# CAGNet: Content-Aware Guidance for Salient Object Detection

Sina Mohammadi[a,1,*], Mehrdad Noori[a,1,*], Ali Bahri[a], Sina Ghofrani Majelan[a], Mohammad Havaei[b]

[a]*School of Electrical Engineering, Iran University of science and Technology, Tehran, Iran*
[b] *Imagia Inc., Montreal, Canada*

## Abstract

Beneficial from Fully Convolutional Neural Networks (FCNs), saliency detection methods have achieved promising results. However, it is still challenging to learn effective features for detecting salient objects in complicated scenarios, in which i) non-salient regions may have "salient-like" appearance; ii) the salient objects may have different-looking regions. To handle these complex scenarios, we propose a Feature Guide Network which exploits the nature of low-level and high-level features to i) make foreground and background regions more distinct and suppress the non-salient regions which have "salient-like" appearance; ii) assign foreground label to different-looking salient regions. Furthermore, we utilize a Multi-scale Feature Extraction Module (MFEM) for each level of abstraction to obtain multi-scale contextual information. Finally, we design a loss function which outperforms the widely used Cross-entropy loss. By adopting four different pre-trained models as the backbone, we prove that our method is very general with respect to the choice of the backbone model. Experiments on six challenging datasets demonstrate that our method achieves the state-of-the-art performance in terms of different evaluation metrics. Additionally, our approach contains fewer parameters than the existing ones, does not need any post-processing, and runs fast at a real-time speed of 28 FPS when processing a $480 \times 480$ image.

*Keywords:* Saliency detection, Fully convolutional neural networks, Attention guidance

## 1. Introduction

Salient object detection aims at localizing the most interesting and prominent parts of an image. Moreover, it is an effective pre-processing step for numerous computer vision tasks such as image classification [1], image segmentation [2, 3, 4], video segmentation [5], image editing [6, 7] and object tracking [8].

Traditional approaches are mostly based on low-level cues and hand-crafted features. For example, the method proposed in [9] uses color feature to detect salient objects. Some other methods use center prior to

---

*Corresponding authors.
Email addresses:* `sina.mhm93@gmail.com` (Sina Mohammadi), `me.noori.1994@gmail.com` (Mehrdad Noori)
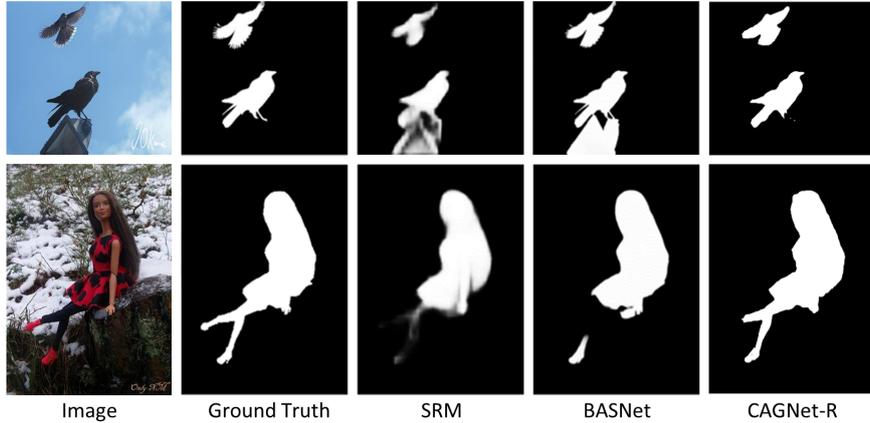[1]Both the authors contributed equally.

Figure 1: Examples of complicated scenarios in salient object detection. In the first row, the triangular object has "salient-like" appearance. In the second row, the appearance of the feet of the doll is different from the rest of the doll. While both scenarios have caused confusion for two recent methods (BASNet [15] and SRM [16]), our method (denoted as CAGNet-R) is capable of handling these complicated scenarios and generating more accurate saliency prediction.

improve the performance of salient object detection [10, 11]. Because of the lack of semantic information, these methods have limited ability to detect the whole structure of salient objects in complex scenes. In recent years, the methods based on the Fully Convolutional Neural Networks (FCNs), such as [12, 13, 14], have been widely used for saliency detection owing to their high capacity of modeling high-level semantics. Even though these methods have achieved promising results, there are still some challenges due to the complicated scenarios of some images. The learned features by these methods usually lack the ability to i) suppress the non-salient regions that have "salient-like" appearance as depicted in the first row of Figure 1, ii) detect salient objects that have different-looking regions as depicted in the second row of Figure 1.

To address the above-mentioned challenges, we propose the Guide Module which takes advantage of the nature of the high-level and low-level features. By adopting this module, high-level features, which do not contain the fine spatial details of low-level features, can exploit the nature of low-level features as a guidance to make foreground and background regions more distinct, and thus it can suppress the non-salient regions that have "salient-like" appearance. For example, as illustrated in the first row of Figure 1, although the triangular object has "salient-like" appearance, it should not be labeled as salient object, since it is not the most interesting and prominent part of the image. From Figure 1, we can see that our method (denoted as GAGNet-R) is able to completely suppress the whole triangular object. Furthermore, by adopting the Guide Module, high-level features, which have the ability of category recognition of image regions because of containing high semantic information, can guide the selection of low-level features. By inspiring from the Channel Attention Block (CAB) proposed in [17], we give our model the ability to guide the selection of

low-level features, which equips our network with the power of assigning foreground label to different-looking salient regions. As illustrated in the second row of Figure 1, the appearance of the feet of the doll is different from the rest of the doll, but as it can be seen, our method is able to highlight the whole doll as the salient object. Thus, by benefiting from the content-aware guidance provided by our Guide Modules, our method is able to handle these complicated scenarios.

Some previous salient object detection methods [18, 19, 20] utilize subsequent single-scale convolutional and max pooling layers to produce deep features. Since salient objects have large variations in scale and location, the learned features by these methods might not be able to handle these complicated variations due to the limited field of view. To extract multi-scale contextual information, some methods [21, 22] apply several parallel dilated convolutions with different rates inspired by structures such as Atrous Spatial Pyramid Pooling (ASPP) [23]. However, the dilated convolution inserts "holes" in the convolution kernels to enlarge the receptive field, which would cause the loss of local information, especially when the dilation rate increases. This problem is called the "gridding issue" which was explored in [24]. To address these problems, different from the methods based on dilated convolution and ASPP, we introduce the Multi-scale Feature Extraction Module (MFEM) which is capable of capturing multi-scale contextual information by enabling densely connections within the multi-scale regions in the feature map. For each level of abstraction (i.e., stage) of the pre-trained backbone, we perform convolutions by adopting a $3 \times 3$ trivial convolutional layer and Global Convolutional Networks (GCNs) [25] with different kernel sizes. Then, the resulting feature maps are stacked to form multi-scale features. GCNs enable densely connections within a large $k \times k$ region in the feature map and thus can alleviate the "gridding issue".

In this paper, we propose a Content-Aware Guidance Network, which we refer to as CAGNet, consisting of three networks: (i) Feature Extraction Network (FEN), (ii) Feature Guide Network (FGN), (iii) Feature Fusion Network (FFN).

The FEN produces multi-scale features at multiple levels of abstraction by adopting the MFEM at each level of a pre-trained backbone. The FGN takes the extracted multi-scale features of the FEN as its input and guides the features in order to use by the FFN. Then, by using multiple add operations and Residual Refinement Modules (RRMs) in the FFN, the guided features are fused effectively. Our proposed RRM is a residual block with spatial attention, which refines the features with the ability of focusing on salient regions and avoiding distractions in the non-salient regions. In summary, the FEN, FGN, and FFN in our proposed architecture work collaboratively to generate more accurate saliency prediction. Additionally, while most saliency detection methods in the literature use the Cross-entropy loss to learn the salient objects, we design

a loss function that outperforms the Cross-entropy by a large margin.

In this paper, by conducting experiments on various backbones, we prove the robustness of our method. Furthermore, our method contains a lower number of parameters in comparison with the previous state-of-the-art methods. It is worth mentioning that since salient object detection is a pre-processing step for many computer vision tasks, it is important to evaluate the performance in terms of the running speed. Our method is capable of running at a real-time speed of 28 FPS, which guarantees that our network can be practically adopted as a pre-processing step for computer vision tasks.

In short, our main contributions are summarized as follows:

- We propose the Feature Guide Network to equip our model with the power of i) making the foreground and background regions more distinct and suppressing the non-salient regions that have "salient-like" appearance; ii) detecting salient objects that have different-looking regions.

- To extract powerful multi-scale features, different from the methods based on dilated convolution and ASPP, we propose a Multi-scale Feature Extraction Module which adopts GCNs to enable densely connections within large regions in the feature map. This module helps the model to alleviate the "gridding issue".

- We design a loss function that outperforms the widely used Cross-entropy loss by a large margin.

- Our method achieves great performance under different backbones, which shows that our proposed framework is very general with respect to the choice of the backbone model. It is interesting to note that while most methods in the saliency detection literature adopt a single backbone in their framework, we evaluate our framework on four different backbones to prove the generalization capability of our method.

- The proposed method achieves the state-of-the-art on several challenging saliency detection datasets. Furthermore, our method contains a lower number of parameters compared to the previous state-of-the-art methods and can run at a real-time speed of 28 FPS.

## 2. Related work

### 2.1. Video Salient Object Detection

Video salient object detection has been widely applied in video segmentation, video compression, video captioning, weakly supervised attention, and autonomous driving. For example, Wang et al. [26] propose an unsupervised method that incorporates geodesic distance into saliency-aware video object segmentation.

To detect salient object in videos, Wang et al. [27] propose a deep video saliency model which employs FCNs for pixel-wise prediction. This model is composed of two components, namely static saliency network and dynamic saliency network, which give the model the ability to capture spatial and temporal statistics of dynamic scenes. Shao et al. [28] introduce a method that identifies foreground regions in videos by using object proposals. This method detects salient objects by ranking and selecting the salient proposals based on object-level saliency cues. Fan et al. [29] propose a model for video salient object detection that is composed of two main components. The first component is Pyramid Dilated Convolution (PDC) which aims at robust static saliency representation learning. The second one is Saliency-Shift-aware convLSTM (SSLSTM) that is able to capture video saliency dynamics through modeling human visual attention-shift behavior. Lu et al. [30] propose a novel CO-attention Siamese Network (COSNet) which uses a co-attention mechanism to capture the temporal correlation across frames.

