional message passing model for better feature selection and integration. MINet [20] focused on scale variation and class imbalance challenges by utilizing multi-level and multi-scale features. PSGLoss [21] introduced an additional progressive self-guided loss and multi-scale feature aggregation module with branch-wise attention to detect salient objects completely and effectively. Some methods introduce an additional boundary-aware branch or a boundary-aware loss function for fine object boundaries. C2SNet [22] presented a contour-to-saliency transferring model that predicts contours and saliency maps simultaneously. BASNet [23] proposed a boundary-aware model and designed hybrid loss to make full use of boundary information. EGNet [24] focused on the complementary information modeling between salient edge and salient object to improve boundaries and localization. BANet [7] designed a boundary-aware model with successive dilation from the perspective of selectivity and invariance. AFNet [25] proposed a multi-scale attentive feedback model and Boundary-Enhanced Loss to predict salient objects with entire structure and exquisite boundaries. PurNet [26] utilized promotion and rectification attention to purify salient objects and introduced a structural similarity loss to restore the complex or fine structures of salient objects. In addition, some models are used for accurate segmentation of different scenes, such as automatic driving [27] and 360 panoramic scenes [28]. However, these methods have brought huge parameters and model complexity for better performance, resulting in slow inference speed.

### B. Efficient SOD

SOD serves as a preprocessing technique for many downstream vision tasks, so accuracy and efficiency are both important factors when building a SOD network. Recently, some methods have been proposed for accurate and fast saliency detection. For example, RAS [29] predicted a global saliency map and proposed reverse attention to guide side-output residual features recursively. PoolNet [30] fully exploited the pooling operations based on the FPN structure for real-time salient object detection. CPD [31] discarded features of shallow layers for acceleration and proposed a cascade partial decoder that utilizes attention mechanism to refine high-level features. ITSDNet [10] proposed an interactive two-stream decoder to explore multiple cues, including saliency, contour and their correlation. DCN [8] leveraged skeleton and edge information to model interiors and boundaries of salient objects together in decomposition and completion framework. CII [13] designed a centralized information interaction strategy to encode the cross-scale information into the learnable filters and modeled information with relative global calibration module. In addition to these resnet-based methods, there are methods using lightweight backbone to improve efficiency. For example, EDNLite [11], adopting MobileNet V2 as backbone model, designed an extremely-downsampled block to learn a global view of the whole image and constructed a scale-correlated pyramid convolution for effective feature fusion in the decoder. In order to design and achieve the ultimate lightweight model, SAMNet [32] abandoned the wildly used pre-trained Convolutional Neural Network (CNN) [33] backbones thoroughly and rebuilt a lightweight encoder-decoder architecture with stereoscopically attentive multi-scale module, which adopted stereoscopic attention mechanism for effective and efficient multi-scale learning. Although smaller and faster than the previous large models, these methods cannot achieve comparable performance.
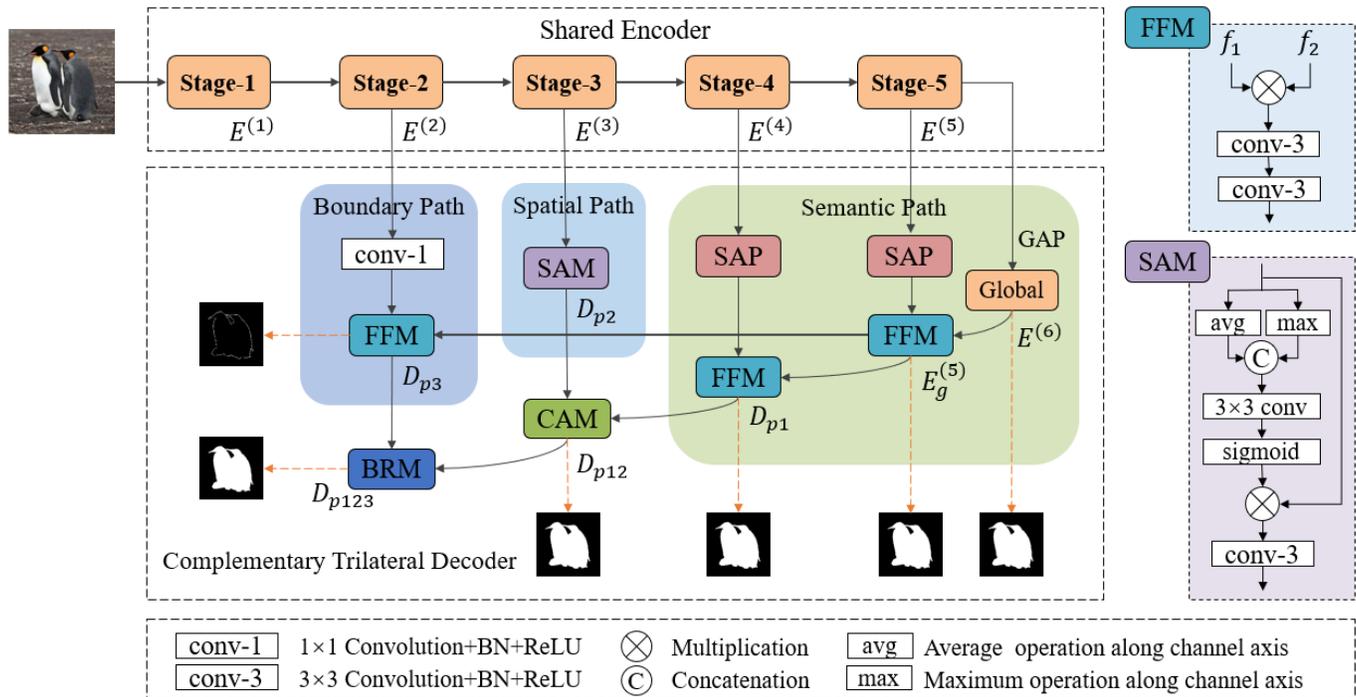
Fig. 4. The framework of our proposed Complementary Trilateral Decoder (CTD) Network taking CTD-M as an example with three branches: Semantic Path, Spatial Path and Boundary Path, which treats semantic context, spatial detail and boundary information separately in the decoder part. The three parts share the same encoder and are derived from different stages of the encoder. The three branches are complementary to each other and we design three specific fusion modules to gradually merge them according to "coarse-fine-fine" strategy

## C. Encoder-Decoder Structure

Recently, many researches follow the U-shape structure to effectively combine low-level and high-level features for saliency detection. TDBU [34] learnt top-down and bottom-up saliency inference in a cooperative and iterative manner. MINet [20] focused on scale variation and class imbalance challenges by utilizing multi-level and multi-scale feature information. DASNet [35] proposed a depth-aware framework to improve the segmentation performance with depth constraints. PFSNet [36] proposed to aggregate adjacent feature nodes in pairs through layer by layer shrinkage, which can fuse details and semantics effectively, and discard interference information. However, these methods based on U-shape structure have brought high model complexity to achieve better performance, resulting in slow inference speed.

Different from these methods, we design a less complex structure via network depth-width tradeoff. Apart from adopting efficient backbones as encoder, we introduce complementary trilateral decoder with cross aggregation module, boundary refinement module and scale-adaptive pooling structure to maintain satisfying competitive accuracy while ensuring the lightness of the model.

## III. METHODOLOGY

### A. Motivation and Framework

As mentioned above, conventional SOD methods adopt strong backbone models to encode deep semantic information (i.e., category), which inherits a large number of parameters

and calculations. However, SOD is a category-insensitive task that focuses on segmenting the salient objects or regions in an image, so it may be unnecessary to extract rich category features from very deep networks. In addition, SOD serves as a preprocessing technique for many downstream vision tasks, so accuracy and efficiency are both important factors when building a SOD network. To this end, we explore the potential of shallow and narrow models, and then specially design a lightweight framework while maintaining satisfying accuracy. The key to the design concept lies in the following aspects: 1) We rely on shallow and narrow backbone models to encode hierarchical features fast while reducing computational burden; 2) We prune the number of feature channels for less redundant information and faster inference speed; 3) We take the advantages of pooling operation to enhance multi-scale representation ability of model; 4) We propose a novel framework for accurate and efficient saliency detection.

