# DuAT: Dual-Aggregation Transformer Network for Medical Image Segmentation

Feilong Tang[1], Qiming Huang[1], Jinfeng Wang[1], Xianxu Hou[2], Jionglong Su[✉1], and Jingxin Liu[✉1]

[1]School of AI and Advanced Computing, Xi'an Jiaotong-Liverpool University, Suzhou, China
[2]School of Computer Science and Software Engineering, Shenzhen University, Shenzhen, China

## Abstract

*Transformer-based models have been widely demonstrated to be successful in computer vision tasks by modelling long-range dependencies and capturing global representations. However, they are often dominated by features of large patterns leading to the loss of local details (e.g., boundaries and small objects), which are critical in medical image segmentation. To alleviate this problem, we propose a Dual-Aggregation Transformer Network called DuAT, which is characterized by two innovative designs, namely, the Global-to-Local Spatial Aggregation (GLSA) and Selective Boundary Aggregation (SBA) modules. The GLSA has the ability to aggregate and represent both global and local spatial features, which are beneficial for locating large and small objects, respectively. The SBA module is used to aggregate the boundary characteristic from low-level features and semantic information from high-level features for better preserving boundary details and locating the re-calibration objects. Extensive experiments in six benchmark datasets demonstrate that our proposed model outperforms state-of-the-art methods in the segmentation of skin lesion images, and polyps in colonoscopy images. In addition, our approach is more robust than existing methods in various challenging situations such as small object segmentation and ambiguous object boundaries.*

## 1. Introduction

Medical image segmentation is a computer-aided automatic procedure for extracting the region of interest (RoI), e.g., tissues, lesions, and body organs. It can assist clinicians by making diagnostic and treatment processes more efficient and precise. For instance, colonoscopy is the gold standard for detecting colorectal lesions, and accurately locating early polyps is of great significance for clinical prevention of rectal cancer [17]. Likewise, melanoma skin cancer is one of the most rapidly increasing cancers worldwide.
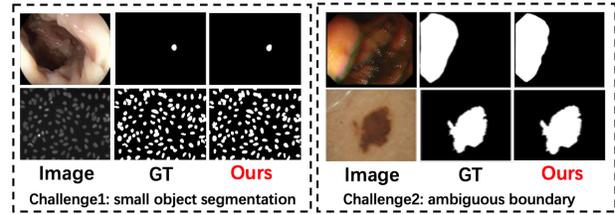


Figure 1. Main challenges of medical images. **Challenge1**: Small object segmentation are difficult to segment due to their low contrast and strong camouflage. **Challenge2**: The object boundary for medical image is ambiguous due to image acquisition influence.

Segmentation of skin lesions from dermoscopic images is a critical step in skin cancer diagnosis and treatment planning [33]. However, it is impractical to manually annotate these structures in clinical practice due to the tedious, time-consuming, and error-prone process. There is a growing need for automatic as well as accurate image segmentation. With the rapid development of deep learning, an increasing number of deep convolutional neural networks (DCNNs) [16, 18, 36, 54, 59] are proposed for medical image segmentation. The limitations of the receptive field in DCNNs make it difficult to capture the global representation. To solve this problem, dilated convolution [11, 51, 53] is proposed for semantic segmentation task. Furthermore, attention models [47, 21, 52] are developed to better capture long-range context information. These alternatives achieve promising results in semantic segmentation.

Transformer-based methods [1, 8, 9, 19, 40, 58, 46] have been proposed and achieved comparable performance to state-of-the-art results. Transformer is originally used for sequence-to-sequence predictive modelling in natural language processing (NLP) tasks with pure attention structure that is good at capturing long-range dependencies [42]. Vision Transformer-based methods [1, 8, 9, 19, 40, 58, 46] have also been proposed and demonstrate promising performance. ViT [15] constructs a vector sequence by dividing each image into fixed-size patches, subsequently applies a

multi-head self-attention (MHSA) and the Multilayer Perceptron (MLP) structure which demonstrate the advanced learning ability for long-distance feature dependence. The recent works [14, 50, 45] demonstrate that the pyramid structure is applicable in Transformers and more suitable for various downstream tasks. Unfortunately, capturing long-range dependencies destroys part of local features, which could result in overly smooth predictions for small objects and blurred boundaries between objects.

Therefore, building a model that retains both local and global features remains challenging. Li *et al.* [26] explore the local context for the aggregated long-range relationship to be more accurately distributed in local regions. Recently, some hybrid architectures of Transformer and CNN are proposed [10, 35, 58], which aim to combine the advantages of both models. Wang *et al.* [45] propose progressive locality decoder to emphasize local features and restrict attention dispersion. Chen *et al.* [10] propose TransUnet, which utilizes the underlying features of CNNs and subsequently uses the transformers to model global interactions to strengthen local features. However, feeding local information directly into the transformer cannot precisely handle local context relationships, resulting in the local information being overwhelmed by the dominant global context. Ultimately it leads to inferior results in the medical image segmentation of small objects.

Considering the problem of unclear object boundary information in semantic segmentation. Previous studies [27, 12, 14] explore fusing low-scale boundary information and high-scale semantic information to better preserve boundary details. Moreover, Zhang *et al.* [56] incorporate semantic information into low-level features while embedding more spatial information into high-level features to make up for the semantic and resolution gap between feature maps. Takikawa *et al.* [39] and Zhen *et al.* [57] design a boundary stream and couple the task of boundary and semantics modeling. Inspired by [56], we exploit the boundary information to selectively fuse it with the high-level semantic features, which will further enhance the semantic representations rather than simply combining them.

In this paper, we propose a novel pyramid transformer for medical image segmentation, referred to as the Dual-Aggregation Transformer Network (DuAT), consisting of the Global-to-Local Spatial Aggregation (*GLSA*) to combine local and global features, as well as the Selective Boundary Aggregation (*SBA*) module to enhance the boundary information and locating the re-calibration objects. We believe that global spatial features help locate large objects, and local spatial features are crucial for identifying small ones. Finally, boundary information is aggregated to fine-tune object boundaries and re-calibrate coarse predictions, Specifically, the *GLSA* module is responsible for extracting and fusing local and global spa-

tial features from the backbone. We separate the channels, one for global representation extracted by Global context (GC) block [7], and the other for local information extracted by multiple depth-wise convolutions. The *SBA* module aims to simulate the biological visual perception process, distinguishing objects from background. Specifically, it incorporates shallow- and deep-level features to establish the relationship between body areas and boundary, enhancing the boundaries characteristics. Our experimental results demonstrate three advantages of our model: more lightweight, better learning ability and improved generalization capability than previous state-of-the-art approaches.

In summary, the main contribution of this paper is threefold.

- We propose a novel framework, named Dual-Aggregation Transformer Network (**DuAT**), for medical images segmentation, which adapts the pyramid vision transformer as encoder to extract more robust features than the existing CNN-based methods.