### 2.2. Stereoscopic Image Salient Object Detection

Some of the saliency detection methods are proposed to detect the salient regions in a stereoscopic three-dimensional (3D) image [31, 32]. Niu et al. [31] introduce a method for saliency detection, which is based on the global disparity contrast in a pair of stereo images. Wang et al. [32] propose a stereo saliency detection algorithm that considers stereoscopic information and the relevancy between the two views of a stereo pair. Different from these methods, in this paper, we focus on salient object detection in RGB Images.

### 2.3. RGB Image Salient Object Detection

Over the past years, numerous methods have been proposed for RGB image saliency detection. Traditional methods predict the saliency score based on hand-crafted features. Most of these methods utilize heuristic priors such as center prior [10, 11], boundary background [33], and color contrast [34]. Aytekin et al. [35] propose a probabilistic framework to encode the boundary connectivity saliency cue and smoothness constraints into a global optimization problem. Wang et al. [36] propose a saliency transfer method to benefit from the existing large annotated datasets for recognizing the primary and smooth connected salient regions from an image. Shan et al. [37] propose a graph-based approach and use background weight map to provide seeds for manifold ranking. Furthermore, they design a third-order smoothness framework to enhance the performance of manifold ranking. These methods, which are based on the traditional approaches, fail to capture semantic and high-level information of the objects.

Recently, deep Convolutional Neural Networks (CNNs) have shown their capabilities in extracting powerful features at multiple levels of abstraction. The CNN features can acquire a richer representation compared

to the traditional hand-crafted features, and thus would result in performance improvement. In recent years, a vast number of methods have adopted CNNs for saliency detection task. For example, Li et al. [38] extract multi-scale features from a CNN and estimate the saliency score for each image super-pixel. Wang et al. [19] employ two CNNs to combine local estimation of super-pixels and global proposal searching to predict saliency maps. Zhao et al. [39] propose multi-context CNNs for exploiting local and global context for salient object detection. Although these CNN-based methods have shown better performance than the traditional methods, they are time-consuming because of taking image patches as input. Moreover, these methods fail to consider important spatial information of the whole image.

To overcome the above-mentioned problems, several methods have utilized FCNs to generate a pixel-wise prediction over the whole image directly. For instance, Li et al. [40] propose a multi-scale FCN to explore the semantic properties and visual contrast information of salient objects. Hou et al. [41] introduce short connections to combine features in different layers. Zhang et al. [13] propose a resolution-based feature combination module to integrate multi-level feature maps into multiple resolutions, which captures spatial details and semantic information simultaneously. Then, by fusing the predicted saliency maps in each resolution, the final saliency map is obtained. Zhang et al. [21] design a bi-directional message passing architecture to pass messages between multi-level features. Wang et al. [42] propose to locate the salient objects globally and then refine them by taking advantage of local context information. Zhang et al. [43] employ a hyper-densely hierarchical feature fusion network to fuse the local and global multi-scale feature maps.

Most of the recent methods focus on using both high-level and low-level features for salient object detection. However, naively using these features may result in confusion for the network, and there needs to be an effective approach to use these features constructively. In this paper, we propose the Feature Guide Network which guides multi-level features to produce more effective features.

To obtain multi-scale features, some previous methods adopt parallel networks and feed them with re-scaled images [44] or multi-context super-pixels [38]. Different from these methods, we propose Multi-scale Feature Extraction Module (MFEM) to extract multi-scale features.

## 3. Our method

In this section, we first explain our proposed Content-Aware Guidance Network (CAGNet) containing three networks: (i) Feature Extraction Network which extracts multi-scale context information, (ii) Feature Guide Network which guides the extracted features by taking advantage of the spatial details of low-level
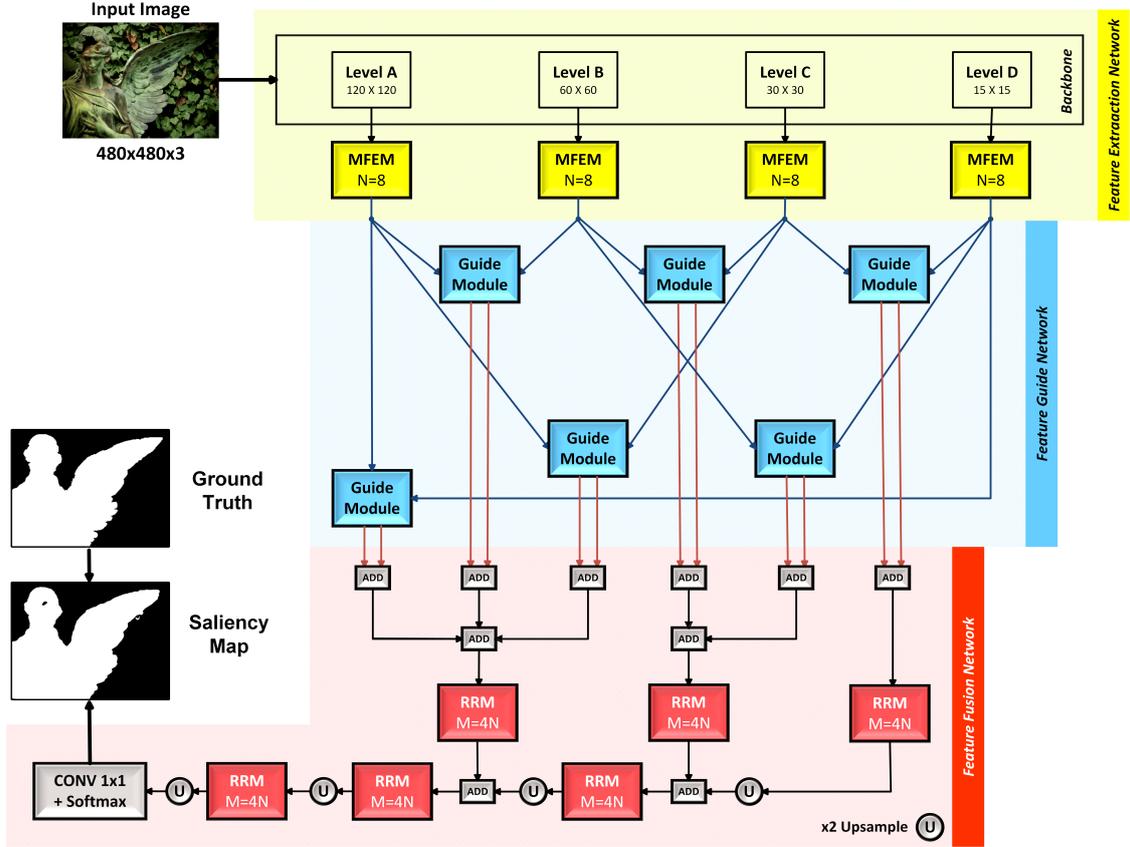
Figure 2: The overall architecture of our proposed Content-Aware Guidance Network (CAGNet). CAGNet consists of three networks: (i) Feature Extraction Network which captures multi-scale contextual features, (ii) Feature Guide Network which guides the extracted features by taking advantage of the nature of high-level and low-level features, (iii) Feature Fusion Network which fuses the guided features effectively to generate the saliency map.

features and the semantic information of high-level features, (iii) Feature Fusion Network which integrates the guided features effectively to generate the saliency map. The architecture of the proposed CAGNet is illustrated in Figure 2. Finally, we describe our designed loss function that has better performance than the widely used Cross-entropy loss.

### 3.1. Feature Extraction Network

Feature Extraction Network consists of a pre-trained backbone that takes the input image and produces multi-level feature maps, and Multi-scale Feature Extraction Modules (MFEMs) which we apply them to multi-level feature maps to capture multi-scale contextual features.

Table 1: Selected layers for different levels of abstraction in the adopted backbones. Note that we take the output of these layers for feature extraction. The size of the feature maps are shown in parentheses. The Level D in CAGNet-V (which is denoted as The added layer) is obtained by adding 1024 kernels of size $3 \times 3$ after the last max pooling layer of the original VGG-16.

| Backbone | Level A | Level B | Level C | Level D |
|---|---|---|---|---|
| VGG-16 [45] | Conv3-3 | Conv4-3 | Conv5-3 | The added layer |
| | $(120 \times 120 \times 256)$ | $(60 \times 60 \times 512)$ | $(30 \times 30 \times 512)$ | $(15 \times 15 \times 1024)$ |
| ResNet50 [46] | Conv2-x | Conv3-x | Conv4-x | Conv5-x |
| | $(120 \times 120 \times 256)$ | $(60 \times 60 \times 512)$ | $(30 \times 30 \times 1024)$ | $(15 \times 15 \times 2048)$ |
| NASNet-Mobile [47] | 1st reduction cell | 4th normal cell | 8th normal cell | 12th normal cell |
| | $(120 \times 120 \times 44)$ | $(60 \times 60 \times 264)$ | $(30 \times 30 \times 528)$ | $(15 \times 15 \times 1056)$ |
| NASNet-Large [47] | 1st reduction cell | 6th normal cell | 12th normal cell | 18th normal cell |
| | $(120 \times 120 \times 168)$ | $(60 \times 60 \times 1008)$ | $(30 \times 30 \times 2016)$ | $(15 \times 15 \times 4032)$ |

### 3.1.1. Pre-trained backbone

In this study, we examine different pre-trained models in our CAGNet as the backbone model, namely VGG-16 [45], ResNet50 [46], NASNet-Mobile [47], and NASNet-large [47], which are denoted as CAGNet-V, CAGNet-R, CAGNet-M, and CAGNet-L, respectively. These backbones are used to produce features at different levels of abstraction. To fit the need of saliency detection task, we remove all fully connected layers in these backbones. In VGG-16, the features after the last max pooling layer cannot introduce a new level of abstraction. Thus, we use a convolutional layer with 1024 kernels of size $3 \times 3$ after the last max pooling layer in VGG-16 to produce a new level.

