

# Rethinking Boundary Detection in Deep Learning Models for Medical Image Segmentation

Yi Lin<sup>†</sup>, Dong Zhang<sup>†</sup>, Xiao Fang, Yufan Chen, Kwang-Ting Cheng, and Hao Chen<sup>✉</sup>

The Hong Kong University of Science and Technology, Hong Kong, China  
jhc@cse.ust.hk

**Abstract.** Medical image segmentation is a fundamental task in the community of medical image analysis. In this paper, a novel network architecture, referred to as Convolution, Transformer, and Operator (CTO), is proposed. CTO employs a combination of Convolutional Neural Networks (CNNs), Vision Transformer (ViT), and an explicit boundary detection operator to achieve high recognition accuracy while maintaining an optimal balance between accuracy and efficiency. The proposed CTO follows the standard encoder-decoder segmentation paradigm, where the encoder network incorporates a popular CNN backbone for capturing local semantic information, and a lightweight ViT assistant for integrating long-range dependencies. To enhance the learning capacity on boundary, a boundary-guided decoder network is proposed that uses a boundary mask obtained from a dedicated boundary detection operator as explicit supervision to guide the decoding learning process. The performance of the proposed method is evaluated on six challenging medical image segmentation datasets, demonstrating that CTO achieves state-of-the-art accuracy with a competitive model complexity.

**Keywords:** Medical Image Segmentation · CNNs · Vision Transformer · Boundary Detection · Network Architecture.

## 1 Introduction

Medical Image Segmentation (MISeg) aims to locate pixel-level semantic lesion areas and/or human organs of the given image, which is one of the fundamental yet challenging tasks in the community of medical image analysis [33,27]. In the past few years, this task has been extensively studied and applied to a wide range of downstream applications, *e.g.*, robotic surgery [16], cancer diagnosis [29], and treatment design [38]. To achieve a desired MISeg result, it is critical to extract a set of rich and discriminative image feature representations.

Recently, thanks to the successful utilization of Vision Transformer (ViT) on computer vision tasks [12], ViT-based methods have greatly promoted the accuracy of medical image analysis [21]. For example, the state-of-the-art methods for some medical image analysis tasks (*e.g.*, diagnosis [39], segmentation [8], and detection [35]) are based on the ViT framework [12]. Compared to CNNs-based methods, ViT has a

<sup>†</sup> Equal contribution; <sup>✉</sup> corresponding author.

The source code is available at <https://github.com/xiaofang007/CTO>

stronger capacity to capture long-range dependencies, which have been shown to be beneficial for visual recognition [41]. For a canonical ViT-based MISeg model, it first partitions the input image into image patches. Then, these patches are treated as tokens for interactions via a multi-head self-attention layer, where the positional embedding is used for capturing the relative spatial information if needed. Finally, a normalization strategy and feature regulation operations are used to generate the output. The above processes are connected to form a basic transformer block, and such a block is repeated to encode semantic representations for the MISeg head network.

Despite that ViT-based methods have achieved preliminary success, they inherently suffer from two potential problems, *i.e.*, lack of translation invariance and weakness in local features [8]. To address these two problems, the advanced CNNs-ViT hybrid architectures were proposed for MISeg, *e.g.*, TransUNet [8], UNETR [21], Swin UNETR [20]. These attempts add convolutional operations in a ViT framework for local feature interactions, and strategies can improve the model convergence are also used. Particularly, the CNNs-ViT hybrid methods for MISeg are mainly based on UNet [33] and add transformer blocks in the backbone networks [21], and skip connections [23].

The explicit boundary also matters - although this information is usually overlooked in the deep learning era. Compared to the implicit learning manner (*e.g.*, CNNs, and ViT), an explicit learning model provides an immediate feature learning pattern, which has remarkable advantages of simple implementation, high efficiency, and purposeful objective. In the recent past, boundary operators are gradually valued in some pixel-level tasks and have been used to explicitly enhance the learning capacity on localization [10,13,28]. For MISeg, we believe that the boundary operator should play a more important role. Because, intuitively, a lesion region can be regarded as a kind of noise compared to normal regions. Besides, empirically, the explicit learning strategy can help the implicit feature learning model improve its representation capacity.

We propose a new network architecture, called CTO (Convolution, Transformer, and Operator), for MISeg that combines CNNs, ViT, and boundary detection operators to leverage both local semantic information and long-range dependencies in the learning process. CTO follows the canonical encoder-decoder segmentation paradigm, where the encoder network is composed of a CNNs backbone and an assistant lightweight ViT branch. To enhance boundary learning capacity, we introduce a boundary-guided decoder network that uses a self-generated boundary mask extracted by boundary detection operators as explicit supervisions to guide the decoding learning process. Our CTO architecture has higher recognition accuracy and achieves a better trade-off between accuracy and efficiency compared to the advanced MISeg architectures. We evaluate CTO on six representative yet challenging MISeg datasets, *i.e.*, two ISIC datasets [19,11], PH2 [31], CoNIC [17], LiTS17 [4], and BTCV [24]. Experimental results demonstrate that our CTO can achieve: 1) a new high accuracy on these datasets; 2) a superior performance to state-of-the-art methods; 3) and with competitive model complexity and efficiency.