- We design dual aggregation modules, Global-to-Local Spatial Aggregation (**GLSA**) module and Selective Boundary Aggregation (**SBA**) module. Their purpose is to solve two challenges that are intuitively demonstrated in Figure 1. Specifically, the GLSA module simultaneously extract local spatial detail information and global spatial semantic information, which reduces incorrect information in the high-level features. The SBA module aims to fine-tune object boundaries, which can well address the "ambiguous" problem of boundaries.

- Extensive experiments on five polyp datasets (ETIS [43], CVC-ClinicDB [37], CVC-ColonDB [3], EndoScene-CVC300 [38], Kvasir [22]), skin lesion dataset (ISIC-2018 [13]) and 2018 Data Science Bowl [6] demonstrate that the proposed DuAT methods advances the state-of-the-art (SOTA) performance.

## 2. Related Work

### 2.1. Vision Transformer

Transformer has dominated the field of NLP with its MHSA layer to capture the pure attention structure of long-range dependencies. Different from convolutional layer, MHSA layer has dynamic weight and global receptive field, which makes it more flexible and effective. Dosovitskiy *et al.* propose a vision transformer (ViT) [15], which is an end-to-end model using the Transformer structure for image recognition task. Specifically, it divides an image into fixed-size patches, which are sequentially fed to multiple Transformer encoder blocks to model the patches. In addition, previous work has proved that the pyramid structure in
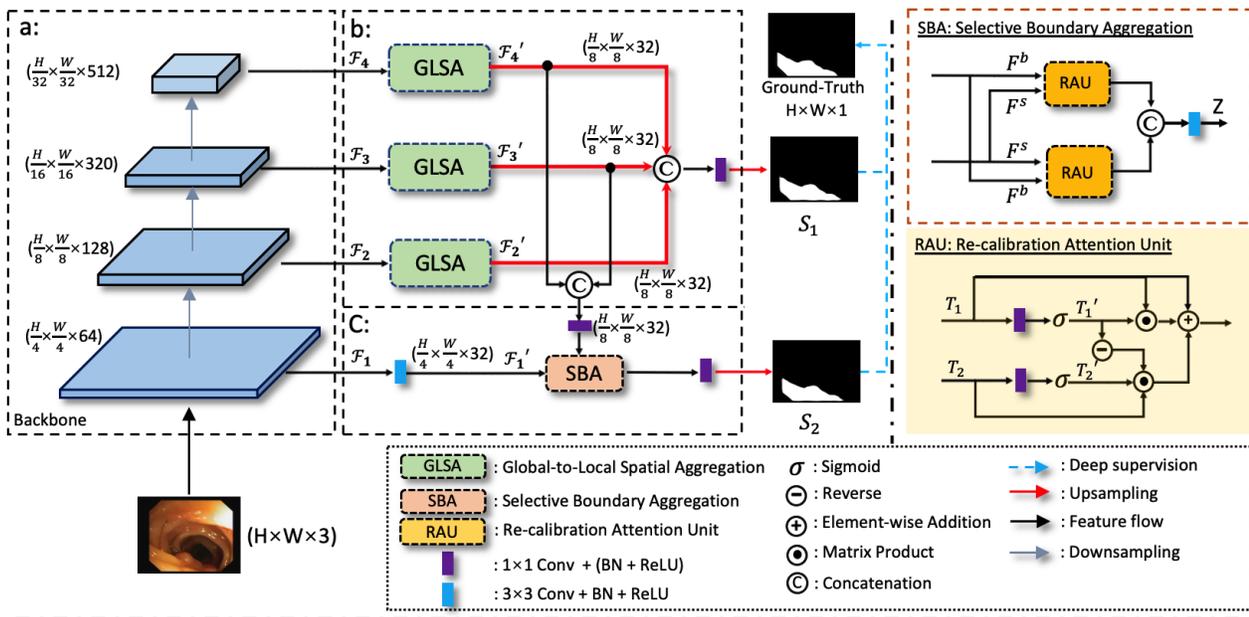
Figure 2. The overall architecture of Dual-Aggregation Transformer Network (DuAT). The entire model is divided into three parts: (a) pyramid vision transformer (PVT) as backbone; (b) pyramid Global-to-Local Spatial Aggregation (GLSA) Module; (c) Selective Boundary Aggregation (SBA) module and it shown on the red box.

convolutional networks is also suitable for Transformer and various downstream tasks, such as Swin Transformer [28], PVT [46], Segformer [50], etc. PVT requires less computation than ViT and adopts the classical Semantic-FPN to deploy the task of semantic segmentation.

In medical image segmentation, TransUNet [10] demonstrates that Transformer can be used as powerful encoders for medical image segmentation. TransFuse [55] is proposed to improve efficiency for global context modeling by fusing transformers and CNNs. Furthermore, to train the model effectively on medical images, Polyp-PVT [14] introduces Similarity Aggregation Module based on graph convolution domain [31]. Inspired by these approaches, we propose a new transformer-based medical segmentation framework, which can accurately locate small objects.

### 2.2. Image Boundary Segmentation

Recently, learning additional boundary information has shown superior performance in many image segmentation tasks. In the early research on FCN-based semantic segmentation, Bertasius *et al.* [4] and Chen *et al.* [11] use boundaries for post-starting to refine the result at the end of the network. Recently, several approaches explicitly model boundary detection as an independent sub-task in parallel with semantic segmentation for sharper result. Ma *et al.* [32] explicitly exploit the boundary information for context aggregation to further enhance the semantic representation of the model. Li *et al.* [25] point out that the object boundary and body parts correspond to the high-frequency and low-frequency information of an image, respectively, based on which they decouple the body and edge with diverse supervisions. Ji *et al.* [24] fuse the low-level edge-aware features and constraint it with the explicit edge supervision. Different from the above works, we propose a novel aggregation method to achieve more accurate localisation and boundary delineation of objects.

## 3. Method

As given in Figure 2, our DuAT model consists of the following:- a pyramid vision transformer (PVT) encoder, a *SBA* module, and *GLSA* module.