The output feature maps of all backbones are re-scaled by a factor of 32 with respect to the input image. We take feature maps at four levels from each backbone. Given an input image with size $W \times H$, these feature maps have spatial sizes of $W/2^n \times H/2^n$ with $n = 2, 3, 4, 5$. The details of selected layers for different levels of abstraction in each backbone are shown in Table 1.

### 3.1.2. Multi-scale Feature Extraction Module

Salient objects have large variations in scale and location in different images. Due to the variability of scale, using single scale convolution may not capture the right size. Moreover, due to the variability of location, using pyramid pooling as a multi-scale feature extractor, as proposed in [16], would cause the loss of important local information because of the large scale of pooling. Another approach to implement a multi-scale feature extractor is to use dilated convolutions like [21], which enlarges the receptive field by inserting "holes" in the convolution kernels, and thus would result in the loss of local information because of sparse connections. This problem, which is called the "gridding issue", was explored in [24].

Based on above observation, we find the Global Convolutional Networks (GCNs) [25] effective to address
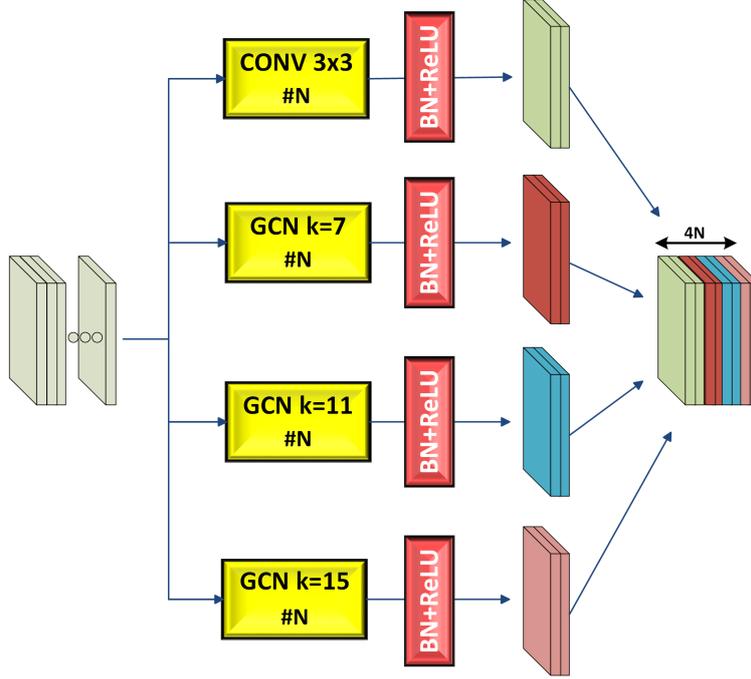
Figure 3: Multi-scale Feature Extraction Module (MFEM). MFEM adopts the $3 \times 3$ trivial convolution and GCNs with $k = 7, 11, 15$ to extract multi-scale features. The '#' symbol denotes the number of layer filters. This figure shows MFEM with N=2.

the "gridding issue" challenge. To avoid sparse connections and enable densely connections within a large $k \times k$ region in the feature map, GCN utilizes a combination of $k \times 1 + 1 \times k$ and $1 \times k + k \times 1$ convolutions to implement the $k \times k$ convolution effectively with a lower number of parameters compared to the trivial $k \times k$ convolution. More details about the GCN can be found in [25]. Furthermore, to obtain multi-scale contextual information, by taking advantage of GCNs, we propose the Multi-scale Feature Extraction Module (MFEM). This module contains GCNs with different kernel sizes and can learn multi-scale context information at multiple levels of abstraction.

As illustrated in Figure 3, in the MFEM, we perform convolutions by utilizing the $3 \times 3$ trivial convolution and GCNs with $k = 7, 11, 15$. Then, the resulting feature maps are concatenated to form multi-scale features.

*3.2. Feature Guide Network*

We employ four MFEMs in the Feature Extraction Network to extract multi-scale features at four different levels of abstraction. These different levels have different recognition information. High-level features have semantic and global information because of the large field of view. Thus, these features can help the category recognition of image regions. Low-level features have spatial and local information due to

the small field of view. Therefore, the information of low-level features can help to better locate the salient regions.

Based on above observation, we propose the Feature Guide Network to better exploit the diverse recognition abilities of different levels. Feature Guide Network is composed of multiple Guide Modules which help to produce more powerful features for saliency detection. As illustrated in Figure 4, Guide Module consists of Low-level Guide and High-level Guide branches. This module takes low-level and high-level features as its inputs and outputs guided low-level and guided high-level features.

In saliency detection, some non-salient regions may have "salient-like" appearance. As shown in the first row of Figure 1, the triangular object at the bottom of the image, which has "salient-like" appearance, may cause confusion for saliency prediction. To address this challenge, we take advantage of the nature of the lower levels to guide the higher levels. In the lower levels, the Feature Extraction Network captures finer spatial information because of its smaller field of view compared to the higher levels. Thus, by applying a $1 \times 1$ convolution on concatenated high-level and low-level features, spatial weights are produced to weight the spatial information of high-level features. With this design, high-level features, which lack the low-level cues, can exploit the fine spatial details of low-level features as a guidance to make salient and non-salient regions more distinguishable. Therefore, by guiding the spatial information of high-level features, our network is able to enhance the distinction of salient and non-salient regions and suppress the non-salient regions with "salient-like" appearance.

In some complicated scenarios, salient regions may have different appearances. As illustrated in the second row of Figure 1, the appearance of the feet of the doll is different from the rest of the doll. Assigning foreground label to these different-looking regions is challenging. To address this challenge, by inspiring from the Channel Attention Block (CAB) proposed in [17], we use the nature of high-level features to guide low-level features in our Feature Guide Network. High-level features have higher semantic information due to the large receptive field. By applying an architecture like Squeeze and Excitation Networks [48] on concatenated high-level and low-level features, channel weights are generated to weight the channels of low-level feature maps. In this way, by utilizing high-level semantic information, the low-level features are guided to produce more attentive features. Thus, Guide Modules provide content-aware guidance for multi-level features, which would result in more accurate saliency prediction.

### 3.3. Feature Fusion Network

By adopting Feature Extraction Network and Feature Guide Network, guided multi-scale features at different levels of abstractions are obtained. To integrate these features effectively, we propose the Feature
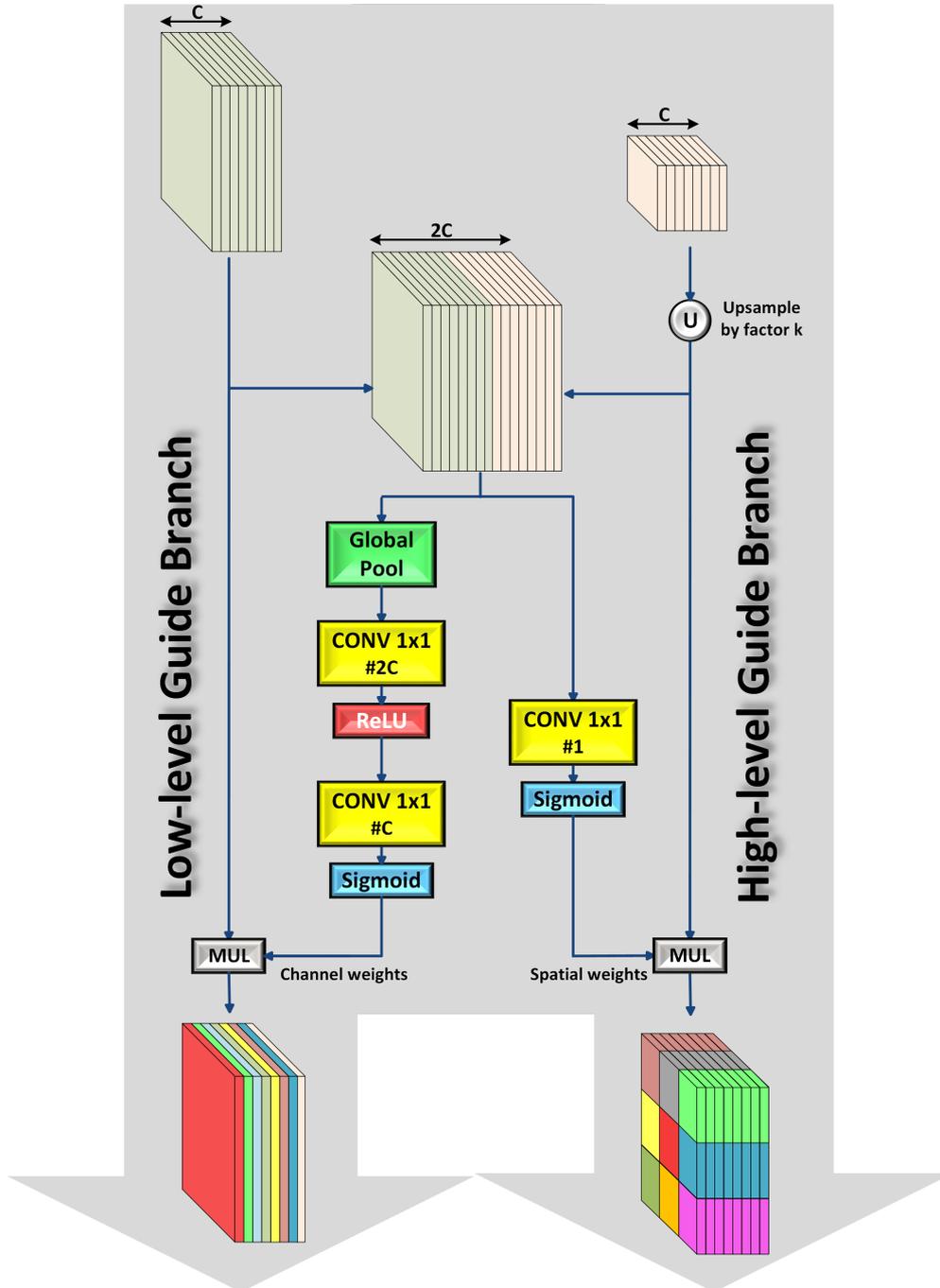
Figure 4: The illustration of the Guide Module. This module consists of High-level Guide and Low-level Guide branches and is adopted to guide the features of the different levels. Note that C shows the number of the channels of the input feature maps, and the '#' symbol denotes the number of layer filters.
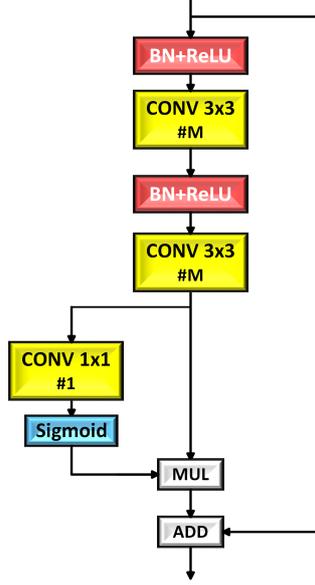
Figure 5: Residual Refinement Module (RRM). RRM is a residual block with spatial attention and is adopted to refine the features effectively. The '#' symbol denotes the number of layer filters.