Taking into account the drawbacks of the U-shape structure mentioned above, we propose a novel framework that treats semantic context, spatial detail and boundary contour separately in the decoder. The whole architecture of our proposed Complementary Trilateral Decoder (CTD) is shown in Fig. 4.

Existing backbones, like ResNet and MobileNet, can encode images in multiple stages and output corresponding feature maps, which can be expressed as $\{E^1, E^2, E^3, E^4, E^5\}$ for convenience. In order to make full use of the feature maps of different stages in the decoding process, we decouple the decoder of the U-shape structure into three branches: Semantic Path, Spatial Path and Boundary Path, which are
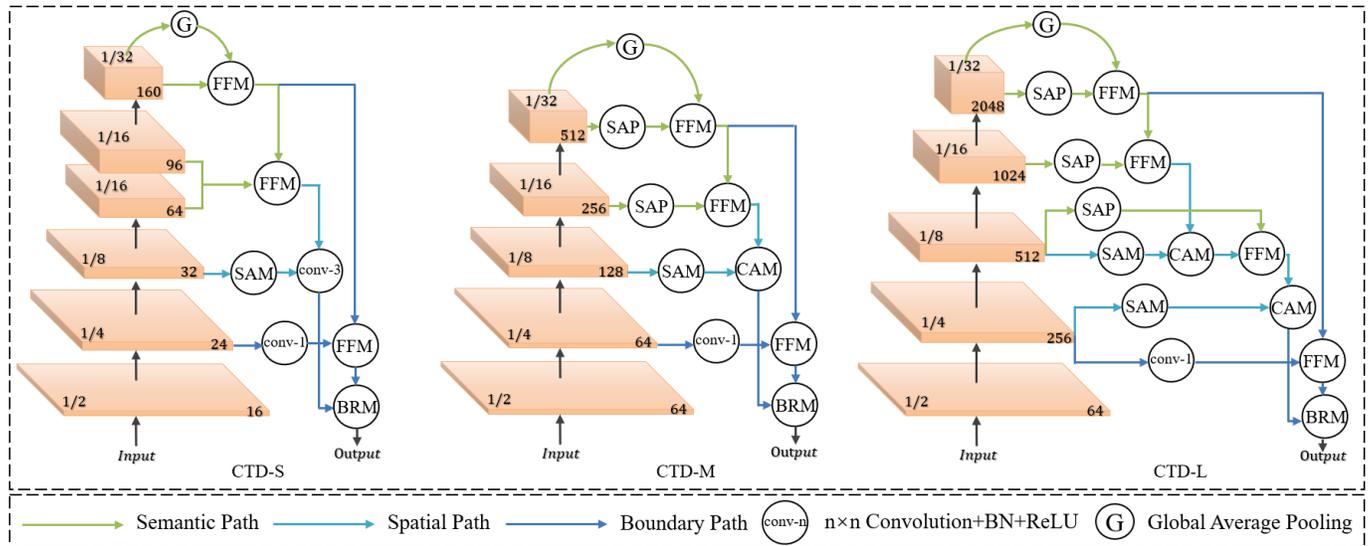
Fig. 5. Comparison of framework details on CTD-S, CTD-M and CTD-L, which adopt MobileNet-V2, ResNet-18 and ResNet-50 as backbone respectively. The number of feature channels and the scaling compared to the image are marked on the feature maps directly.

devised to confront the dilution of semantic context, loss of spatial structure and absence of boundary detail, respectively. These three branches share the same encoder and are derived from different levels of the encoder. In order to make full use of the characteristics and complementarity of these three branches, we propose the Cross Aggregation Module (CAM) and Boundary Refinement Module (BRM) to gradually merge them according to "coarse-fine-finer" strategy, which significantly improves the region accuracy and boundary quality. Considering the multi-scale representation of high-level features and efficient computation simultaneously, we take the advantages of pooling operation and propose a novel Scale-Adaptive Pooling (SAP) structure. Benefiting from the unique framework and these lightweight designs, our method can generate more accurate segmentation results with fewer parameters and faster speed.

### B. Rethink of depth and width

Following the backbone settings of the previous models, we adopt ResNet-50 as the encoder of CTDNet, which is regarded as a standard version CTD-L. Then taking into account the three aspects of accuracy, parameters and speed, we design CTD-S and CTD-M with MobileNet-V2 and ResNet-18 from the perspective of width and depth respectively to explore the potential of efficient networks, which can respond to resource-constrained environments and the demand for fast processing of massive data.

The features extracted by different stages of backbone models are a mixture of low-level, mid-level and high-level features. Larger resolution features contribute less to performance but cost more computations, so we only use features of the last four stages $\{E^2, E^3, E^4, E^5\}$ that have strides of $\{4, 8, 16, 32\}$ with respect to the input image. The comparison of framework details on CTD-S, CTD-M and CTD-L is shown in Fig. 5. Specifically, compared with the conference version, our new proposed CTD-M based on ResNet-18 introduces the

SAP structure in the decoder. which replace the original conv-1 module in the semantic path and improve the accuracy with little impact on the model parameters and speed. As to CTD-S which adopts MobileNet-V2 as backbone, the penultimate stage in the encoder contains two different blocks with the same spatial resolution. Since a narrow network cannot extract complex semantic features like deep networks, we remove the SAP used to extract semantic information in the decoder of CTD-S, and replace the downstream CAM with a common conv-3 module to further reduce the amount of model parameters. In contrast to CTD-S, CTD-L is able to obtain more complex features during encoding with its ResNet-50 backbone, so we can further mine $E^3$ and $E^2$ with SAP and SAM for extra semantic and spatial information respectively. In the following, we introduce each path and module by taking the framework of CTD-M as an example.

### C. Complementary Trilateral Decoder

The proposed Complementary Trilateral Decoder (CTD) framework includes three branches: Semantic Path, Spatial Path and Boundary Path, which are derived from different levels of the shared encoder and are devised to confront the dilution of semantic context, loss of spatial structure and absence of boundary detail, respectively.

**Semantic Path:** Both semantic context and global context are helpful to locate salient objects accurately. However, one of the problems of the U-shape structure is that semantic information of high-level features will be gradually diluted when they are transmitted to lower layers in the top-down path. In addition, the receptive field of the lightweight networks is not large enough to capture global context information.

To solve these issues, we propose the Semantic Path to capture rich semantic context and global context, which can produce initial coarse saliency maps with accurate locations of salient objects. First, we embed a Global Average Pooling (GAP) layer on the tail of the backbone network, which

can provide the maximum receptive field with the strongest global context. The output of global pooling is up-sampled and represented as $E^6$. Then we apply the proposed SAP structure to the latter two stages $E^4$ and $E^5$, which can enlarge the receptive field and capture multi-scale information efficiently. Finally, we design a simple Feature Fusion Module (FFM) to effectively fuse $E^4$, $E^5$ and $E^6$, which forms a partial U-shape structure (see Fig. 3). The Semantic Path can be formulized as:

$$E_g^5 = FFM_1(SAP(E^5), Up(GAP(E^5))), \qquad (1)$$

$$\mathcal{D}_{p1} = FFM_2(SAP(E^4), Up(E_g^5)), \qquad (2)$$

where $SAP$ and $Up$ represent the proposed SAP structure and up-sampling operations, respectively. $\mathcal{D}_{p1}$ represents the output of the Semantic Path.