### 3.1. Transformer Encoder

Some recent studies [5, 50] report that vision transformers [15, 46] have stronger performance and robustness to input disturbances (e.g., noise) than CNNs. Inspired by this, we use the Transformer based on pyramid structure as the encoder. Specifically, the pyramid vision transformer (PVT) [46] is utilized as the encoder module for multi-level feature maps $\{\mathcal{F}_i | i \in (1, 2, 3, 4)\}$ extraction. Among these feature maps, $\mathcal{F}_1$ gives detailed boundary information of target, and $\mathcal{F}_2$, $\mathcal{F}_3$ and $\mathcal{F}_4$ provide high-level features.
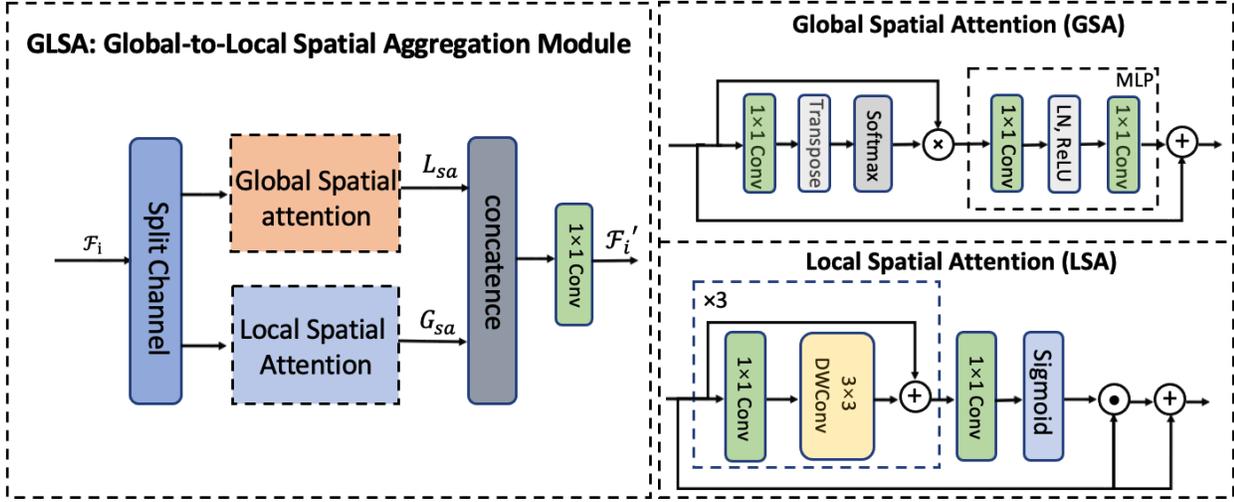
Figure 3. Overview of the Global-to-Local Spatial Aggregation Module *GLSA*, it is composed of global spatial attention (*GSA*) and local spatial attention *(LSA)*.

## 3.2. Selective Boundary Aggregation

As observed in [56, 25], shallow- and deep-layer features complement each other. The shallow layer has less semantics but is rich in details, with more distinct boundaries and less distortion. Furthermore, the deep level contains a rich semantic information. Therefore, directly fusing low-level features with high-level ones may result in redundancy and inconsistency. To address this, we propose the SBA module, which selectively aggregate the boundary information and semantic information to depict more fine-grained contour of objects and the location of re-calibrate objects.

Different from previous fusion methods, we design a novel Re-calibration attention unit (RAU) block that adaptively picks up mutual representations from tow inputs $(F^s, F^b)$ before fusion. As given in Figure 2, the shallow - and deep-level information is fed into the two RAU blocks by different ways to make up for the missing spatial boundary information of the high-level semantic features and the missing semantic information of low-level features. Finally, the outputs of two RAU blocks are concatenated after a $3 \times 3$ convolution. This aggregation strategy realizes the robust combination of different features and refines the rough features. The RAU block function $PAU(\cdot, \cdot)$ process can be expressed as:

$$T_1' = W_\theta(T_1), T_2' = W_\phi(T_2) \tag{1}$$
$$PAU(T_1, T_2) = T_1' \odot T_1 + T_2' \odot T_2 \odot (\ominus(T_1')) + T_1, \tag{2}$$

where $T_1, T_2$ are the input features, two linear mapping and sigmoid functions $W_\theta(\cdot), W_\phi(\cdot)$ are applied to the input features to reduce the channel dimension to 32 and obtain fea-

ture maps $T_1'$ and $T_2'$. $\odot$ is Point-wise multiplication. $\ominus(\cdot)$ is the reverse operation by subtracting the feature $T_1'$, refining the imprecise and coarse estimation into an accurate and complete prediction map [16]. We take a convolutional operation with a kernel size of $1 \times 1$ as the linear mapping process. As a result, the process of SBA can be formulated as:

$$Z = C_{3\times3}(\text{Concat}(PAU(F^s, F^b), PAU(F^b, F^s))), \tag{3}$$

where $C_{3\times3}(\cdot)$ is a $3 \times 3$ convolution with a batch normalization and a ReLU activation layer. $F^s \in \mathbb{R}^{\frac{H}{8} \times \frac{W}{8} \times 32}$ contains deep-level semantic information after fusing the third and fourth layers from the encoder, $F^b \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 32}$ is the first layer with rich boundary details from the backbone. $Concat(\cdot)$ is the concatenation operation along the channel dimension. $Z \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 32}$ is the output of the SBA module.

## 3.3. Global-to-Local Spatial Aggregation

The attention mechanism strengthens the information related to the optimization goal and suppresses irrelevant information. In order to capture both global and local spatial features, we propose the GLSA module, which fuses the results of two separate local and global attention units. As demonstrated in Figure 3, this dual-stream design effectively preserves both local and non-local modeling capabilities. Moreover, we use separating channels to balance the accuracy and computational resources. Specifically, the feature map $\{\mathcal{F}_i | i \in (2, 3, 4)\}$ with 64 channels is split evenly into two feature map groups $\mathcal{F}_i^1, \mathcal{F}_i^2 (i \in (2, 3, 4))$ and separately fed into Global Spatial attention (GSA) module and
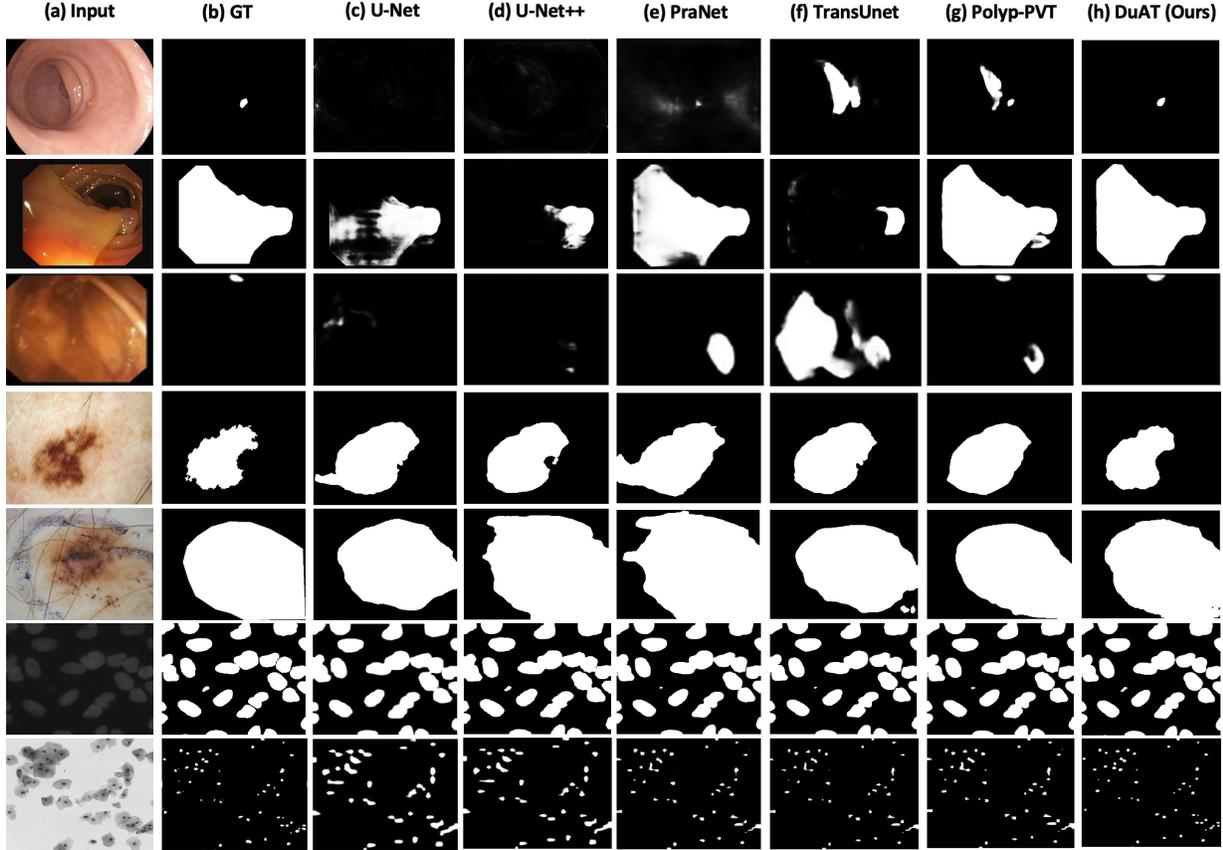
Figure 4. Qualitative results of different methods. (a) Inputs images, (b) GT, which stands for the ground truths, (h) semantic segmentation maps produced by our method, (c) U-Net [36], (d) U-Net++ [59], (e) PraNet [16], (f) TransUnet [10], (g) Polyp-PVT [14].