Fusion Network. In this network, we use add operations to combine different feature maps. In order to refine the features effectively, we introduce Residual Refinement Module (RRM), which is schematically depicted in Figure 5. RRM is a residual block [46, 49] with spatial attention. This module is used to refine the features and has the ability of focusing on salient regions and avoiding distractions in the non-salient regions. Several works such as [17] and [25] have used the standard residual block as their refinement module. Following their works, we have used the full pre-activation version of the residual block along with a spatial attention branch.

By adopting multiple RRM modules and add operations in Feature Fusion Network, finally the saliency map is obtained by utilizing a convolutional layer with two $1 \times 1$ kernels with softmax activation.

### 3.4. Our designed loss function for learning the salient objects

In saliency detection literature, Cross-entropy loss function is widely used for learning the salient objects. However, the networks trained with Cross-entropy loss often differentiate boundary pixels with low confidence, which would result in performance degradation. In this paper, we design a loss function that leads to better results compared to the Cross-entropy loss, as shown in the ablation analysis section. Let $I = \{I_m, m = 1, ..., M\}$, $S_m$, and $G_m$ denote the training images, saliency map for the $m$-th training image,

and ground truth for the *m-th* training image, respectively. Our designed loss is formulated as:

$$L = \alpha_1 L_P + \alpha_2 L_R + \alpha_3 L_{MAE} \tag{1}$$

where $\alpha_1$ , $\alpha_2$ and $\alpha_3$ are the balance parameters. We empirically set $\alpha_1 = 1$, $\alpha_2 = 0.5$, and $\alpha_3 = 1$. $L_P$ and $L_R$ are computed as:

$$L_P = 1 - \frac{1}{M} \sum_{m=1}^{M} P(S_m, G_m) \tag{2}$$

$$L_R = 1 - \frac{1}{M} \sum_{m=1}^{M} R(S_m, G_m) \tag{3}$$

where $P(S.G)$ and $R(S, G)$ are calculated similar to Precision and Recall:

$$P(S, G) = \frac{\sum_n s_n g_n}{\sum_n s_n + \epsilon} \tag{4}$$

$$R(S, G) = \frac{\sum_n s_n g_n}{\sum_n g_n + \epsilon} \tag{5}$$

where $s_n \in S$ and $g_n \in G$ , and $\epsilon$ is a regularization constant. $L_{MAE}$ calculates the discrepancy between the predicted saliency map $S$ and the ground truth $G$:

$$L_{MAE} = \frac{1}{M} \sum_{m=1}^{M} MAE(S_m, G_m) \tag{6}$$

where $MAE(S, G)$ is computed as :

$$MAE(S, G) = \frac{1}{N} \sum_n | s_n - g_n | \tag{7}$$

where $N$ denotes the total number of pixels. In ablation analysis section, we demonstrate that our designed loss function outperforms the Cross-entropy loss function.

## 4. Experiments

### 4.1. Datasets and evaluation metrics

The proposed method is evaluated on six public salient object detection datasets. ECSSD [50] contains 1,000 semantically meaningful and complex images with multiple objects of different sizes. DUT-OMRON [33] consists of 5,168 challenging images with high variety of content, each of which has complex background and one or two salient objects. HKU-IS [38] contains 4447 images with low color contrast. Images in this dataset are selected to include multiple foreground objects or objects touching the image boundary. DUTS [51] dataset is currently the largest salient object detection dataset and comprised of 10,553 images in the training set and 5,019 images in the test set. Both training and test sets have very challenging scenarios. The PASCAL-S [52] dataset has 850 natural images chosen from the PASCAL VOC 2010 [53] segmentation dataset. In SOC [54] dataset, each salient image is accompanied by attributes that reflect common challenges in real-world scenarios. We use the validation set of this dataset for testing.

We use six metrics to evaluate the performance of our method as well as previous state-of-the-art saliency detection methods, namely Precision-Recall (PR) curves, F-measure curves, Average F-measure (denoted as avgF) score, weighted F-measure (denoted as wF) score, E-measure (denoted as E) score, and Mean Absolute Error (denoted as MAE) score. More detailed descriptions about these metrics can be found in [55, 56, 57]. Furthermore, the code for computing these metrics in Python can be found at https://github.com/Mehrdad-Noori/Saliency-Evaluation-Toolbox.

Precision is the fraction of correct salient pixels in the predicted saliency map, and Recall is defined as the fraction of correct salient pixels in the ground truth. To calculate Precision and Recall, the binarized saliency map is compared against the ground truth mask. The threshold is varied from 0 to 1 to generate a sequence of binary masks. These binary masks are used to calculate (Precision, Recall) pairs and (F-measure, threshold) pairs to plot the PR curves and the F-measure curves.

The Average F-measure score is calculated by using the thresholding method suggested in [58]. This threshold is used to generate binary maps for computing the F-measure which is defined as:

$$F_\beta = \frac{(1 + \beta^2) \cdot Precision \cdot Recall}{\beta^2 \cdot Precision + Recall} \tag{8}$$

where $\beta^2$ is set to 0.3 to weight precision more than recall. The weighted F-measure score [56] and E-measure score [57] are also adopted to evaluate the performance. Finally, the MAE score is calculated as the average pixel-wise absolute difference between the ground truth mask and the predicted saliency map.
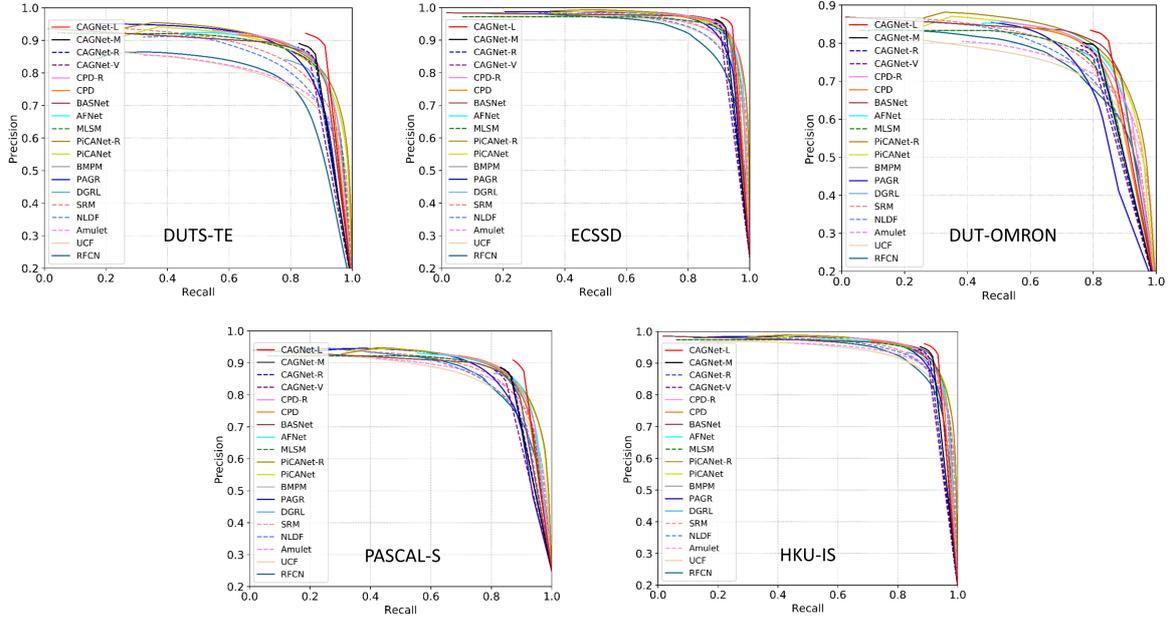
Figure 6: The PR curves of the proposed method and previous state-of-the-art methods.

## 4.2. Implementation details

We develop our proposed method in Keras [59] framework using TensorFlow [60] backend. The backbone models (i.e., VGG-16 [45], ResNet-50 [46], NASNet Mobile [47], and NASNet Large [47]) are initialized with ImageNet [61] weights. In our experiments, the input image is resized into $480 \times 480$ pixels for training and testing. To reduce overfitting, two types of data augmentation are randomly employed: horizontal filliping and rotation (range of 0-12 degrees). We do not use validation set and train the model until its training loss converges. All the experiments are performed using the stochastic gradient descent with a momentum coefficient 0.9 and an initial learning rate of $8e$-3 which is divided by 10 if no improvement in training loss is seen for 10 epochs. We perform our experiments on an NVIDIA 1080 Ti GPU. The code, the trained models, and the saliency maps of our method can be found at https://github.com/Mehrdad-Noori/CAGNet.