After that, we introduce the detailed structure of FFM. To be specific, FFM receives two inputs $f_1, f_2$ and we adopt the multiplication operation to fuse these two features. Compared with addition and concatenation, the multiplication operation can avoid redundant information and suppress background noise. The fused features pass through two $3 \times 3$ convolution layers to obtain more robust feature representation. The above process can be described as:

$$FFM(f_1, f_2) = \mathcal{F}_{3\times3}(\mathcal{F}_{3\times3}(f_1 \otimes f_2)), \qquad (3)$$

where $\mathcal{F}_{3\times3}$ represents $3 \times 3$ convolution. Note that each convolution is followed by a batch normalization and a ReLU activation function.

**Spatial Path:** The Semantic Path captures rich semantic context and global context, while the Spatial Path is designed to preserve more spatial details. Spatial information is helpful to supplement structural details and generate complete segmentation results. However, spatial information is seriously lost after multiple down-samplings and cannot be recovered perfectly by integrating the hierarchical features from the encoder. Therefore, we propose the Spatial Path to learn more discriminative feature representation from spatial dimension.

The Spatial Path is drawn from mid-level features $E^3$ with large resolution (1/8 of the input size), which is beneficial to encode affluent spatial details. Specifically, we design a Spatial Attention Module (SAM) to refine features effectively (see Fig. 4). We first apply average channel pooling and maximum channel pooling that are performed pooling operations along the channel axis. These two generated single-channel spatial maps $S_{avg}$ and $S_{max}$ are concatenated. Then we compute a spatial attention map $M_{sa}$ by a $5 \times 5$ convolution and sigmoid function. The spatial attention map $M_{sa}$ can re-weight the features $E^3$ from spatial dimension by element-wise multiplication. Finally, the weighted features $E_{sa}^3$ are fed into a $3 \times 3$ convolution layer to squeeze the number of channels to 64. The Spatial Path can be formulized as:

$$S_{avg} = CP_{avg}(E^3) = avg_{i\in[0,n-1]}(E_i^3), \qquad (4)$$

$$S_{max} = CP_{max}(E^3) = max_{i\in[0,n-1]}(E_i^3), \qquad (5)$$

$$M_{sa} = \sigma(\mathcal{F}_{5\times5}(Concat(S_{avg}, S_{max}))), \qquad (6)$$

$$\mathcal{D}_{p2} = \mathcal{F}_{3\times3}(M_{sa} \otimes E^3) = \mathcal{F}_{3\times3}(E_{sa}^3), \qquad (7)$$

where $CP_{avg}$ and $CP_{max}$ represent average channel pooling and maximum channel pooling operations, respectively. $E_i^3$ and $n$ denote the i-th channel of the feature map $E^3$ and the number of channels. $\mathcal{F}_{5\times5}$ and Concat represent $5 \times 5$ convolution and concatenation. $\sigma$ and $\otimes$ denote sigmoid function and element-wise multiplication. $\mathcal{D}_{p2}$ represents the output of the Spatial Path.

**Boundary Path:** The boundary contour can be regarded as the demarcation between the salient regions and the background. Boundary information is also helpful to segment salient objects accurately. However, we observe that saliency maps produced by many existing SOD methods based on the U-shape structure suffer from coarse boundaries. Therefore, we propose the Boundary Path to improve boundary quality by utilizing boundary contour explicitly. The Boundary Path is drawn from low-level features $E^2$, which preserves more boundary information due to larger resolution (1/4 of the input size). However, it is likely to bring noise and interference, such as the boundaries of non-salient regions. Therefore, we exploit high-level location information as guidance to help enhance salient boundary features and suppress non-salient boundary features with an extra edge supervision (see Fig. 4). Specifically, we first apply a $1 \times 1$ convolution to low-level features $E^2$ for channel compression. Then we up-sample the high-level features $E_g^5$ (see Eq. (1)) to the same size as $E^2$ by bilinear interpolation. Finally, we use the proposed FFM to fuse the local information and location information efficiently. In addition, we add an extra salient edge supervision to supervise the Boundary Path explicitly. The Boundary Path can be formulized as:

$$\mathcal{D}_{p3} = FFM_3(\mathcal{F}_{1\times1}(E^2), Up(E_g^5)), \qquad (8)$$

where $\mathcal{F}_{1\times1}$ represents $1 \times 1$ convolution. Note that each convolution is followed by a batch normalization and a ReLU activation function. $\mathcal{D}_{p3}$ represents the output of the Boundary Path. From that we can know that Boundary Path can generate the boundary contours of salient objects, so our method can also be used for edge detection, which is incidental to our main task. Some visual examples can be found in Fig. 4.

### D. Interaction between Three Branches

In order to make full use of the characteristics and complementarity of these three branches, we propose the Cross Aggregation Module (CAM) and Boundary Refinement Module (BRM) to gradually merge them according to "coarse-fine-finer" strategy, which significantly improves the region accuracy and boundary quality.

**Cross Aggregation Module:** The output of the Semantic Path $\mathcal{D}_{p1}$ contains rich semantic information with global context, which can produce initial coarse saliency maps with accurate locations of salient objects (see Fig. 6(a)). In contrast, the output of the Spatial Path $\mathcal{D}_{p2}$ preserves more spatial details. Both paths are complementary to each other, so we design a novel fusion module CAM to merge these two branches effectively, which can produce relative fine saliency maps with precise structures of salient objects (see Fig. 8(d)). As Fig. 6(a) shows, the two inputs of CAM have different resolutions:
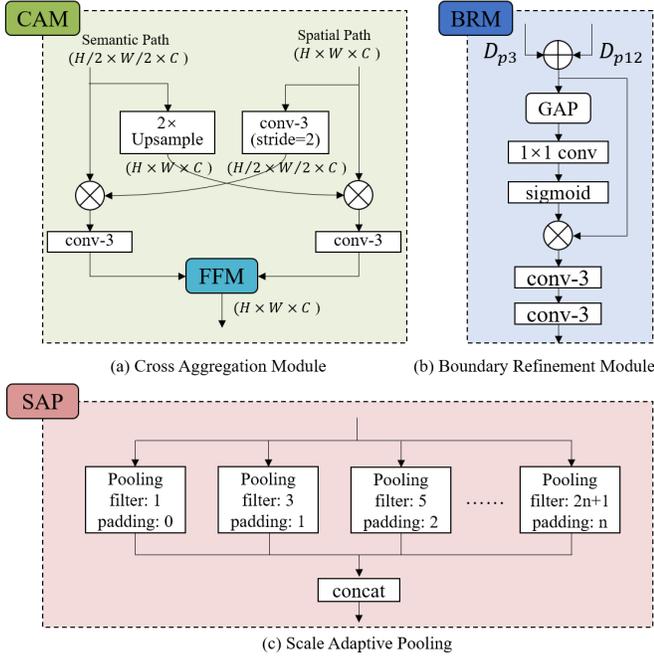
Fig. 6. The detailed structure of CAM, BRM and SAP