Local Spatial attention (LSA) module. The outputs of those two attention units are finally concatenated following by a $1 \times 1$ convolution layer. We formulate such a process as

$$\mathcal{F}_i^1, \mathcal{F}_i^2 = \text{Split}(\mathcal{F}_i) \qquad (4)$$

$$\mathcal{F}_i^{'} = C_{1\times1}(\text{Concat}(G_{sa}(\mathcal{F}_i^1), L_{sa}(\mathcal{F}_i^2))). \qquad (5)$$

where $G_{sa}$ denotes the global spatial attention and $L_{sa}$ denotes the local spatial attention. $\mathcal{F}_i^{'} \in \mathbb{R}^{\frac{H}{8} \times \frac{W}{8} \times 32}$ is the output features. We will introduce LSA and GSA module in detail in the following.

(1) GSA module: The GSA emphasizes the long-range relationship of each pixel in the spatial space and can be used as a supplement to local spatial attention. Many efforts [2], [7] claim that the long-range interaction can make the feature more powerful. Inspired by the manners of extracting long-range interaction in [2], we simply generate global spatial attention map ($G_{sa} \in \mathbb{R}^{\frac{H}{8} \times \frac{W}{8} \times 32}$) and $\mathcal{F}_i^1$ as

input as following:

$$Att_G(\mathcal{F}_i^1) = Softmax(Transpose(C_{1\times1}(\mathcal{F}_i^1))), \qquad (6)$$

$$G_{sa}(\mathcal{F}_i^1) = MLP(Att_G(\mathcal{F}_i^1) \otimes \mathcal{F}_i^1) + \mathcal{F}_i^1. \qquad (7)$$

where $Att_G(\cdot)$ is the attention operation, $C_{1\times1}$ means $1 \times 1$ convolution. $\otimes$ denotes matrix multiplication. $MLP(\cdot)$ consists of two fully-connection layers with a ReLU non-linearity and normalization layer. The first layer of MLP transforms its input to a higher-dimensional space which the expansion ratio is two, while the second layer restores the dimension to be the same as the input.

(2) LSA module: The LSA module extracts the local features of the region of interest effectively in the spatial dimension of the given feature map, such as small objects. In short, we compute local spatial attention response ($L_{sa} \in \mathbb{R}^{\frac{H}{8} \times \frac{W}{8} \times 32}$) and $\mathcal{F}_i^2$ as input as follow:

$$Att_L(\mathcal{F}_i^2) = \sigma(C_{1\times1}(\mathcal{F}_c(\mathcal{F}_i^2)) + \mathcal{F}_i^2)), \qquad (8)$$

$$L_{sa} = Att_L(\mathcal{F}_i^2) \odot \mathcal{F}_i^2 + \mathcal{F}_i^2. \qquad (9)$$

Table 1. Quantitative comparison of different methods on Kvasir, ClinicDB, ISIC-2018 and 2018-DSB datasets (seen datasets) to validate our model's learning ability. ↑ denotes higher the better and ↓ denotes lower the better.

| Methods | Kvasir | | | ClinicDB | | | ISIC-2018 | | | 2018-DSB | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | mDice↑ | mIou↑ | MAE↓ | mDice↑ | mIou↑ | MAE↓ | mDice↑ | mIou↑ | MAE↓ | mDice↑ | mIou↑ | MAE↓ |
| U-Net [36] | 0.818 | 0.746 | 0.055 | 0.823 | 0.755 | 0.019 | 0.855 | 0.785 | 0.045 | 0.908 | 0.831 | 0.040 |
| UNet++ [23] | 0.821 | 0.743 | 0.048 | 0.794 | 0.729 | 0.022 | 0.809 | 0.729 | 0.041 | 0.911 | 0.837 | 0.039 |
| PraNet [16] | 0.898 | 0.840 | 0.030 | 0.899 | 0.849 | 0.009 | 0.875 | 0.787 | 0.037 | 0.912 | 0.838 | 0.036 |
| CaraNet [30] | 0.918 | 0.865 | 0.023 | 0.936 | 0.887 | 0.007 | 0.870 | 0.782 | 0.038 | 0.910 | 0.835 | 0.037 |
| TransUNet [10] | 0.913 | 0.857 | 0.028 | 0.935 | 0.887 | 0.008 | 0.880 | 0.809 | 0.036 | 0.915 | 0.845 | 0.033 |
| TransFuse [55] | 0.920 | 0.870 | 0.023 | 0.942 | 0.897 | 0.007 | 0.901 | 0.840 | 0.035 | 0.916 | 0.855 | 0.033 |
| UCTransNet [44] | 0.918 | 0.860 | 0.023 | 0.933 | 0.860 | 0.008 | 0.905 | 0.83 | 0.035 | 0.911 | 0.835 | 0.035 |
| Polyp-PVT [14] | 0.917 | 0.864 | 0.023 | 0.937 | 0.889 | 0.006 | 0.913 | 0.852 | 0.032 | 0.917 | 0.859 | 0.030 |
| **DuAT (Ours)** | **0.924** | **0.876** | **0.023** | **0.948** | **0.906** | **0.006** | **0.923** | **0.867** | **0.029** | **0.926** | **0.870** | **0.027** |