## 4.3. Comparison with the state-of-the-art

We compare our method with 16 previous state-of-the-art methods, namely MDF [38], RFCN [18], UCF [20], Amulet [13], NLDF [12], DSS [41], BMPM [21], PAGR [62], PiCANet [63], SRM [16], DGRL [42], MLMS [64], AFNet [65], CapSal [66], BASNet [15], and CPD [67]. For a fair comparison, we use the saliency maps provided by the authors.
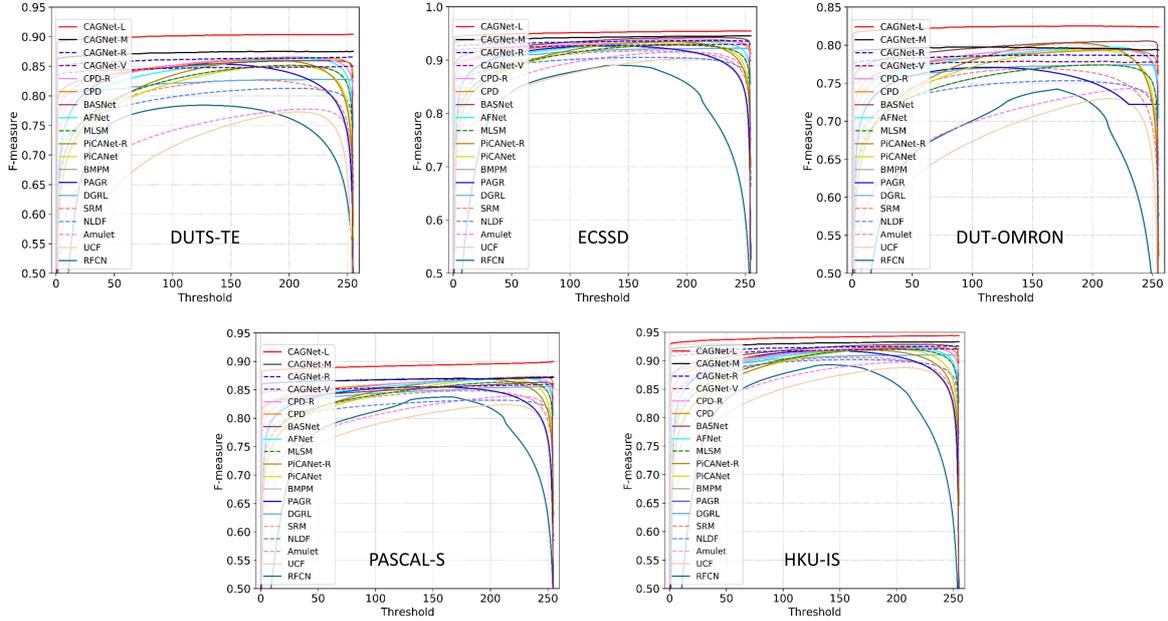
Figure 7: The F-measure curves of the proposed method and previous state-of-the-art methods.

*Quantitative Evaluation.* P-R curves and F-measure curves on the five datasets are shown in Figure 6 and Figure 7, respectively. We can see that our proposed method performs favorably against other methods in all cases. Especially, it is obvious that our CAGNet-L performs better than all other methods by a relatively large margin. Moreover, we compare our method with other previous state-of-the-art methods in terms of avgF, wF, E-measure, and MAE score on five benchmark datasets in Table 2. As seen from this table, our method ranks first in most cases. We also evaluate attributes-based performance on the challenging SOC dataset. The results of the recent state-of-the-art methods on nine attributes of SOC and their average are shown in Table 3. We can see that our method has achieved great performance and is capable of handling the complex scenarios of SOC dataset. It is interesting to note that our method contains fewer parameters than the existing ones, is end-to-end, and does not need any post-processing steps such as CRF [68]. Another interesting thing about our method is that although our CAGNet-M has significantly fewer parameters than the other networks (only 5.57 million parameters), it has shown outstanding performance. This functionality is desirable for the applications in which we have limitation in terms of the memory. Furthermore, our CAGNet-V has a real-time speed of 28 FPS when processing a $480 \times 480$ image, and therefore it can be practically adopted as a preprocessing step for computer vision tasks.

16

Table 2: Comparison of the proposed method and other 16 methods on five salient object detection datasets in terms of avgF, wF, E-measure, and MAE scores. CAGNet with VGG-16, ResNet50, NASNet-Mobile, and NASNet-Large backbones, are denoted as CAGNet-V, CAGNet-R, CAGNet-M, and CAGNet-L, respectively. The best score and the second best score under each setting are shown in red and blue, respectively, and the best score under all settings is underlined. The unit of the total number of parameters (denoted as #Par) is million. Note that the authors of [62] did not release the code, and they just provided the saliency maps, and thus reporting the total number of parameters of this method is not possible.

| Dataset | Year | Backbone | #Par | DUTS-TE [51] | | | | ECSSD [50] | | | | DUT-O [33] | | | | PASCAL-S [52] | | | | HKU-IS [38] | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metric | | | | avgF | wF | E | MAE | avgF | wF | E | MAE | avgF | wF | E | MAE | avgF | wF | E | MAE | avgF | wF | E | MAE |
| VGG [45] | | | | | | | | | | | | | | | | | | | | | | | |
| MDF [38] | 2015 | VGG16 | 56.86 | 0.669 | 0.588 | 0.813 | 0.093 | 0.807 | 0.705 | 0.853 | 0.105 | 0644 | 0.564 | 0.802 | 0.092 | 0.711 | 0.590 | 0.759 | 0.146 | 0.784 | 0.564 | 0.871 | 0.129 |
| RFCN [18] | 2016 | VGG16 | 134.69 | 0.711 | 0.586 | 0.840 | 0.090 | 0.834 | 0.698 | 0.877 | 0.107 | 0.627 | 0.524 | 0.779 | 0.110 | 0.754 | 0.636 | 0.810 | 0.132 | 0.835 | 0.680 | 0.906 | 0.089 |
| UCF [20] | 2017 | VGG16 | 23.98 | 0.631 | 0.596 | 0.770 | 0.112 | 0.844 | 0.806 | 0.895 | 0.069 | 0.621 | 0.573 | 0.768 | 0.120 | 0.738 | 0.700 | 0.809 | 0.116 | 0.823 | 0.779 | 0.904 | 0.062 |
| Amulet [13] | 2017 | VGG16 | 33.15 | 0.678 | 0.658 | 0.803 | 0.085 | 0.868 | 0.840 | 0.912 | 0.059 | 0.647 | 0.626 | 0.784 | 0.098 | 0.771 | 0.741 | 0.831 | 0.099 | 0.841 | 0.817 | 0.914 | 0.051 |
| NLDF [12] | 2017 | VGG16 | 35.49 | 0.739 | 0.710 | 0.855 | 0.065 | 0.878 | 0.839 | 0.912 | 0.063 | 0.684 | 0.634 | 0.817 | 0.080 | 0.782 | 0.742 | 0.842 | 0.101 | 0.873 | 0.838 | 0.929 | 0.048 |
| DSS [41] | 2017 | VGG16 | 62.23 | 0.716 | 0.702 | 0.845 | 0.065 | 0.873 | 0.836 | 0.915 | 0.062 | 0.674 | 0.643 | 0.820 | 0.074 | 0.776 | 0.728 | 0.836 | 0.103 | 0.856 | 0.821 | 0.926 | 0.050 |
| PAGR [62] | 2018 | VGG19 | — | 0.784 | 0.724 | 0.883 | 0.055 | 0.894 | 0.833 | 0.917 | 0.061 | 0.711 | 0.622 | 0.843 | 0.071 | 0.808 | 0.738 | 0.854 | 0.095 | 0.886 | 0.820 | 0.939 | 0.047 |
| BMPM [21] | 2018 | VGG16 | 22.09 | 0.745 | 0.761 | 0.863 | 0.049 | 0.868 | 0.871 | 0.916 | 0.045 | 0.692 | 0.681 | 0.839 | 0.064 | 0.771 | 0.785 | 0.847 | 0.075 | 0.871 | 0.859 | 0.938 | 0.039 |
| PiCANet [63] | 2018 | VGG16 | 32.85 | 0.749 | 0.747 | 0.865 | 0.054 | 0.885 | 0.865 | 0.926 | 0.046 | 0.710 | 0.691 | 0.842 | 0.068 | 0.804 | 0.781 | 0.862 | 0.079 | 0.870 | 0.847 | 0.938 | 0.042 |
| MLMS [64] | 2019 | VGG16 | 74.38 | 0.745 | 0.761 | 0.863 | 0.049 | 0.868 | 0.871 | 0.916 | 0.044 | 0.692 | 0.681 | 0.839 | 0.064 | 0.771 | 0.785 | 0.847 | 0.075 | 0.871 | 0.859 | 0.938 | 0.039 |
| AFNet [65] | 2019 | VGG16 | 21.08 | 0.793 | 0.785 | 0.895 | 0.046 | 0.908 | 0.886 | 0.941 | 0.042 | 0.738 | 0.717 | 0.859 | 0.057 | 0.828 | 0.804 | 0.887 | 0.071 | 0.888 | 0.869 | 0.947 | 0.036 |
| CPD [67] | 2019 | VGG16 | 29.32 | 0.813 | 0.801 | 0.908 | 0.043 | 0.914 | 0.895 | 0.943 | 0.040 | 0.745 | 0.715 | 0.869 | 0.057 | 0.832 | 0.806 | 0.884 | 0.074 | 0.895 | 0.879 | 0.950 | 0.033 |
| **CAGNet-V** | - | VGG16 | 20.98 | 0.823 | 0.797 | 0.904 | 0.044 | 0.915 | 0.893 | 0.939 | 0.042 | 0.744 | 0.718 | 0.860 | 0.057 | 0.831 | 0.799 | 0.881 | 0.077 | 0.906 | 0.886 | 0.947 | 0.033 |
| ResNet [46] | | | | | | | | | | | | | | | | | | | | | | | |
| SRM [16] | 2017 | ResNet50 | 43.74 | 0.753 | 0.722 | 0.867 | 0.059 | 0.892 | 0.853 | 0.927 | 0.054 | 0.707 | 0.658 | 0.843 | 0.069 | 0.803 | 0.762 | 0.861 | 0.087 | 0.874 | 0.835 | 0.938 | 0.046 |
| DGRL [42] | 2018 | ResNet50 | 126.35 | 0.794 | 0.774 | 0.899 | 0.050 | 0.906 | 0.891 | 0.946 | 0.041 | 0.733 | 0.709 | 0.856 | 0.062 | 0.827 | 0.802 | 0.891 | 0.073 | 0.890 | 0.875 | 0.949 | 0.036 |
| PiCANet-R [63] | 2018 | ResNet50 | 37.02 | 0.759 | 0.755 | 0.873 | 0.051 | 0.886 | 0.867 | 0.927 | 0.046 | 0.717 | 0.695 | 0.848 | 0.065 | 0.804 | 0.782 | 0.862 | 0.078 | 0.870 | 0.840 | 0.940 | 0.043 |
| CapSal [66] | 2019 | ResNet101 | 91.09 | 0.755 | 0.689 | 0.867 | 0.063 | — | — | — | — | — | — | — | — | 0.827 | 0.791 | 0.878 | 0.074 | 0.841 | 0.780 | 0.905 | 0.058 |
| BASNet [15] | 2019 | ResNet34 | 87.06 | 0.791 | 0.803 | 0.884 | 0.047 | 0.880 | 0.904 | 0.921 | 0.037 | 0.756 | 0.751 | 0.869 | 0.056 | 0.781 | 0.800 | 0.853 | 0.077 | 0.895 | 0.889 | 0.946 | 0.032 |
| CPD-R [67] | 2019 | ResNet50 | 47.85 | 0.805 | 0.795 | 0.904 | 0.043 | 0.917 | 0.898 | 0.949 | 0.037 | 0.747 | 0.719 | 0.873 | 0.056 | 0.831 | 0.803 | 0.887 | 0.072 | 0.891 | 0.875 | 0.950 | 0.034 |
| **CAGNet-R** | - | ResNet50 | 26.06 | 0.838 | 0.817 | 0.914 | 0.040 | 0.921 | 0.903 | 0.944 | 0.037 | 0.753 | 0.729 | 0.862 | 0.054 | 0.847 | 0.820 | 0.896 | 0.067 | 0.910 | 0.893 | 0.950 | 0.030 |
| NASNet [47] | | | | | | | | | | | | | | | | | | | | | | | |
| **CAGNet-M** | - | Mobile | 5.57 | 0.852 | 0.832 | 0.923 | 0.037 | 0.933 | 0.916 | 0.953 | 0.034 | 0.764 | 0.743 | 0.864 | 0.052 | 0.846 | 0.819 | 0.893 | 0.069 | 0.919 | 0.904 | 0.956 | 0.028 |
| **CAGNet-L** | - | Large | 89.42 | 0.886 | 0.871 | 0.944 | 0.029 | 0.943 | 0.932 | 0.963 | 0.026 | 0.798 | 0.779 | 0.889 | 0.047 | 0.877 | 0.858 | 0.921 | 0.053 | 0.932 | 0.921 | 0.965 | 0.024 |