$\mathcal{D}_{p1} \in \mathbb{R}^{\frac{H}{2} \times \frac{W}{2} \times C}$ and $\mathcal{D}_{p2} \in \mathbb{R}^{H \times W \times C}$. First, we perform the multi-scale transformation on each input. Specifically, we up-sample $\mathcal{D}_{p1}$ to the same size as $\mathcal{D}_{p2}$ by bilinear interpolation and down-sample $\mathcal{D}_{p2}$ to the same size as $\mathcal{D}_{p1}$ by a $3 \times 3$ convolution with stride 2, obtaining the corresponding features $\mathcal{D}'_{p1} \in \mathbb{R}^{H \times W \times C}$ and $\mathcal{D}'_{p2} \in \mathbb{R}^{\frac{H}{2} \times \frac{W}{2} \times C}$. Second, we perform cross aggregation on each scale by the multiplication operation and then apply a $3 \times 3$ convolution respectively to adapt them, which can capture multi-scale information and promote interaction between two branches. Note that each convolution is followed by a batch normalization and a ReLU activation function. Finally, these two features $C_1 \in \mathbb{R}^{\frac{H}{2} \times \frac{W}{2} \times C}$ and $C_2 \in \mathbb{R}^{H \times W \times C}$ are fed into the proposed FFM to construct a comprehensive and powerful feature representation. The whole process can be described as:

$$\mathcal{D}'_{p1} = Up(\mathcal{D}_{p1}), \mathcal{D}'_{p2} = Down_{3\times3}(\mathcal{D}_{p2}), \quad (9)$$

$$C_1 = \mathcal{F}_{3\times3}(\mathcal{D}_{p1} \otimes \mathcal{D}'_{p2}), C_2 = \mathcal{F}_{3\times3}(\mathcal{D}_{p2} \otimes \mathcal{D}'_{p1}), \quad (10)$$

$$\mathcal{D}_{p12} = FFM_4(Up(C_1), C_2), \quad (11)$$

where $Down_{3\times3}$ denotes down-sampling operation using $3 \times 3$ convolution with stride 2. $\mathcal{D}_{p12}$ represents the final output of the CAM.

**Boundary Refinement Module:** Although we obtain relatively fine saliency maps after merging the Semantic Path $\mathcal{D}_{p1}$ and the Spatial Path $\mathcal{D}_{p2}$, we can leverage the salient boundary information provided by the Boundary Path to further refine boundary. Therefore, we propose a fusion module BRM to merge $\mathcal{D}_{p12}$ and the Boundary Path $\mathcal{D}_{p3}$, which can generate final finer saliency maps with clear boundaries of salient objects (see Fig. 8(e)).

As Fig. 6(b) shows, we first concatenate the $\mathcal{D}_{p12}$ and $\mathcal{D}_{p3}$. Then we pool the fused features $B_f$ to generate a feature

vector and compute an attention vector to guide the feature learning by a $1 \times 1$ convolution and sigmoid function. This weight vector can re-weight the $B_f$ for feature selection and refinement by multiplication operation. Finally, the refined features $B_r$ are combined with $B_f$ and then pass through two $3 \times 3$ convolution layers to further enhance feature representation. Note that each $3 \times 3$ convolution is followed by a batch normalization and a ReLU activation function. The above process can be described as:

$$B_f = Up(\mathcal{D}_{p12}) + \mathcal{D}_{p3}, \quad (12)$$

$$B_r = B_f \otimes \sigma(\mathcal{F}_{1\times1}(GAP(B_f))), \quad (13)$$

$$\mathcal{D}_{p123} = \mathcal{F}_{3\times3}(\mathcal{F}_{3\times3}(B_r + B_f)), \quad (14)$$

$\mathcal{D}_{p123}$ represents the final output of the BRM.

### E. Scale-Adaptive Pooling

The receptive filed is of great significance for CNN, because the size of receptive filed can roughly indicate how much we can use context information. However, the empirical receptive fields of CNN are much smaller than the theoretical ones especially for high-level features. On the other hand, the multi-scale representation can help perceiving multi-scale objects, which is also indispensable for SOD. Some approaches have been proposed to enlarge the receptive field and capture multi-scale features, such as Pyramid Pooling Module (PPM) [37] and Atrous Spatial Pyramid Pooling (ASPP) [38], but these structures are computation demanding and memory consuming.

To this end, we propose a novel structure named Scale-Adaptive Pooling (SAP), as shown in Fig. 6(c). The scale-adaptive includes two aspects: receptive field scale-adaptive and spatial scale-adaptive. To be specific, we use $N(N = 4)$ parallel pooling operations with stride $s = 1$, padding $p = n$ and filter $f = 2n + 1$, where $n = \{0, 1, , N - 1\}$. On the one hand, we use pooling of different kernel sizes to obtain multiple receptive fields, which can capture multi-scale context information and enlarge the receptive field to a certain extent. On the other hand, these pooling operations take a fixed stride of 1 and different paddings to keep the spatial size the same as input. Finally, the $N$ parallel branches are concatenated as the final output of SAP module.

Compared with PPM and ASPP, our proposed SAP has the following advantages: 1) Unlike ASPP that uses dilated convolutions with different dilation rates, SAP only uses pooling operations without increasing learnable parameters; 2) Different from PPM that requires frequent down-sampling and up-sampling, SAP does not change the spatial size of feature maps and reduces spatial information loss to a certain extent; 3) SAP can capture multi-scale context information from receptive fields of different sizes, which is more efficient and much faster.

### F. Loss Function

In this paper, we adopt three loss functions: BCE loss [39], IoU loss [40] and $L_1$ loss [41]. BCE loss calculates the error
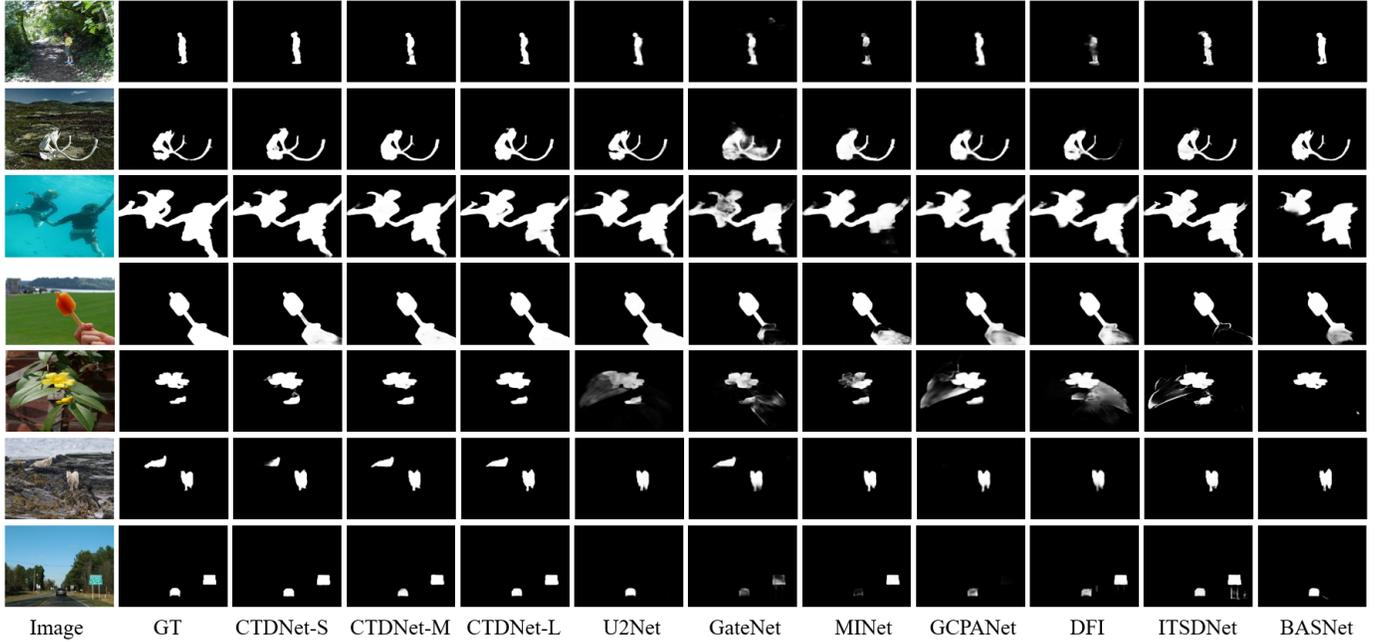
Fig. 7. Qualitative comparison of our model with existing state-of-the-art SOD models in some challenging scenarios

for each pixel between the prediction mask and the ground truth, which is formulated as:

$$\ell_{bce}(P, G) = -\sum_{i=1}^{H}\sum_{j=1}^{W}[G(i,j)log(P(i,j)) \\ +(1-G(i,j))log(1-P(i,j))], \quad (15)$$

where $P(i,j)$ and $G(i,j)$ represent the pixel of prediction mask (P) and the ground truth (G) at location $(i,j)$ in an image. H and W are the height and width of the image, respectively. IoU loss is used to measure the similarity of structure instead of focusing on single pixel. We adopt the following form:

$$\ell_{iou}(P, G) = 1 - \frac{\sum_{i=1}^{H}\sum_{j=1}^{W}G(i,j)P(i,j)}{\sum_{i=1}^{H}\sum_{j=1}^{W}[G(i,j)+P(i,j)-G(i,j)p(i,j)]}, \quad (16)$$