Table 2. Quantitative comparison of different methods on ColonDB, ETIS and EndoScene datasets (unseen datasets) to validate our model's generalization capability. ↑ denotes higher the better and ↓ denotes lower the better.

| Methods | ColonDB | | | ETIS | | | EndoScene | | |
|---|---|---|---|---|---|---|---|---|---|
| | mDice↑ | mIou↑ | MAE↓ | mDice↑ | mIou↑ | MAE↓ | mDice↑ | mIou↑ | MAE↓ |
| U-Net [36] | 0.512 | 0.444 | 0.061 | 0.398 | 0.335 | 0.036 | 0.710 | 0.627 | 0.022 |
| UNet++ [23] | 0.483 | 0.410 | 0.064 | 0.401 | 0.344 | 0.035 | 0.707 | 0.624 | 0.018 |
| PraNet [16] | 0.712 | 0.640 | 0.043 | 0.628 | 0.567 | 0.031 | 0.851 | 0.797 | 0.010 |
| CaraNet [30] | 0.773 | 0.689 | 0.042 | 0.747 | 0.672 | 0.017 | 0.903 | 0.838 | 0.007 |
| TransUNet [10] | 0.781 | 0.699 | 0.036 | 0.731 | 0.824 | 0.021 | 0.893 | 0.660 | 0.009 |
| TransFuse [55] | 0.781 | 0.706 | 0.035 | 0.737 | 0.826 | 0.020 | 0.894 | 0.654 | 0.009 |
| SSformer [45] | 0.772 | 0.697 | 0.036 | 0.767 | 0.698 | 0.016 | 0.887 | 0.821 | 0.007 |
| Polyp-PVT [14] | 0.808 | 0.727 | 0.031 | 0.787 | 0.706 | 0.013 | 0.900 | 0.833 | 0.007 |
| **DuAT (Ours)** | **0.819** | **0.737** | **0.026** | **0.822** | **0.746** | **0.013** | **0.901** | **0.840** | **0.005** |

where $\mathcal{F}_c(\cdot)$ denotes cascading three $1 \times 1$ convolution layers and $3 \times 3$ depth-wise convolution layers. The number of channels is adjusted to 32 in the $\mathcal{F}_c$. $Att_L(\cdot)$ is the local attention operation, $\sigma(\cdot)$ is the sigmoid function, $\odot$ is point-wise multiplication. This structural design can efficiently aggregate local spatial information using fewer parameters.

## 3.4. Loss function

[34, 49] report that combining multiple loss functions with adaptive weights at different levels can improve the performance of the network with better convergence speed. Therefore, we use binary cross-entropy loss ($\mathcal{L}_{BCE}^{\omega}(\cdot)$) and the weighted IoU loss ( $\mathcal{L}_{Iou}^{\omega}(\cdot)$) for supervision. Our loss function is formulated in Eq.10, where $S$ is the two side-outputs $(i,e., S_1, S_2)$ and $G$ is the ground truth, respectively. $\lambda_1$ and $\lambda_2$ are the weighting coefficients.

$$\mathcal{L}(S, G) = \lambda_1 \mathcal{L}_{IoU}^{\omega}(S, G) + \lambda_2 \mathcal{L}_{BCE}^{\omega}(S, G) \quad (10)$$

Therefore, the total loss $\mathcal{L}_{total}$ for the proposed DuAT can be formulated as:

$$\mathcal{L}_{total} = \mathcal{L}(S_1, G) + \mathcal{L}(S_2, G). \quad (11)$$

## 4. Experiments

### 4.1. Datasets

In the experiment, we evaluate our proposed model on three different kinds of medical image sets: colonoscopy (Colon) images, dermoscopic (Derm) images, and microscopy (Micro) images, so as to assess the learning ability and generalization capability of our model. The detailed statistics of each dataset are shown in Table 3.

**Colonoscopy polyp images**: Experiments are conducted on five polyp segmentation datasets (ETIS [43], CVC-ClinicDB (ClinicDB) [37], CVC-ColonDB (ColonDB) [3], EndoScene-CVC300 (EndoScene) [38], Kvasir-SEG (Kvasir) [22]). We follow the same training/testing protocols in [14, 16], i.e., the images from the Kvasir and ClinicDB are randomly split into 80% for training, 10% for validation, and 10% for testing (seen data). And test on the out-of-distribution datasets which are ColonDB with 380 images, EndoScene with 60 images and ETIS with 196 images (unseen data). Since the resolutions of images are not uniform, we resize them to 352×352 resolution.

**ISIC-2018 Dataset**: The dataset comes from ISIC-2018 challenge [13] [41] and is useful for skin lesion analysis. It includes 2596 images and the corresponding annotations,

Table 3. Statistics on polyp, ISIC-2018 and 2018-DSB datasets.

| Dataset | Imaging | Images | Shape | Train | Valid | Test |
|---|---|---|---|---|---|---|
| Kvasir | Colon | 1000 | Variable | 800 | 100 | 100 |
| ClinicDB | Colon | 612 | $384 \times 288$ | 488 | 62 | 62 |
| ColonDB | Colon | 380 | $574 \times 500$ | - | - | 380 |
| ETIS | Colon | 196 | $1225 \times 966$ | - | - | 196 |
| EndoScene | Colon | 60 | $574 \times 500$ | - | - | 60 |
| ISIC-2018 | Derm | 2596 | $512 \times 384$ | 2078 | 259 | 259 |
| 2018-DSB | Micro | 670 | $256 \times 256$ | 536 | 67 | 67 |

Table 4. Number of Parameters and FLOPs on a $352 \times 352$ input image. Note that "N/A" means the result is not available.

| Method | Type | Params(M) | FLOPs(G) |
|---|---|---|---|
| U-Net [36] | CNN | 43.93 | 43.47 |
| UNet++ [23] | CNN | 9.04 | 64.03 |
| PraNet [16] | CNN | 30.50 | 13.01 |
| CaraNet [30] | CNN | 44.59 | 21.65 |
| TransUNet [10] | Transformer | 93.19 | 60.84 |
| SSformer [45] | Transformer | 26.70 | 28.07 |
| TransFuse [55] | Transformer | N/A | N/A |
| Polyp-PVT [14] | Transformer | 25.08 | 10.00 |
| **Ours** | Transformer | 24.92 | 9.88 |

which are resized to $512 \times 384$ resolution. The images are randomly split into 80% for training, 10% for validation, and 10% for testing.

**2018 Data Science Bowl (2018-DSB)**: The dataset comes from 2018 Data Science Bowl challenge [6] and is used to find the nuclei in divergent images, including 670 images and the corresponding annotations, which are resized to $256 \times 256$ resolution. The images are randomly split into 80% for training, 10% for validation, and 10% for testing.