Table 3: Attributes-based performance on the challenging SOC dataset. The best score and the second best score are shown in red and blue, respectively,

| Attribute | AC | | BO | | CL | | HO | | MB | | OC | | OV | | SC | | SO | | average | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metric | avgF | E | avgF | E | avgF | E | avgF | E | avgF | E | avgF | E | avgF | E | avgF | E | avgF | E | avgF | E |
| Amulet [13] | 0.682 | 0.782 | 0.749 | 0.509 | 0.668 | 0.744 | 0.689 | 0.774 | 0.700 | 0.809 | 0.668 | 0.738 | 0.746 | 0.772 | 0.640 | 0.749 | 0.513 | 0.675 | 0.673 | 0.728 |
| NLDF [12] | 0.697 | 0.812 | 0.638 | 0.448 | 0.679 | 0.738 | 0.706 | 0.795 | 0.693 | 0.807 | 0.658 | 0.749 | 0.718 | 0.755 | 0.670 | 0.790 | 0.554 | 0.725 | 0.665 | 0.736 |
| DSS [41] | 0.663 | 0.778 | 0.682 | 0.450 | 0.633 | 0.718 | 0.677 | 0.772 | 0.706 | 0.806 | 0.630 | 0.725 | 0.689 | 0.727 | 0.641 | 0.766 | 0.540 | 0.720 | 0.651 | 0.718 |
| SRM [16] | 0.734 | 0.834 | 0.785 | 0.567 | 0.714 | 0.781 | 0.737 | 0.818 | 0.793 | 0.867 | 0.687 | 0.772 | 0.781 | 0.800 | 0.681 | 0.801 | 0.585 | 0.752 | 0.722 | 0.777 |
| BMPM [21] | 0.712 | 0.815 | 0.530 | 0.392 | 0.657 | 0.745 | 0.704 | 0.802 | 0.734 | 0.846 | 0.670 | 0.769 | 0.729 | 0.771 | 0.670 | 0.803 | 0.568 | 0.757 | 0.664 | 0.745 |
| BASNet [15] | 0.740 | 0.844 | 0.503 | 0.433 | 0.679 | 0.766 | 0.702 | 0.796 | 0.805 | 0.868 | 0.654 | 0.760 | 0.719 | 0.764 | 0.655 | 0.802 | 0.590 | 0.766 | 0.672 | 0.755 |
| CPD-R [67] | 0.765 | 0.860 | 0.824 | 0.662 | 0.741 | 0.808 | 0.766 | 0.845 | 0.812 | 0.880 | 0.741 | 0.813 | 0.799 | 0.831 | 0.727 | 0.837 | 0.635 | 0.796 | 0.757 | 0.815 |
| **CAGNet-V** | 0.769 | 0.852 | 0.772 | 0.611 | 0.732 | 0.795 | 0.775 | 0.851 | 0.813 | 0.886 | 0.729 | 0.809 | 0.791 | 0.817 | 0.739 | 0.841 | 0.630 | 0.794 | 0.750 | 0.806 |
| **CAGNet-R** | 0.741 | 0.848 | 0.799 | 0.675 | 0.726 | 0.791 | 0.764 | 0.835 | 0.849 | 0.899 | 0.727 | 0.803 | 0.781 | 0.807 | 0.747 | 0.856 | 0.633 | 0.790 | 0.752 | 0.812 |
| **CAGNet-M** | 0.765 | 0.856 | 0.784 | 0.649 | 0.743 | 0.814 | 0.795 | 0.861 | 0.825 | 0.899 | 0.747 | 0.833 | 0.801 | 0.828 | 0.752 | 0.851 | 0.647 | 0.800 | 0.762 | 0.821 |
| **CAGNet-L** | 0.785 | 0.873 | 0.827 | 0.719 | 0.754 | 0.827 | 0.816 | 0.876 | 0.879 | 0.925 | 0.749 | 0.829 | 0.828 | 0.852 | 0.780 | 0.881 | 0.692 | 0.829 | 0.790 | 0.846 |

*Qualitative Evaluation.* Some qualitative results are shown in Figure 8. Thanks to the proposed modules, it can be seen that our model is capable of highlighting the inner part of foreground regions in various complicated scenes. Furthermore, our model is able to suppress the background regions which are incorrectly
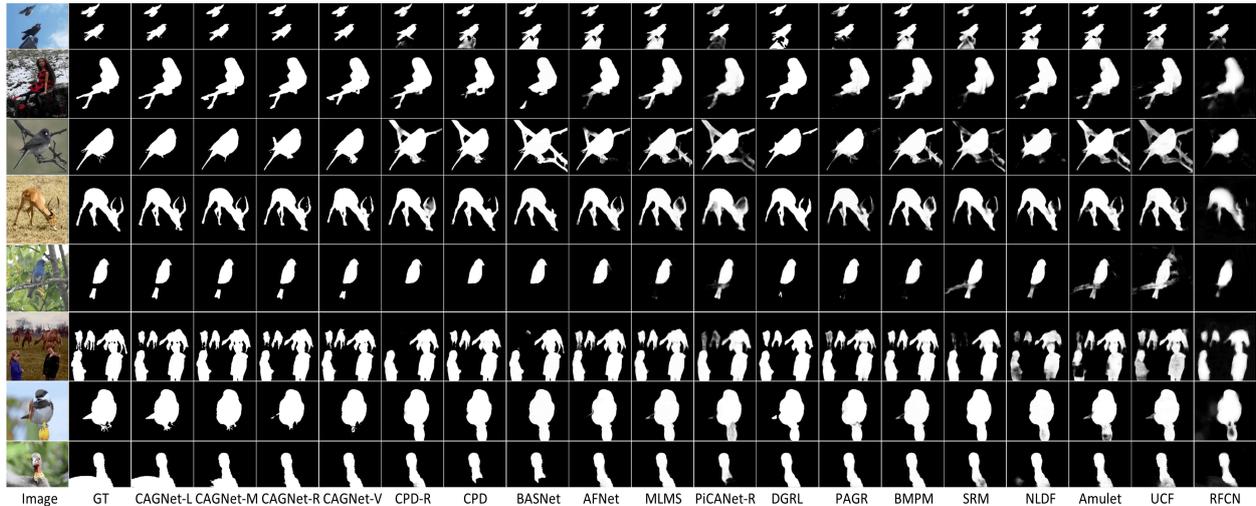
Figure 8: Qualitative comparisons with previous state-of-the-art methods. As it can be seen, our method is capable of predicting saliency maps that are closer to the ground truth compared to the other methods.

labeled by other saliency detection methods. Thus, by taking advantage of different proposed modules, our method is able to handle various complex scenarios.

## 4.4. Ablation analysis

Our proposed CAGNet consists of three modules, namely the Multi-scale Feature Extraction Module (MFEM), the Guide Module, and the Residual Refinement Module (RRM). We perform the ablation analysis on CAGNet-V by using three challenging large-scale datasets, namely DUTS-TE [51], DUT-O [33], and HKU-IS [38]. In order to investigate the effectiveness of each module, we gradually add them to our base network. Our base network is obtained by applying the following modifications to the CAGNet: i) replacing the MFEM modules with $1 \times 1$ convolutions with the same number of filters, ii) removing Guide Modules from the network (which means that the multi-level features are not multiplied by the channel weights and the spatial weights), iii) removing the RRM modules from the model.

We perform ablation analysis by adding each module to our base network in a stepwise manner. The results are shown in Table 4. In this table, the base network is denoted as Base.

*The effectiveness of Guide Module.* We add the High-level Guide branch, Low-level Guide branch, and both High-level and Low-level Guide branches (i.e., the Guide Module) to the base network, which are denoted as HG, LG, and GM, respectively in Table 4. As seen from this table, the performance improves, which shows the beneficial effect of using our Guide Module. Using this module results in i) making salient and non-salient regions more distinct and suppressing the non-salient regions that have "salient-like"

18

Table 4: Ablation analysis of our proposed method with different settings. The best results are shown in **red**.