$L_1$ loss is used to measure the minimum absolute error for each pixel between the prediction mask and the ground truth, which is formulated as:

$$\ell_1(P, G) = \sum_{i=1}^{H}\sum_{j=1}^{W}|G(i,j)-P(i,j)|, \quad (17)$$

As described above, our model is deeply supervised with six outputs. All outputs pass through a $3 \times 3$ convolution and sigmoid function to convert the feature maps to the corresponding single-channel prediction masks. For $\mathcal{D}_{p123}$, $\mathcal{D}_{p12}$, $\mathcal{D}_{p1}$, $E_g^5$ and $E^6$, we use three loss functions together to supervise these five saliency maps (see Eq. (17)), while for $\mathcal{D}_{p3}$, we use BCE loss and $L_1$ loss to supervise the boundary prediction mask ($P_b$). Note that the ground truth of salient boundary ($G_b$) can be easily obtained from the ground truth of salient objects.

$$\ell_s(P, G) = \ell_{iou}(P, G) + \beta\ell_{bce}(P, G) + \gamma\ell_1(P, G), \quad (18)$$

$$\ell_b(P_b, G_b) = \frac{1}{2}(\ell_{bce}(P_b, G_b) + l_1(P_b, G_b)), \quad (19)$$

where $\beta$ and $\gamma$ are hyperparameters to balance the weight between the three loss functions, so that the network can achieve better performance. In our paper, the parameter $\beta$ is set to 0.6. The total loss function is denoted as follows:

$$\mathcal{L}(P, P_b, G, G_b) = l_b(P_b, G_b) + \sum_{k=1}^{5}\alpha_k l_s(P^k, G), \quad (20)$$

where $\alpha_k$ denotes the weight of the $k-th$ loss term.

## IV. EXPERIMENTS

### A. Datasets

We conduct experiments on five standard benchmark datasets: ECSSD [42], PASCAL-S [43], DUTS [44], HKU-IS [45] and DUT-OMRON [46], and the detailed introduction is provided as follows: ECSSD contains 1,000 natural images with many semantically meaningful and complex structures in their ground-truth segmentation. PASCAL-S consists of 850 natural images that are carefully selected from the PASCAL VOC dataset. DUTS is currently the largest SOD dataset, which contains 10,553 images for training (DUTS-TR) and 5,019 images for testing (DUTS-TE). HKU-IS includes 4,447 challenging images and most of them have multiple disconnected salient objects, overlapping image boundaries or low color contrast. DUT-OMRON has 5,168 high quality images, which have one or more salient objects and relatively cluttered background.

### B. Evaluation Metrics

To quantitatively evaluate the performance, we adopt three evaluation metrics: Mean Absolute Error (MAE), F-measure and E-measure [47]. MAE measures the pixel-wise average

TABLE I
QUANTITATIVE COMPARISONS WITH STATE-OF-THE-ART SOD MODELS ON FIVE BENCHMARKS IN TERMS OF PARAMETERS, SPEED, $mF_\beta$, MAE, $E_m$. THE SPEED DATA MARKED WITH † ARE UNIFORMLY RE-MEASURED ON TITAN XP GPU IN THE SAME ENVIRONMENT. THE BEST TWO RESULTS ARE SHOWN IN RED AND GREEN, RESPECTIVELY.

| Method | Params (M) | Speed (FPS) | ECSSD $mF_\beta$ | MAE | $E_m$ | PASCAL-S $mF_\beta$ | MAE | $E_m$ | DUTS-TE $mF_\beta$ | MAE | $E_m$ | HKU-IS $mF_\beta$ | MAE | $E_m$ | DUT-OMRON $mF_\beta$ | MAE | $E_m$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **VGG/ResNet-based Models** | | | | | | | | | | | | | | | | | |
| **C2SNet[18]** | 158.86 | 30 | .864 | .055 | .914 | .758 | .080 | .839 | .716 | .063 | .846 | .851 | .048 | .927 | .683 | .072 | .829 |
| **RAS[18]** | 21.23 | 35 | .889 | .056 | .914 | .777 | .101 | .829 | .751 | .059 | .861 | .871 | .045 | .929 | .713 | .062 | .846 |
| **PiCANet[18]** | 38.32 | 7 | .885 | .046 | .910 | .789 | .077 | .828 | .749 | .054 | .852 | .870 | .042 | .934 | .710 | .068 | .834 |
| **BMPM[18]** | 75.07 | 22 | .868 | .045 | .914 | .758 | .073 | .836 | .745 | .049 | .860 | .871 | .039 | .937 | .692 | .064 | .837 |
| **PAGE[19]** | 47.40 | 25 | .906 | .042 | .920 | .806 | .075 | .841 | .777 | .052 | .869 | .882 | .037 | .940 | .736 | .062 | .853 |
| **AFNet[19]** | 35.99 | 26 | .908 | .042 | .918 | .820 | .070 | .850 | .793 | .046 | .879 | .888 | .036 | .942 | .738 | .057 | .853 |
| **BANet[19]** | 56.02 | 13 | .923 | .035 | .928 | .823 | .069 | .852 | .815 | .040 | .892 | .900 | .032 | .950 | .746 | .059 | .860 |
| **EGNet[19]** | 111.78 | 8 | .920 | .037 | .927 | .817 | .073 | .848 | .815 | .039 | .891 | .901 | .031 | .950 | .755 | .053 | .867 |
| **SCRN[19]** | 25.32 | 32 | .918 | .038 | .926 | .826 | .064 | .857 | .809 | .040 | .888 | .896 | .034 | .949 | .746 | .056 | .863 |
| **PoolNet[19]** | 68.16 | 18 | .915 | .039 | .924 | .815 | .074 | .848 | .809 | .040 | .889 | .899 | .032 | .949 | .747 | .056 | .863 |
| **CPD[19]** | 47.97 | 62 | .917 | .037 | .925 | .820 | .070 | .849 | .805 | .043 | .887 | .891 | .034 | .944 | .747 | .056 | .866 |
| **BASNet[19]** | 87.03 | 25 | .880 | .037 | .921 | .771 | .075 | .846 | .791 | .048 | .884 | .895 | .032 | .946 | .756 | .056 | .869 |
| **GateNet[20]** | - | - | .916 | .040 | .924 | .819 | .067 | .851 | .807 | .040 | .889 | .899 | .033 | .949 | .746 | .055 | .862 |
| **U2Net[20]** | 46.21 | 30 | .892 | .033 | .924 | .770 | .073 | .842 | .792 | .045 | .886 | .896 | .031 | .948 | .761 | .054 | .871 |
| **DFI[20]** | 29.57 | 57 | .920 | .038 | .924 | .830 | .064 | .855 | .814 | .039 | .892 | .901 | .031 | .951 | .752 | .055 | .865 |
| **GCPANet[20]** | 67.05 | 50 | .919 | .035 | .920 | .827 | .061 | .847 | .817 | .038 | .891 | .898 | .031 | .949 | .748 | .056 | .860 |
| **ITSDNet[20]** | 26.55 | 43 | .895 | .035 | .927 | .785 | .071 | .850 | .804 | .041 | .895 | .899 | .031 | .952 | .756 | .061 | .863 |
| **MINet[20]** | 162.38 | 31 | .924 | .033 | .927 | .829 | .063 | .851 | .828 | .037 | .898 | .909 | .029 | .953 | .755 | .055 | .865 |
| **CII[21]** | 24.48 | 50 | .929 | .033 | .926 | .841 | .061 | .857 | .834 | .037 | .900 | .915 | .029 | .953 | .768 | .054 | .872 |
| **EDN[22]** | 42.85 | 51.7 | .932 | .032 | .929 | .847 | .061 | .864 | .851 | .035 | .908 | .919 | .026 | .956 | .783 | .048 | .876 |
| **CTD-L (Ours)** | 26.48 | 84 | .931 | .032 | .925 | .845 | .059 | .860 | .863 | .032 | .914 | .922 | .025 | .957 | .783 | .049 | .878 |
| **Lightweight Models** | | | | | | | | | | | | | | | | | |
| **SAMNet[21]** | 1.33 | 28† | .891 | .050 | .911 | .778 | .090 | .823 | .745 | .058 | .849 | .871 | .045 | .934 | .717 | .065 | .840 |
| **CII (ResNet-18)[21]** | 11.89 | 100† | .915 | .039 | .921 | .820 | .067 | .848 | .814 | .043 | .890 | .906 | .032 | .949 | .745 | .058 | .860 |
| **EDN-Lite[22]** | 1.8 | 44† | .916 | .042 | .919 | .820 | .072 | .853 | .809 | .045 | .888 | .901 | .034 | .945 | .748 | .057 | .857 |
| **CTD-S (Ours)** | 1.7 | 125 | .915 | .038 | .924 | .826 | .066 | .851 | .816 | .041 | .890 | .907 | .029 | .951 | .758 | .056 | .863 |
| **CTD-M (Ours)** | 12.6 | 158 | .923 | .035 | .922 | .831 | .065 | .858 | .835 | .038 | .900 | .915 | .028 | .954 | .767 | .054 | .869 |