### 4.2. Evaluation Metrics and Implementation Details

**Evaluation Metrics.** We employ three widely-used metrics i.e., mean Dice (mDice), mean IoU (mIoU) and mean absolute error (MAE) to evaluate the model performances. Mean Dice and IoU are the most commonly used metrics and mainly emphasise the internal consistency of segmentation results. MAE is used to evaluate the pixel-level accuracy representing the average absolute error between the prediction and true values.

**Implementation Details.** We use rotation and horizontal flip for data augmentation. Considering the differences in the sizes and color of each polyp image, we adopt a multi-scale training [16, 20] and the color exchange [48]. The network is trained end-to-end by AdamW [29] optimizer. The learning rate is set to 1e-4 and the weight decay is adjusted to 1e-4 too. The batch size is set at 16. We use PyTorch framework for implementation with an NVIDIA RTX 3090 GPU. We will provide the source code after the paper is published.

### 4.3. Results

**Learning Ability**. We first evaluate our proposed DuAT model for its segmentation performance on seen datasets. As summarized in Table 1, our model is compared to six recently published models: U-Net [36], UNet++ [23], PraNet [16], CaraNet [30], TransUNet [10], TransFuse [55], TransFuse [55] and Polyp-PVT [14] . It can be observed that our DuAT model outperforms all other models, and achieving 0.924 and 0.948 mean Dice scores on Kvasir and ClinicDB segmentation respectively. For ISIC-2018 dataset, our DuAT model achieves a 1.0% improvement in terms of mDice and 1.5% of mIoU over SOTA method. For 2018-DSB, DuAT achieves a mIoU of 0.87, mDice of 0.926 and

0.027 of MAE, which are 1.1%, 1.0%, 0.03% higher than the best performing Polyp-PVT. These results demonstrate that our model can effectively segment polyps.

**Generalization Capabilities.** We further evaluate the generalisation capability of our model on unseen datasets (ETIS, ColonDB, EndoScene). These three datasets have their own specific challenges and properties. For example, ColonDB is a small-scale database that contains 380 images from 15 short colonoscopy sequences. ETIS consists of 196 polyp images for early diagnosis of colorectal cancer. EndoScene is a re-annotated branch with associated polyp and background (mucosa and lumen). As seen in Table 2, our model outperforms the existing medical segmentation baselines on all unseen datasets for all metrics. Moreover, our DuAT is able to achieve an average dice of 82.2 % on the most challenging ETIS dataset, 3.5% higher than Polyp-PVT.

**Visual Results.** We also demonstrate qualitative the performance of our model on five benchmarks, as given in Figure 4. On ETIS (the first and second row), DuAT is able to accurately capture the target object's boundary and detect a small polyp while other methods fail to detect. On ISIC-2018 (third row), all methods are able to segment the lesion skin, but our method show the most similar results compared to the ground truth. On 2018-DSB (the fourth row), we can observe that our DuAT is better able to capture the presence of nuclei and obtain better segmentation predictions. More qualitative results can be found in the supplementary material.

**Computational Efficiency.** Table 4 presents the number of parameters and floating-point operations for different methods. As our proposed DuAT and Polyp-PVT [14] adopt the same backbone, they have similar model size (Params). DuAT uses 24.92M of parameter and 9.88G of FLOPs, which is more lightweight and compact than CNN-based neural network and Transformer-based methods.

**Small Object Segmentation Analysis.** To demonstrate the detection ability of our model for small objects, the ra-
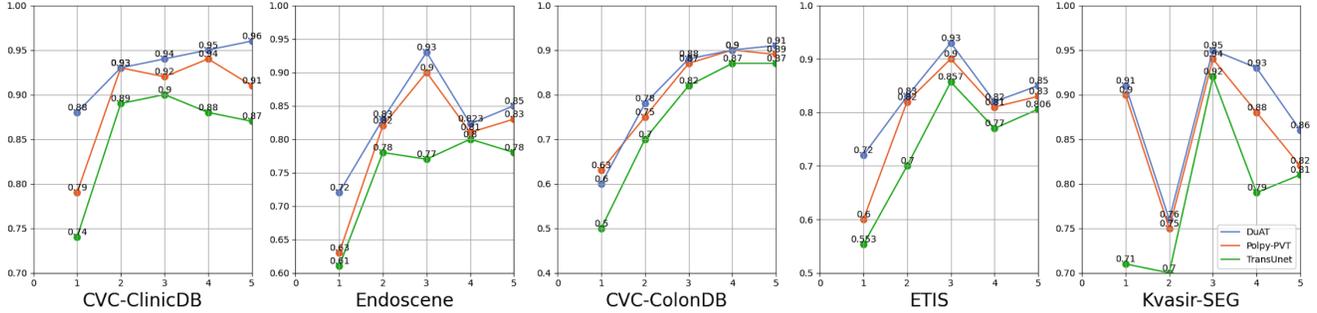
Figure 5. Performance vs. Size on the five polyp datasets. The x-axis is the proportion size % of polyp and y-axis is the averaged mDice coefficient. Blue is for our DuAT, orange is for the Polyp-PVT, green is for the TransUnet.

Table 5. Ablation study for DuAT on the Kvasir, ETIS and ISIC-2018 datasets. ↑ denotes higher the better and ↓ denotes lower the better.

| Methods | Kvasir-SEG (seen) | | | ETIS (unseen) | | | ISIC-2018 (seen) | | |
|---|---|---|---|---|---|---|---|---|---|
| | mDice↑ | mIou↑ | MAE↓ | mDice↑ | mIou↑ | MAE↓ | mDice↑ | mIou↑ | MAE↓ |
| Baseline | 0.910 | 0.856 | 0.030 | 0.759 | 0.668 | 0.035 | 0.877 | 0.783 | 0.040 |
| + GSA | 0.912 | 0.860 | 0.029 | 0.772 | 0.675 | 0.030 | 0.887 | 0.803 | 0.038 |
| + LSA | 0.916 | 0.863 | 0.028 | 0.785 | 0.690 | 0.027 | 0.900 | 0.839 | 0.035 |
| + GSA + LSA (Serial) | 0.914 | 0.863 | 0.028 | 0.786 | 0.695 | 0.025 | 0.909 | 0.845 | 0.034 |
| + LSA + GSA (Serial) | 0.910 | 0.860 | 0.029 | 0.799 | 0.713 | 0.021 | 0.910 | 0.852 | 0.033 |
| +GLSA | 0.917 | 0.864 | 0.025 | 0.814 | 0.723 | 0.016 | 0.916 | 0.816 | 0.031 |
| w/o SBA | 0.917 | 0.864 | 0.025 | 0.814 | 0.723 | 0.016 | 0.916 | 0.816 | 0.031 |
| w/o GLSA | 0.915 | 0.863 | 0.026 | 0.790 | 0.696 | 0.023 | 0.901 | 0.800 | 0.033 |
| **SBA + GLSA (Ours)** | **0.924** | **0.876** | **0.023** | **0.822** | **0.746** | **0.013** | **0.923** | **0.867** | **0.029** |

tio of the number of pixels in the object to the number of pixels in the entire image is used to account for the size of the object. We then evaluate the performance of the segmentation model based on the size of the object. We set the area with a proportion less than 5%. For the segmentation model, we first obtain the mean Dice coefficient of the five polyp datasets. Similar to computing the histogram, we calculate the average of mean Dice of test data whose size values fall into each interval. For the small object segmentation analysis, we compare our DuAT with Polyp-PVT and TransUnet, and the results are given in Figure 5. The overall accuracy of DuAT is higher than TransUnet [10] and Polyp-PVT [46] on samples with small size polyps.