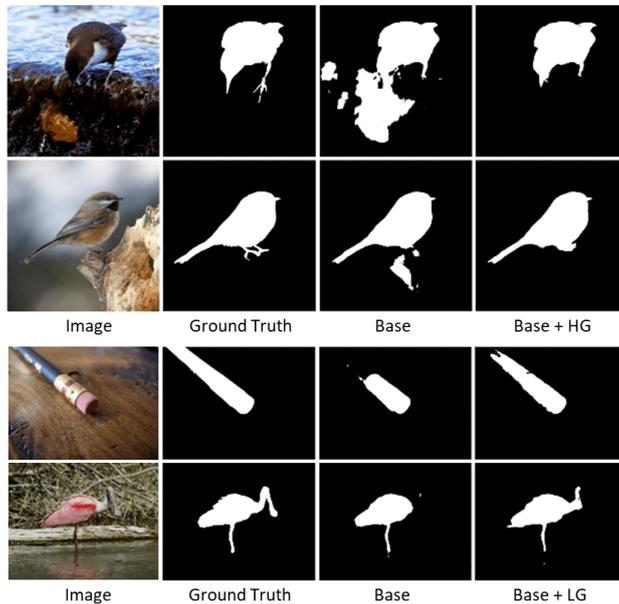| Dataset | DUTS-TE [51] | | | DUT-O [33] | | | HKU-IS [38] | | |
|---|---|---|---|---|---|---|---|---|---|
| Metric | avgF | wF | MAE | avgF | wF | MAE | avgF | wF | MAE |
| Base | 0.7452 | 0.7099 | 0.0627 | 0.6481 | 0.6027 | 0.0862 | 0.8579 | 0.8298 | 0.0482 |
| Base + HG | 0.7598 | 0.7204 | 0.0588 | 0.6571 | 0.6068 | 0.0819 | 0.8700 | 0.8398 | 0.0450 |
| Base + LG | 0.7650 | 0.7268 | 0.0588 | 0.6649 | 0.6178 | 0.0817 | 0.8714 | 0.8419 | 0.0447 |
| Base + GM | 0.7707 | 0.7335 | 0.0558 | 0.6687 | 0.6209 | 0.0780 | 0.8743 | 0.8451 | 0.0432 |
| Base + GM + MFEM | 0.8003 | 0.7779 | 0.0481 | 0.7256 | 0.6971 | 0.0616 | 0.8960 | 0.8776 | 0.0346 |
| **Base + GM + MFEM + RRM (= CAGNet-V)** | **0.8226** | **0.7971** | **0.0445** | **0.7444** | 0.7179 | **0.0571** | **0.9056** | **0.8858** | 0.0332 |
| CE Loss Function | 0.7591 | 0.7517 | 0.0524 | 0.7017 | 0.6793 | 0.0652 | 0.8783 | 0.8558 | 0.0398 |
| Dilated Convolution | 0.8214 | 0.7961 | 0.0457 | 0.7414 | 0.7152 | 0.0583 | 0.9029 | 0.8811 | 0.0343 |
| Trivial Convolution | 0.8166 | 0.7940 | 0.0458 | 0.7439 | **0.7203** | 0.0581 | 0.9041 | 0.8857 | **0.0331** |



Figure 9: Visual Comparison for two branches of our Guide Module. The first and second rows show the comparison for High-level Guide branch (denoted as HG). The third and fourth rows show the comparison for Low-level Guide branch (denoted as LG).

appearance, ii) assigning foreground label to different-looking salient regions. To further investigate the effectiveness of our guide branches, we show a visual comparison for each branch in Figure 9. As it can be seen, when we add the High-level Guide branch to the base network (Base+HG), the non-salient regions that have "salient-like" appearance are suppressed. Furthermore, when we add the Low-level Guide branch to the base network (Base+LG), different-looking salient regions (head of the pencil and the rest of the pencil, head of the bird and the rest of the bird) are labeled as salient.

*The effectiveness of MFEM.* Based on the aforementioned architecture, we replace the $1 \times 1$ convolutions with the MFEM modules. As seen in Table 4, our proposed MFEM has a beneficial effect on saliency

Table 5: The results of CAGNet-V with different settings for the parameter N. The best results are shown in <span style="color:red">red</span>. The unit of the total number of parameters (denoted as #Par) is million.

| Dataset | DUTS-TE [51] | | | DUT-O [33] | | | HKU-IS [38] | | | #Par |
|---------|------|------|------|------|------|------|------|------|------|------|
| Metric | avgF | wF | MAE | avgF | wF | MAE | avgF | wF | MAE | |
| N=1 | 0.8068 | 0.7777 | 0.0474 | 0.7254 | 0.6932 | 0.0611 | 0.8963 | 0.8734 | 0.0355 | 19.61 |
| N=2 | 0.8100 | 0.7821 | 0.0469 | 0.7339 | 0.7044 | 0.0601 | 0.9020 | 0.8811 | 0.0337 | 19.79 |
| N=4 | 0.8155 | 0.7903 | 0.0461 | 0.7412 | 0.7145 | 0.0576 | 0.9010 | 0.8807 | 0.0336 | 20.17 |
| **N=8** | **0.8226** | 0.7971 | **0.0445** | 0.7444 | 0.7179 | 0.0571 | **0.9056** | **0.8858** | 0.0332 | 20.98 |
| N=16 | 0.8189 | 0.7935 | 0.0451 | 0.7474 | 0.7205 | 0.0567 | 0.9050 | 0.8848 | **0.0328** | 22.86 |
| N=32 | 0.8218 | **0.7973** | 0.0452 | 0.7497 | 0.7255 | 0.0576 | 0.9048 | 0.8843 | 0.0332 | 27.61 |
| N=64 | 0.8189 | 0.7932 | 0.0469 | **0.7523** | **0.7278** | **0.0561** | 0.8980 | 0.8773 | 0.0354 | 41.04 |

detection and improves the results, which shows extracting multi-scale features can help to detect salient objects with different scales and locations.

*The effectiveness of RRM.* To reveal the effect of the RRMs, we add them to the aforementioned architecture. From Table 4, it can be observed that using our refinement module is helpful for saliency detection and improves the performance.

*The effectiveness of our designed loss.* To demonstrate the effectiveness of our designed loss function, we train our CAGNet-V with Cross-entropy, which is denoted as CE Loss Function in Table 4. As seen in this table, our designed loss outperforms the cross-entropy loss by a significant margin.

To further prove the effectiveness of our MFEM, we implement the MFEMs in CAGNet-V by adopting dilated convolutional layers (kernel size=3, dilation rates=1, 3, 5, 7), denoted as Dilated Convolution in Table 4. We can see that the performance degrades, which shows that our proposed MFEM can capture more powerful multi-scale features by enabling densely connections within a large $k \times k$ region in the feature map. We also implement the MFEMs in CAGNet-V by adopting trivial convolutional layers with kernel size=3, 7, 11, 15, denoted as Trivial Convolution in Table 4. As seen from this table, the performance gets worse compared to our CAGNet-V with the proposed MFEM. It is interesting to note that CAGNet-V with our proposed MFEM contains fewer parameters than the CAGNet-V with the MFEM implemented by adopting trivial convolutional layers (20.98 million vs. 27.03 million), which is due to the architectural design of GCNs.

We perform another experiment on CAGNet-V and train it with different settings for the parameter $N$. The results are shown in Table 5. In this paper, by considering the trade-off between the performance and the number of parameters, we have chosen $N = 8$ for our method.
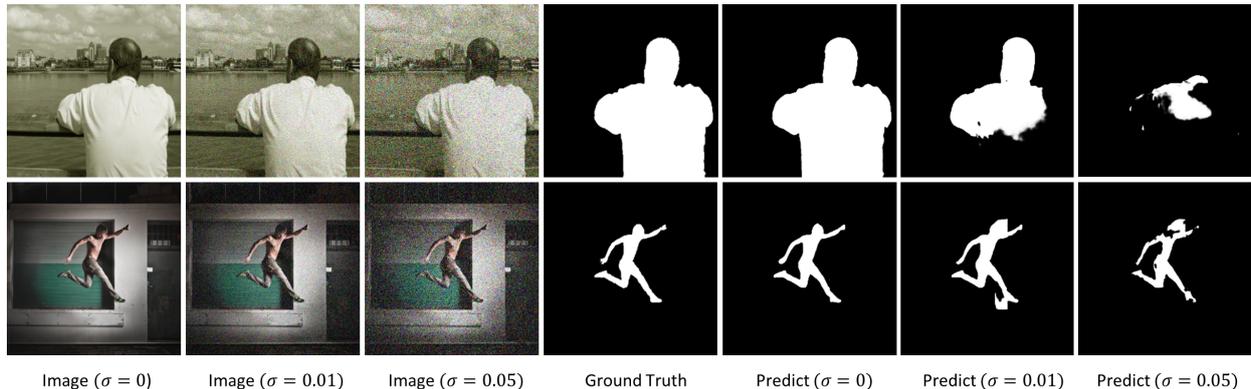
Figure 10: The predictions of our method for two images corrupted by Additive White Gaussian Noise (AWGN) with two variance values ($\sigma = 0.01$ and $\sigma = 0.05$).

## 5. Conclusion and future work

In this paper, we propose a novel end-to-end framework that has the power of i) making the foreground and background regions more distinct and suppressing the non-salient regions that have "salient-like" appearance; ii) detecting salient objects that have different-looking regions. Our proposed model is also capable of capturing multi-scale contextual information effectively. The attentive guided multi-scale features learned by our method and the great results of our deigned loss function prove that a promising approach for saliency detection is introduced in this paper. Experimental evaluations over six datasets demonstrate that our proposed method outperforms the previous state-of-the-art methods under different evaluation metrics.

Based on the great performance and the real-time speed of our approach and its superiority over previous approaches, we plan to use our saliency detector in industrial object-related applications, such as object-based surveillance and object tracking. However, in real-world scenarios, images are affected by noise, which would lead to performance degradation of the most recently introduced saliency detectors [69] including ours. Figure 10 shows the predicted saliency maps for images corrupted by Additive White Gaussian Noise (AWGN). As seen, our model fails to output accurate saliency predictions in the presence of noise. This motivates us to plan on enhancing the robustness of our method by handling noise in an end-to-end approach.