absolute difference between the prediction mask and ground truth, which is formulated as:

$$MAE = \frac{1}{W \times H} \sum_{i=1}^{H} \sum_{j=1}^{W} |P(i,j) - G(i,j)|, \quad (21)$$

where P and G represent the prediction mask and the corresponding ground truth, respectively. H, W are the height and width of the image. The smaller MAE indicates better performance. Another metric F-measure ($F_\beta$) takes both precision and recall into account, which is defined as:

$$F_\beta = \frac{(1 + \beta^2) \cdot Precision \cdot Recall}{\beta^2 \cdot Precision + Recall} \quad (22)$$

where $\beta^2$ is set to 0.3 to emphasize the precision over recall. Larger F-measure indicates better performance. We choose the max F-measure ($mF_\beta$) in our paper. In addition, E-measure ($E_m$) combines local pixel values with the image-level mean value to jointly evaluate the similarity between the prediction mask and the ground truth.

### C. Implementation Details

We implement our proposed method by PyTorch and conduct experiments on one NVIDIA 1080Ti GPU. We use MobileNet-V2, ResNet-18 and ResNet-50 pre-trained on ImageNet as backbone networks, respectively. We choose DUTS-TR as training dataset as used in [18], [48] and evaluate our model on other datasets. In the training period, all training images are resized to 352×352 with random cropping and random horizontal flipping to feed into the proposed model. We use stochastic gradient descent (SGD) optimizer with momentum of 0.9 and weight decay of 5e-4 to train our model. The batch size is set to 32 and maximum epoch is set to 48. We adopt the warm-up and linear decay learning rate strategy with the maximum learning rate 5e-3 for pre-trained backbone and 5e-2 for the rest of network. During the inference period, each image is simply resized to 352×352 and then fed into our model to predict saliency map without any post-processing (e.g., CRF [49]).

### D. Comparison results

To prove the effectiveness of our method, we compare with 18 state-of-the-art SOD models, including C2SNet [22], RAS

TABLE II

THE ABLATION STUDY OF OUR PROPOSED COMPONENTS. THE BACKBONE NETWORK TAKES RESNET-18 AS AN EXAMPLE. BY ADDING EACH MODULE GRADUALLY, OUR MODEL ACHIEVES THE BEST PERFORMANCE.

| Base | Semantic Path | | | Spatial Path | | Boundary Path | | Params | Speed | DUTS-TE | | | DUT-OMRON | | |
|------|------|-----|--------|------|------|------|------|--------|-------|---------|-----|-------|---------|-----|-------|
| | FFM | SAP | Global | SAM | CAM | FFM | BRM | (M) | (FPS) | $mF_\beta$ | MAE | $E_m$ | $mF_\beta$ | MAE | $E_m$ |
| ✓ | | | | | | | | 11.352 | 270 | .797 | .046 | .880 | .724 | .064 | .843 |
| ✓ | ✓ | | | | | | | 11.389 | 260 | .803 | .046 | .882 | .733 | .064 | .847 |
| ✓ | ✓ | ✓ | | | | | | 12.127 | 240 | .806 | .044 | .882 | .742 | .06 | .855 |
| ✓ | ✓ | ✓ | ✓ | | | | | 12.234 | 220 | .811 | .043 | .886 | .744 | .06 | .855 |
| ✓ | ✓ | ✓ | ✓ | ✓ | | | | 12.300 | 210 | .817 | .042 | .89 | .751 | .059 | .857 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | 12.448 | 190 | .821 | .041 | .891 | .753 | .057 | .863 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | 12.522 | 180 | .826 | .040 | .895 | .754 | .056 | .863 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 12.612 | 160 | .835 | .038 | .900 | .767 | .054 | .869 |



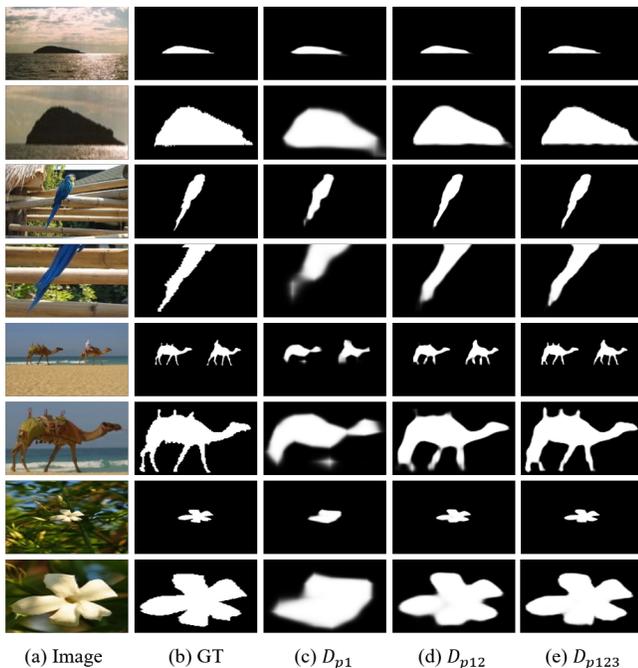|          |        |            |             |              |
|----------|--------|------------|-------------|--------------|
| (a) Image | (b) GT | (c) $D_{p1}$ | (d) $D_{p12}$ | (e) $D_{p123}$ |

Fig. 8. The saliency maps from different locations of our network. Each example contains two rows and the second row of each example denotes the zoom-in view of salient object. The results conform to "coarse-fine-finer" predictions along with the gradual combination of these three branches, which demonstrates the complementarity of these three branches.