## 4.4. Ablation Study

We further conduct ablation study to demonstrate the necessity and effectiveness of each component of our proposed model on three datasets, and we choose mDice, mIoU and MAE for evaluation.

**Effectiveness of SBA and GLSA.** We conduct an experiment to evaluate DuAT without SBA module "(w/o SBA)". The performance without the SBA drops sharply on all three datasets are shown in Table 5. In particular, the mDice is reduced from 0.822 to 0.814 on ETIS. Moreover, we further investigate the contribution of the Global-to-Local Spatial Aggregation by removing it from the overall DuAT and replacing it with convolution operation with a kernel size of 3,

which is denoted as "(w/o GLSA)". The performance of the complete DuAT shows an improvement of 2.2 % and 6.7% in terms of mDice and mIoU respectively on ISIC-2018. After using the two modules (SBA + GLSA), the model's performance is improved again. These results demonstrate that these modules enable our model to distinguish polyp and lesion tissues effectively.

The visual results are given in Figure 6. We observe that the SBA module facilitates the fine-grained of ambiguous boundaries and the GLSA module greatly improves the accuracies of small object detection and target object location.
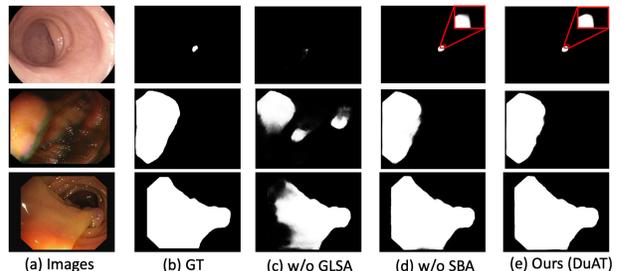


Figure 6. The effectiveness of each component.

**Arrangements of GSA and LSA.** *GSA* and *LSA* represent global spatial attention module and local spatial attention module respectively. We further study the effec-

tiveness and different arrangements of *GSA* and *LSA*. The results tested on Kavsir, ETIS, and ISIC-2018 datasets are shown in Table 5 (the second and sixth row), and all the methods are using the same backbone PVTv2 [46]. *GSA + LSA(Serial)* means first performing *GSA* then *LSA*, while *LSA + GSA(Serial)* is the opposite. Overall, all improve the baseline and our GLSA group achieves more accuracy and reliable results. The GLSA module outperforms the *GSA*, *LSA*, *GSA + LSA (Serial)*, *LSA + GSA (Serial)* by 4.2%, 2.9%, 2.8%, 1.5% in term of mean *Dice* on the ETIS dataset.

## 5. Conclusions

In this work, we propose DuAT to address the issues related to medical image segmentation. Two components, the *GLSA* and *SBA* are proposed. Specifically, the *GLSA* module extracts the global and local spatial features from the encoder and is beneficial for locating the large and small objects. The *SBA* module alleviates the unclear boundary of high-level features and further improves its performance. As a result, DuAT can achieve strong learning, generalization ability, and lightweight segmentation efficiency. Both qualitative and quantitative results demonstrate the superiority of our DuAT over other competing methods. We hope that this research will inspire more ideas to solve the medical image segmentation task and we will extend the proposed model to tackle 3D medical image segmentation task in the future work.

## Acknowledgments

## References

[1] Josh Beal, Eric Kim, Eric Tzeng, Dong Huk Park, Andrew Zhai, and Dmitry Kislyuk. Toward transformer-based object detection. *arXiv preprint arXiv:2012.09958*, 2020.

[2] Irwan Bello, Barret Zoph, Ashish Vaswani, Jonathon Shlens, and Quoc V Le. Attention augmented convolutional networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3286–3295, 2019.

[3] Jorge Bernal, F Javier Sánchez, Gloria Fernández-Esparrach, Debora Gil, Cristina Rodríguez, and Fernando Vilariño. Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians. *Computerized Medical Imaging and Graphics*, 43:99–111, 2015.

[4] Gedas Bertasius, Jianbo Shi, and Lorenzo Torresani. Semantic segmentation with boundary neural fields. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3602–3610, 2016.

[5] Srinadh Bhojanapalli, Ayan Chakrabarti, Daniel Glasner, Daliang Li, Thomas Unterthiner, and Andreas Veit. Understanding robustness of transformers for image classification.

[6] Juan C Caicedo, Allen Goodman, Kyle W Karhohs, Beth A Cimini, Jeanelle Ackerman, Marzieh Haghighi, CherKeng Heng, Tim Becker, Minh Doan, Claire McQuin, et al. Nucleus segmentation across imaging experiments: the 2018 data science bowl. *Nature methods*, 16(12):1247–1253, 2019.

[7] Yue Cao, Jiarui Xu, Stephen Lin, Fangyun Wei, and Han Hu. Gcnet: Non-local networks meet squeeze-excitation networks and beyond. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019.

[8] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020.

[9] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12299–12310, 2021.

[10] Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L Yuille, and Yuyin Zhou. Transunet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306*, 2021.

[11] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan 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):834–848, 2017.

[12] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018.

[13] Noel Codella, Veronica Rotemberg, Philipp Tschandl, M Emre Celebi, Stephen Dusza, David Gutman, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, et al. Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1902.03368*, 2019.

[14] Bo Dong, Wenhai Wang, Deng-Ping Fan, Jinpeng Li, Huazhu Fu, and Ling Shao. Polyp-pvt: Polyp segmentation with pyramid vision transformers. *arXiv preprint arXiv:2108.06932*, 2021.

[15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

[16] Deng-Ping Fan, Ge-Peng Ji, Tao Zhou, Geng Chen, Huazhu Fu, Jianbing Shen, and Ling Shao. Pranet: Parallel reverse attention network for polyp segmentation. In *International*

*conference on medical image computing and computer-assisted intervention*, pages 263–273. Springer, 2020.

[17] Pasqualino Favoriti, Gabriele Carbone, Marco Greco, Felice Pirozzi, Raffaele Emmanuele Maria Pirozzi, and Francesco Corcione. Worldwide burden of colorectal cancer: a review. *Updates in surgery*, 68(1):7–11, 2016.

[18] Huazhu Fu, Jun Cheng, Yanwu Xu, Damon Wing Kee Wong, Jiang Liu, and Xiaochun Cao. Joint optic disc and cup segmentation based on multi-label deep network and polar transformation. *IEEE transactions on medical imaging*, 37(7):1597–1605, 2018.