## References

[1] C. F. Flores, A. Gonzalez-Garcia, J. van de Weijer, B. Raducanu, Saliency for fine-grained object recognition in domains with scarce training data, Pattern Recognition 94 (2019) 62–73.

[2] Z. Li, G. Liu, D. Zhang, Y. Xu, Robust single-object image segmentation based on salient transition region, Pattern recognition 52 (2016) 317–331.

[3] X.-H. Zhi, H.-B. Shen, Saliency driven region-edge-based top down level set evolution reveals the asynchronous focus in image segmentation, Pattern Recognition 80 (2018) 241–255.

[4] Q. Cai, H. Liu, Y. Qian, S. Zhou, X. Duan, Y.-H. Yang, Saliency-guided level set model for automatic object segmentation, Pattern Recognition 93 (2019) 147–163.

[5] W. Wang, J. Shen, F. Porikli, Saliency-aware geodesic video object segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 3395–3402.

[6] Y. Chen, Y. Pan, M. Song, M. Wang, Improved seam carving combining with 3d saliency for image retargeting, Neuro-computing 151 (2015) 645–653.

[7] G. Zhang, Z. Yuan, Q. Tong, M. Zheng, J. Zhao, A novel framework for background subtraction and foreground detection, Pattern Recognition 84 (2018) 28–38.

[8] S. Hong, T. You, S. Kwak, B. Han, Online tracking by learning discriminative saliency map with convolutional neural network, in: International conference on machine learning, 2015, pp. 597–606.

[9] L. Huo, L. Jiao, S. Wang, S. Yang, Object-level saliency detection with color attributes, Pattern recognition 49 (2016) 162–173.

[10] A. Aksac, T. Ozyer, R. Alhajj, Complex networks driven salient region detection based on superpixel segmentation, Pattern Recognition 66 (2017) 268–279.

[11] J. Liang, J. Zhou, L. Tong, X. Bai, B. Wang, Material based salient object detection from hyperspectral images, Pattern Recognition 76 (2018) 476–490.

[12] Z. Luo, A. Mishra, A. Achkar, J. Eichel, S. Li, P.-M. Jodoin, Non-local deep features for salient object detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 6609–6617.

[13] P. Zhang, D. Wang, H. Lu, H. Wang, X. Ruan, Amulet: Aggregating multi-level convolutional features for salient object detection, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 202–211.

[14] X. Xi, Y. Luo, P. Wang, H. Qiao, Salient object detection based on an efficient end-to-end saliency regression network, Neurocomputing 323 (2019) 265–276.

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

[16] T. Wang, A. Borji, L. Zhang, P. Zhang, H. Lu, A stagewise refinement model for detecting salient objects in images, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 4019–4028.

[17] C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, N. Sang, Learning a discriminative feature network for semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 1857–1866.

[18] L. Wang, L. Wang, H. Lu, P. Zhang, X. Ruan, Salient object detection with recurrent fully convolutional networks, IEEE Transactions on Pattern Analysis and Machine Intelligence 41 (7) (2019) 1734–1746.

[19] L. Wang, H. Lu, X. Ruan, M.-H. Yang, Deep networks for saliency detection via local estimation and global search, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 3183–3192.

[20] P. Zhang, D. Wang, H. Lu, H. Wang, B. Yin, Learning uncertain convolutional features for accurate saliency detection, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 212–221.

[21] L. Zhang, J. Dai, H. Lu, Y. He, 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.

[22] M. Kampffmeyer, N. Dong, X. Liang, Y. Zhang, E. P. Xing, Connnet: a long-range relation-aware pixel-connectivity

network for salient segmentation, IEEE Transactions on Image Processing 28 (5) (2018) 2518–2529.

[23] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, 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 40 (4) (2017) 834–848.

[24] P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, G. Cottrell, Understanding convolution for semantic segmentation, in: 2018 IEEE winter conference on applications of computer vision (WACV), IEEE, 2018, pp. 1451–1460.

[25] C. Peng, X. Zhang, G. Yu, G. Luo, J. Sun, Large kernel matters–improve semantic segmentation by global convolutional network, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 4353–4361.

[26] W. Wang, J. Shen, R. Yang, F. Porikli, Saliency-aware video object segmentation, IEEE transactions on pattern analysis and machine intelligence 40 (1) (2017) 20–33.

[27] W. Wang, J. Shen, L. Shao, Video salient object detection via fully convolutional networks, IEEE Transactions on Image Processing 27 (1) (2017) 38–49.

[28] F. Guo, W. Wang, J. Shen, L. Shao, J. Yang, D. Tao, Y. Y. Tang, Video saliency detection using object proposals, IEEE transactions on cybernetics 48 (11) (2017) 3159–3170.

[29] D.-P. Fan, W. Wang, M.-M. Cheng, J. Shen, Shifting more attention to video salient object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2019, pp. 8554–8564.

[30] X. Lu, W. Wang, C. Ma, J. Shen, L. Shao, F. Porikli, See more, know more: Unsupervised video object segmentation with co-attention siamese networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 3623–3632.

[31] Y. Niu, Y. Geng, X. Li, F. Liu, Leveraging stereopsis for saliency analysis, in: 2012 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2012, pp. 454–461.

[32] W. Wang, J. Shen, Y. Yu, K.-L. Ma, Stereoscopic thumbnail creation via efficient stereo saliency detection, IEEE transactions on visualization and computer graphics 23 (8) (2016) 2014–2027.

[33] C. Yang, L. Zhang, H. Lu, X. Ruan, 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.

[34] M.-M. Cheng, N. J. Mitra, X. Huang, P. H. Torr, S.-M. Hu, Global contrast based salient region detection, IEEE Transactions on Pattern Analysis and Machine Intelligence 37 (3) (2014) 569–582.

[35] C. Aytekin, A. Iosifidis, M. Gabbouj, Probabilistic saliency estimation, Pattern Recognition 74 (2018) 359–372.

[36] W. Wang, J. Shen, L. Shao, F. Porikli, Correspondence driven saliency transfer, IEEE Transactions on Image Processing 25 (11) (2016) 5025–5034.

[37] D. Shan, X. Zhang, C. Zhang, Visual saliency based on extended manifold ranking and third-order optimization refinement, Pattern Recognition Letters 116 (2018) 1–7.

[38] G. Li, 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.

[39] R. Zhao, W. Ouyang, H. Li, X. Wang, Saliency detection by multi-context deep learning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 1265–1274.

[40] G. Li, Y. Yu, Deep contrast learning for salient object detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 478–487.

[41] Q. Hou, M.-M. Cheng, X. Hu, A. Borji, Z. Tu, P. H. Torr, Deeply supervised salient object detection with short connections,

in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 3203–3212.

[42] T. Wang, L. Zhang, S. Wang, H. Lu, G. Yang, X. Ruan, 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.

[43] P. Zhang, W. Liu, Y. Lei, H. Lu, Hyperfusion-net: Hyper-densely reflective feature fusion for salient object detection, Pattern Recognition 93 (2019) 521–533.

[44] G. Li, Y. Xie, L. Lin, Y. Yu, Instance-level salient object segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 2386–2395.

[45] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556.

[46] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.

[47] B. Zoph, V. Vasudevan, J. Shlens, Q. V. Le, Learning transferable architectures for scalable image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 8697–8710.

[48] J. Hu, L. Shen, G. Sun, Squeeze-and-excitation networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 7132–7141.

[49] K. He, X. Zhang, S. Ren, J. Sun, Identity mappings in deep residual networks, in: European conference on computer vision, Springer, 2016, pp. 630–645.

[50] Q. Yan, L. Xu, J. Shi, J. Jia, Hierarchical saliency detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2013, pp. 1155–1162.

[51] L. Wang, H. Lu, Y. Wang, M. Feng, D. Wang, B. Yin, 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.

[52] Y. Li, X. Hou, C. Koch, J. M. Rehg, 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.

[53] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The pascal visual object classes (voc) challenge, International journal of computer vision 88 (2) (2010) 303–338.

[54] D.-P. Fan, M.-M. Cheng, J.-J. Liu, S.-H. Gao, Q. Hou, A. Borji, Salient objects in clutter: Bringing salient object detection to the foreground, in: Proceedings of the European conference on computer vision (ECCV), 2018, pp. 186–202.

[55] A. Borji, M.-M. Cheng, H. Jiang, J. Li, Salient object detection: A benchmark, IEEE transactions on image processing 24 (12) (2015) 5706–5722.

[56] R. Margolin, L. Zelnik-Manor, A. Tal, How to evaluate foreground maps?, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2014, pp. 248–255.

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

[58] R. Achanta, S. Hemami, F. Estrada, S. Süsstrunk, Frequency-tuned salient region detection, in: IEEE International Conference on Computer Vision and Pattern Recognition (CVPR 2009), no. CONF, 2009, pp. 1597–1604.

[59] F. Chollet, et al., Keras, `https://keras.io` (2015).

[60] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga,

S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, X. Zheng, TensorFlow: Large-scale machine learning on heterogeneous systems, software available from tensorflow.org (2015).

URL `https://www.tensorflow.org/`

[61] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., Imagenet large scale visual recognition challenge, International journal of computer vision 115 (3) (2015) 211–252.

[62] X. Zhang, T. Wang, J. Qi, H. Lu, G. Wang, Progressive attention guided recurrent network for salient object detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 714–722.

[63] N. Liu, J. Han, 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.

[64] R. Wu, M. Feng, W. Guan, D. Wang, H. Lu, E. Ding, A mutual learning method for salient object detection with intertwined multi-supervision, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 8150–8159.

[65] M. Feng, H. Lu, 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.

[66] L. Zhang, J. Zhang, Z. Lin, H. Lu, Y. He, Capsal: Leveraging captioning to boost semantics for salient object detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 6024–6033.

[67] Z. Wu, L. Su, 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.

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

[69] S. Cai, J. Huang, D. Zeng, X. Ding, J. Paisley, Menet: A metric expression network for salient object segmentation, in: Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18, International Joint Conferences on Artificial Intelligence Organization, 2018, pp. 598–605.