[29], PiCANet [18], BMPM [19], BANet [7], EGNet [24], SCRN [50], PoolNet [30], PAGE [51], AFNet [25], CPD [31], BASNet [23], GateNet [52], DFI [53], ITSDNet [10], GCPANet [9], MINet [20], U2Net [54], CII [13], EDN [11] and SAMNet [32]. For a fair comparison, we use saliency maps released by the authors and evaluate them with the same Matlab code.

**Qualitative Comparison.** To intuitively show the advantages of our three models, we provide some visual examples of various SOD models, as shown in Fig. 7. We can observe that our models CTD-S, CTD-M and CTD-L can generate more complete and more accurate segmentation results than other counterparts. It can handle various challenging scenarios, such as multiple salient objects (row 3, 6 and 7), fine structure (row 2, 4 and 5), cluttered backgrounds (row 1 and 5) and small objects (row 6 and 7). In addition, we do not use any post-processing to obtain these results. Therefore, our model shows its effectiveness and robustness in processing complicated images.

**Quantitative Comparison.** Tab. I shows the quantitative results on five popular datasets in terms of $mF_\beta$, MAE and $E_m$. To facilitate the practical application in different environments, we adopt MobileNet-V2, ResNet-18 and ResNet-50 as backbones respectively and name our model CTD-S, CTD-M and CTD-L accordingly. Compared with VGG/ResNet-based models, CTD-L obtains the best performance under almost all evaluation metrics on five benchmarks. More importantly, our models CTD-S and CTD-M outperform other lightweight models and achieve comparable or even better performance than VGG/ResNet-based models. In addition, we also list the parameters and speed of each method to measure efficiency. The MobileNet-based CTD-S that can run at a 125 FPS speed only has 1.7M parameters, which is more than five times compressed compared to the most existing model parameters; the ResNet-based CTD-M only has 12.6M parameters and runs at a speed of 158 FPS on one GTX 1080Ti GPU for 352×352 input images, which surpasses the existing approaches by a large margin. Moreover, CTD-L runs at a 84 FPS speed with 26.5M parameters, which is much smaller and faster than the existing VGG/ResNet-based SOD methods. It should be noted that although CTD-S has fewer parameters than CTD-M, it is not as fast as CTD-M due to the depthwise separable convolutions in its backbone [55]. In conclusion, our model achieves a favorable trade-off between speed and accuracy, which clearly demonstrates its superiority and efficiency.

### E. Ablation Study

Firstly, we investigate the complementarity of these three branches. Secondly, we verify the effectiveness of our proposed components. Lastly, we validate the improvements of

TABLE III
THE COMPLEMENTARITY OF THESE THREE BRANCHES. $\mathcal{D}_{p1}, \mathcal{D}_{p12}$ AND $\mathcal{D}_{p123}$ DENOTE THE SEMANTIC PATH, THE COMBINATION OF BOTH SEMANTIC PATH AND SPATIAL PATH, THE COMBINATION OF THESE THREE BRANCHES, RESPECTIVELY

| Merge | DUTS-TE | | | DUT-OMRON | | |
|---|---|---|---|---|---|---|
| | $mF_\beta$ | MAE | $E_m$ | $mF_\beta$ | MAE | $E_m$ |
| $\mathcal{D}_{p1}$ | .766 | .049 | .872 | .712 | .063 | .848 |
| $\mathcal{D}_{p12}$ | .819 | .040 | .894 | .754 | .056 | .866 |
| $\mathcal{D}_{p123}$ | .835 | .038 | .900 | .767 | .054 | .869 |

TABLE IV
THE IMPROVEMENTS OF ACCURACY WITH SAP MODULE COMPARED TO THE CONFERENCE VERSION

| Model | DUTS-TE | | | DUT-OMRON | | |
|---|---|---|---|---|---|---|
| | $mF_\beta$ | MAE | $E_m$ | $mF_\beta$ | MAE | $E_m$ |
| CTDNet(conference) | .853 | .034 | .909 | .779 | .052 | .875 |
| CTDNet with SAP | .856 | .033 | .912 | .781 | .050 | .876 |

accuracy with SAP module compared to the conference version. All experiments are conducted on DUTS-TE and DUT-OMRON datasets.

**The complementarity of these three branches.** To demonstrate the complementarity and necessity of these three branches, we conduct experiments both qualitatively and quantitatively. As shown in Tab. III, when merging the Semantic Path $\mathcal{D}_{p1}$ and Spatial Path $\mathcal{D}_{p2}$, the performance can be greatly improved. Moreover, the performance can be further boosted by merging $\mathcal{D}_{p12}$ and $\mathcal{D}_{p3}$, which benefits from the salient boundary information provided by the Boundary Path $\mathcal{D}_{p3}$. In addition, we visualize some examples in Fig. 8. Each example contains two rows and the second row of each example denotes the zoom-in view of salient object. As we can see, the produced saliency maps conform to "coarse-fine-finer" predictions along with the gradual combination of these three branches. Column 3 represents initial coarse saliency maps produced by $\mathcal{D}_{p1}$ with accurate locations of salient objects. Column 4 represents relatively fine saliency maps produced by $\mathcal{D}_{p12}$ with precise structures of salient objects. Column 5 represents final finer saliency maps produced by $\mathcal{D}_{p123}$ with clear boundaries of salient objects. Obviously, experimental results verify the complementarity and necessity of these three branches.

**The effectiveness of our proposed components.** To demonstrate the effectiveness of our proposed components, we conduct ablation experiments by gradually adding them. First, we replace all the proposed fusion modules with simple addition operation followed by the $3 \times 3$ convolution to construct a baseline network, which still maintains three branches in the decoder. Second, we gradually add the FFM, SAP and global context for the Semantic Path and Boundary Path. Then we add the SAM in the Spatial Path. Finally, we use the proposed fusion modules CAM and BRM to merge these three branches. As shown in Tab. II, our method achieves

the best performance when all modules are contained, which demonstrates the effectiveness and necessity of each module.

**The improvements of accuracy with SAP module.** To explore the extent to which the newly proposed module SAP improves the accuracy compared to the conference version, we conduct ablation experiments by only adding SAP modules into the CTDNet of the conference version. As shown in Tab. IV, the SAP module effectively improves the accuracy of the model and helps the model obtain the best performance against state-of-the-art methods. Note that although the newly added module SAP has slightly slowed down the speed compared with conference version, the overall efficiency of our new proposed model is still excellent compared with other counterparts and the accuracy has been effectively improved.

## V. CONCLUSION

In this paper, we first reveal that existing salient object detection methods adopt deeper and wider networks to achieve better performance, resulting in imbalance in accuracy, parameters and speed. To this end, we design a lightweight framework while maintaining satisfying competitive accuracy. Specifically, after analyzing the drawbacks of the U-shape structure, we propose a novel trilateral decoder framework by decoupling the U-shape structure into three complementary branches, which are devised to confront the dilution of semantic context, loss of spatial structure and absence of boundary detail, respectively. Along with the fusion of three branches, the coarse segmentation results are gradually refined in structure details and boundary quality. Taking the advantages of pooling operation, we propose the Scale-Adaptive Pooling to obtain multi-scale receptive filed without additional learnable parameters. On the premise of inheriting this framework, we explore the maximum potential of shallow and narrow network and provide three versions: CTD-S (1.7M, 125FPS), CTD-M(12.6M, 158FPS) and CTD-L (26.5M, 84FPS) to facilitate the practical application in different environments. Experiments demonstrate that our proposed method performs better than state-of-the-art methods on five benchmarks, which achieves a favorable trade-off between efficiency and performance.