[19] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing Xu, and Yunhe Wang. Transformer in transformer. *Advances in Neural Information Processing Systems*, 34, 2021.

[20] Chien-Hsiang Huang, Hung-Yu Wu, and Youn-Long Lin. Hardnet-mseg: a simple encoder-decoder polyp segmentation neural network that achieves over 0.9 mean dice and 86 fps. *arXiv preprint arXiv:2101.07172*, 2021.

[21] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Ccnet: Criss-cross attention for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 603–612, 2019.

[22] Debesh Jha, Pia H Smedsrud, Michael A Riegler, Pål Halvorsen, Thomas de Lange, Dag Johansen, and Håvard D Johansen. Kvasir-seg: A segmented polyp dataset. In *International Conference on Multimedia Modeling*, pages 451–462. Springer, 2020.

[23] Debesh Jha, Pia H Smedsrud, Michael A Riegler, Dag Johansen, Thomas De Lange, Pål Halvorsen, and Håvard D Johansen. Resunet++: An advanced architecture for medical image segmentation. In *2019 IEEE International Symposium on Multimedia (ISM)*, pages 225–2255. IEEE, 2019.

[24] Ge-Peng Ji, Lei Zhu, Mingchen Zhuge, and Keren Fu. Fast camouflaged object detection via edge-based reversible recalibration network. *Pattern Recognition*, 123:108414, 2022.

[25] Xiangtai Li, Xia Li, Li Zhang, Guangliang Cheng, Jianping Shi, Zhouchen Lin, Shaohua Tan, and Yunhai Tong. Improving semantic segmentation via decoupled body and edge supervision. In *European Conference on Computer Vision*, pages 435–452. Springer, 2020.

[26] Xiangtai Li, Li Zhang, Ansheng You, Maoke Yang, Kuiyuan Yang, and Yunhai Tong. Global aggregation then local distribution in fully convolutional networks. *arXiv preprint arXiv:1909.07229*, 2019.

[27] Xiangtai Li, Houlong Zhao, Lei Han, Yunhai Tong, Shaohua Tan, and Kuiyuan Yang. Gated fully fusion for semantic segmentation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 11418–11425, 2020.

[28] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021.

[29] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

[30] Ange Lou, Shuyue Guan, and Murray Loew. Caranet: Context axial reverse attention network for segmentation of small medical objects. *arXiv preprint arXiv:2108.07368*, 2021.

[31] Yi Lu, Yaran Chen, Dongbin Zhao, and Jianxin Chen. Graph-fcn for image semantic segmentation. In *International symposium on neural networks*, pages 97–105. Springer, 2019.

[32] Haoxiang Ma, Hongyu Yang, and Di Huang. Boundary guided context aggregation for semantic segmentation. *arXiv preprint arXiv:2110.14587*, 2021.

[33] Prashant Mathur, Krishnan Sathishkumar, Meesha Chaturvedi, Priyanka Das, Kondalli Lakshminarayana Sudarshan, Stephen Santhappan, Vinodh Nallasamy, Anish John, Sandeep Narasimhan, Francis Selvaraj Roselind, et al. Cancer statistics, 2020: report from national cancer registry programme, india. *JCO Global Oncology*, 6:1063–1075, 2020.

[34] Xuebin Qin, Zichen Zhang, Chenyang Huang, Chao Gao, Masood Dehghan, and Martin Jagersand. Basnet: Boundary-aware salient object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7479–7489, 2019.

[35] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12179–12188, 2021.

[36] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.

[37] Juan Silva, Aymeric Histace, Olivier Romain, Xavier Dray, and Bertrand Granado. Toward embedded detection of polyps in wce images for early diagnosis of colorectal cancer. *International journal of computer assisted radiology and surgery*, 9(2):283–293, 2014.

[38] Nima Tajbakhsh, Suryakanth R Gurudu, and Jianming Liang. Automated polyp detection in colonoscopy videos using shape and context information. *IEEE transactions on medical imaging*, 35(2):630–644, 2015.

[39] Towaki Takikawa, David Acuna, Varun Jampani, and Sanja Fidler. Gated-scnn: Gated shape cnns for semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5229–5238, 2019.

[40] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021.

[41] Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific data*, 5(1):1–9, 2018.

[42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[43] David Vázquez, Jorge Bernal, F Javier Sánchez, Gloria Fernández-Esparrach, Antonio M López, Adriana Romero,

Michal Drozdzal, and Aaron Courville. A benchmark for endoluminal scene segmentation of colonoscopy images. *Journal of healthcare engineering*, 2017, 2017.

[44] Haonan Wang, Peng Cao, Jiaqi Wang, and Osmar R Zaiane. Uctransnet: rethinking the skip connections in u-net from a channel-wise perspective with transformer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 2441–2449, 2022.

[45] Jinfeng Wang, Qiming Huang, Feilong Tang, Jia Meng, Jionglong Su, and Sifan Song. Stepwise feature fusion: Local guides global. *arXiv preprint arXiv:2203.03635*, 2022.

[46] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pvt v2: Improved baselines with pyramid vision transformer. *Computational Visual Media*, pages 1–10, 2022.

[47] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018.

[48] Jun Wei, Yiwen Hu, Ruimao Zhang, Zhen Li, S Kevin Zhou, and Shuguang Cui. Shallow attention network for polyp segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 699–708. Springer, 2021.

[49] Jun Wei, Shuhui Wang, and Qingming Huang. F$^3$net: fusion, feedback and focus for salient object detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 12321–12328, 2020.

[50] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34, 2021.

[51] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.

[52] Li Zhang, Xiangtai Li, Anurag Arnab, Kuiyuan Yang, Yunhai Tong, and Philip HS Torr. Dual graph convolutional network for semantic segmentation. *arXiv preprint arXiv:1909.06121*, 2019.

[53] Li Zhang, Dan Xu, Anurag Arnab, and Philip HS Torr. Dynamic graph message passing networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3726–3735, 2020.

[54] Ruifei Zhang, Guanbin Li, Zhen Li, Shuguang Cui, Dahong Qian, and Yizhou Yu. Adaptive context selection for polyp segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 253–262. Springer, 2020.

[55] Yundong Zhang, Huiye Liu, and Qiang Hu. Transfuse: Fusing transformers and cnns for medical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 14–24. Springer, 2021.

[56] Zhenli Zhang, Xiangyu Zhang, Chao Peng, Xiangyang Xue, and Jian Sun. Exfuse: Enhancing feature fusion for semantic segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 269–284, 2018.

[57] Mingmin Zhen, Jinglu Wang, Lei Zhou, Shiwei Li, Tianwei Shen, Jiaxiang Shang, Tian Fang, and Long Quan. Joint semantic segmentation and boundary detection using iterative pyramid contexts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13666–13675, 2020.

[58] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6881–6890, 2021.

[59] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: A nested u-net architecture for medical image segmentation. In *Deep learning in medical image analysis and multimodal learning for clinical decision support*, pages 3–11. Springer, 2018.