## VI. ACKNOWLEDGEMENT

## REFERENCES

[1] Z. Zhao, C. Xia, C. Xie, and J. Li, "Complementary trilateral decoder for fast and accurate salient object detection," in *Proceedings of the 29th acm international conference on multimedia*, 2021, pp. 4967–4975.
[2] A. Borji, M.-M. Cheng, Q. Hou, H. Jiang, and J. Li, "Salient object detection: A survey," *Computational visual media*, pp. 1–34, 2019.
[3] 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*, 2021.
[4] Y. Gao, M. Shi, D. Tao, and C. Xu, "Database saliency for fast image retrieval," *IEEE Transactions on Multimedia*, vol. 17, no. 3, pp. 359–369, 2015.
[5] S. Hong, T. You, S. Kwak, and B. Han, "Online tracking by learning discriminative saliency map with convolutional neural network," in *International conference on machine learning*, 2015, pp. 597–606.

[6] R. M. Haralick and L. G. Shapiro, "Image segmentation techniques," *Computer vision, graphics, and image processing*, vol. 29, no. 1, pp. 100–132, 1985.

[7] J. Su, J. Li, Y. Zhang, C. Xia, and Y. Tian, "Selectivity or invariance: Boundary-aware salient object detection," in *ICCV*, 2019.

[8] Z. Wu, L. Su, and Q. Huang, "Decomposition and completion network for salient object detection," *IEEE Transactions on Image Processing*, vol. 30, pp. 6226–6239, 2021.

[9] Z. Chen, Q. Xu, R. Cong, and Q. Huang, "Global context-aware progressive aggregation network for salient object detection," *arXiv preprint arXiv:2003.00651*, 2020.

[10] H. Zhou, X. Xie, J.-H. Lai, Z. Chen, and L. Yang, "Interactive two-stream decoder for accurate and fast saliency detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9141–9150.

[11] 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.

[12] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4510–4520.

[13] J.-J. Liu, Z.-A. Liu, P. Peng, and M.-M. Cheng, "Rethinking the u-shape structure for salient object detection," *IEEE Transactions on Image Processing*, vol. 30, pp. 9030–9042, 2021.

[14] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition. 2016 ieee conf comput vispattern recognit. 2016: 770-778 https://doi. org/10.1109." *CVPR*, 2016.

[15] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

[16] C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang, "Bisenet: Bilateral segmentation network for real-time semantic segmentation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 325–341.

[17] C. Xia, J. Li, X. Chen, A. Zheng, and Y. Zhang, "What is and what is not a salient object? learning salient object detector by ensembling linear exemplar regressors," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4142–4150.

[18] N. Liu, J. Han, and M.-H. Yang, "Picanet: Learning pixel-wise contextual attention for saliency detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3089–3098.

[19] L. Zhang, J. Dai, H. Lu, Y. He, and G. Wang, "A bi-directional message passing model for salient object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1741–1750.

[20] Y. Pang, X. Zhao, L. Zhang, and H. Lu, "Multi-scale interactive network for salient object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9413–9422.

[21] S. Yang, W. Lin, G. Lin, Q. Jiang, and Z. Liu, "Progressive self-guided loss for salient object detection," *IEEE Transactions on Image Processing*, vol. 30, pp. 8426–8438, 2021.

[22] X. Li, F. Yang, H. Cheng, W. Liu, and D. Shen, "Contour knowledge transfer for salient object detection," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 355–370.

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

[24] J.-X. Zhao, J.-J. Liu, D.-P. Fan, Y. Cao, J. Yang, and M.-M. Cheng, "Eg-net: Edge guidance network for salient object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 8779–8788.

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

[26] J. Li, J. Su, C. Xia, M. Ma, and Y. Tian, "Salient object detection with purificatory mechanism and structural similarity loss," *IEEE Transactions on Image Processing*, vol. 30, pp. 6855–6868, 2021.

[27] J. Su, C. Xia, and J. Li, "Exploring driving-aware salient object detection via knowledge transfer," in *2021 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2021, pp. 1–6.

[28] J. Li, J. Su, C. Xia, and Y. Tian, "Distortion-adaptive salient object detection in 360° omnidirectional images," *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 1, pp. 38–48, 2019.

[29] S. Chen, X. Tan, B. Wang, and X. Hu, "Reverse attention for salient object detection," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 234–250.

[30] 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 Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3917–3926.

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

[32] Y. Liu, X.-Y. Zhang, J.-W. Bian, L. Zhang, and M.-M. Cheng, "Samnet: Stereoscopically attentive multi-scale network for lightweight salient object detection," *IEEE Transactions on Image Processing*, vol. 30, pp. 3804–3814, 2021.

[33] E. Shelhamer, J. Long, and T. Darrell, "Fully convolutional networks for semantic segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 4, pp. 640–651, 2017.

[34] W. Wang, J. Shen, M.-M. Cheng, and L. Shao, "An iterative and cooperative top-down and bottom-up inference network for salient object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5968–5977.

[35] J. Zhao, Y. Zhao, J. Li, and X. Chen, "Is depth really necessary for salient object detection?" in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 1745–1754.

[36] M. Ma, C. Xia, and J. Li, "Pyramidal feature shrinking for salient object detection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 3, 2021, pp. 2311–2318.

[37] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2881–2890.

[38] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2017.

[39] P.-T. De Boer, D. P. Kroese, S. Mannor, and R. Y. Rubinstein, "A tutorial on the cross-entropy method," *Annals of operations research*, vol. 134, no. 1, pp. 19–67, 2005.

[40] G. Máttyus, W. Luo, and R. Urtasun, "Deeproadmapper: Extracting road topology from aerial images," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3438–3446.

[41] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.

[42] Q. Yan, L. Xu, J. Shi, and J. Jia, "Hierarchical saliency detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 1155–1162.

[43] Y. Li, X. Hou, C. Koch, J. M. Rehg, and A. L. Yuille, "The secrets of salient object segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 280–287.

[44] L. Wang, H. Lu, Y. Wang, M. Feng, D. Wang, B. Yin, and X. Ruan, "Learning to detect salient objects with image-level supervision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 136–145.

[45] G. Li and Y. Yu, "Visual saliency based on multiscale deep features," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 5455–5463.

[46] C. Yang, L. Zhang, H. Lu, X. Ruan, and M.-H. Yang, "Saliency detection via graph-based manifold ranking," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 3166–3173.

[47] D.-P. Fan, C. Gong, Y. Cao, B. Ren, M.-M. Cheng, and A. Borji, "Enhanced-alignment measure for binary foreground map evaluation," *arXiv preprint arXiv:1805.10421*, 2018.

[48] T. Wang, L. Zhang, S. Wang, H. Lu, G. Yang, X. Ruan, and A. Borji, "Detect globally, refine locally: A novel approach to saliency detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3127–3135.

[49] P. Krähenbühl and V. Koltun, "Efficient inference in fully connected crfs with gaussian edge potentials," in *Advances in neural information processing systems*, 2011, pp. 109–117.

[50] Z. Wu, L. Su, and Q. Huang, "Stacked cross refinement network for edge-aware salient object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 7264–7273.

[51] W. Wang, S. Zhao, J. Shen, S. C. Hoi, and A. Borji, "Salient object detection with pyramid attention and salient edges," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1448–1457.

[52] X. Zhao, Y. Pang, L. Zhang, H. Lu, and L. Zhang, "Suppress and balance: A simple gated network for salient object detection," *arXiv preprint arXiv:2007.08074*, 2020.

[53] J.-J. Liu, Q. Hou, and M.-M. Cheng, "Dynamic feature integration for simultaneous detection of salient object, edge and skeleton," *arXiv preprint arXiv:2004.08595*, 2020.

[54] X. Qin, Z. Zhang, C. Huang, M. Dehghan, O. R. Zaiane, and M. Jager-sand, "U2-net: Going deeper with nested u-structure for salient object detection," *Pattern Recognition*, vol. 106, p. 107404, 2020.

[55] M. Orsic, I. Kreso, P. Bevandic, and S. Segvic, "In defense of pre-trained imagenet architectures for real-time semantic segmentation of road-driving images," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 607–12 616.