## 2 Related Work

**Medical Image Segmentation (MISeg).** The existing MISeg methods can be roughly divided into the following three camps: i) CNNs-based methods; ii) ViT-based methods; and iii) CNNs-ViT hybrid methods. One of the most notable commonalities among these methods is that they are mainly based on an encoder-decoder paradigm. In the first camp, there are representative V-Net [32], U-Net [33], Attention-UNet [34].

These methods use CNNs as the backbone to extract image features, and combine some elaborate tricks (*e.g.*, skip connection, multi-scale representation [7], feature interaction [6]) for feature enhancement. However, since convolution is inherently a local operation, methods in this camp may result in the problem of incomplete segmentation mask. In the second camp, there are Swin-UNet [5] and MissFormer [23]. Such methods use a ViT to replace CNNs as encoder/decoder to aggregate long-range feature dependencies. However, due to the limited number of medical images and the small inherent variability, such methods are difficult to optimize and have excessively high computational costs. In the third camp, there are TransUNet [8], UNETR [21] and Swin UNETR [20]. This type of method combines advantages of both CNNs and ViT, *i.e.*, the model can capture not only local information but also long-range feature dependencies. However, an obvious disadvantage is that they are computationally intensive and have high computation complexity. In our work, we propose to use a lightweight ViT as an assistant to help the mainstream CNN capture long-range feature dependencies. Besides, a boundary-enhanced feature, which is generated by an explicit boundary detection operator, is used to guide the decoding learning process.

**Operators in Image Processing.** The operator is a fundamental component in traditional digital image processing, where the boundary detection operator is the most core element. The commonly used boundary detection operators can be divided into: i) the first derivative operator (*e.g.*, Roberts, Prewitt, and Sobel), and ii) the second derivative operators (*e.g.*, Laplacian) [25]. Recently, boundary detection operators have been revived in pixel-level computer vision tasks, such as manipulation detection [10] and camouflaged object detection [13]. In this paper, the boundary detection operator is used as an explicit mask extractor to guide an implicit feature learning model for MISeg. Our contribution is to use feature maps of the intermediate layer to synthesize a high-quality boundary prediction without requiring additional information.

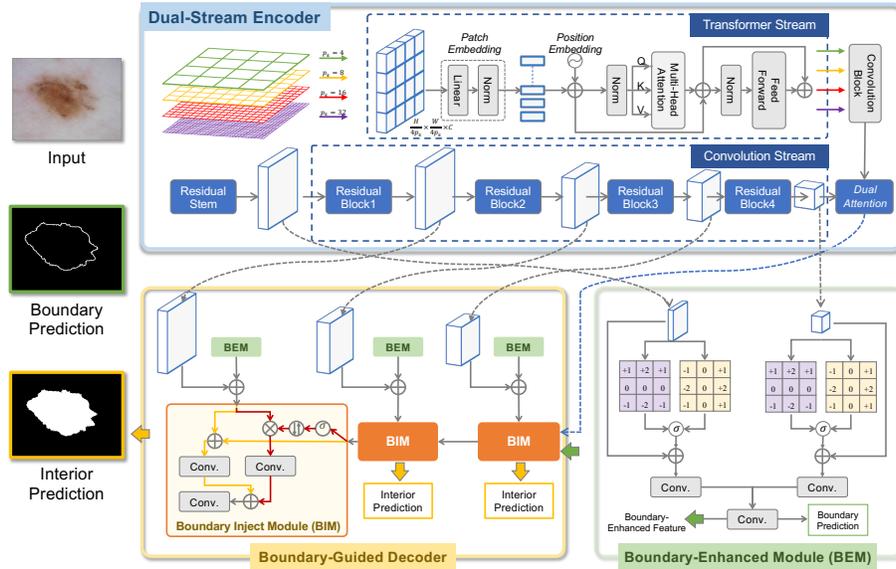
### 3 Convolution, Transformer, and Operator (CTO)

#### 3.1 Overview

The overall architecture of CTO is illustrated in Figure 1. For an input image  $X \in \mathbb{R}^{H \times W \times 3}$  with a spatial resolution of  $H \times W$  and  $C$  channels, we aim to predict a pixel-wise labelmap  $Y$ , where each pixel has been assigned a class label. The whole model follows an encoder-decoder pattern, which also adopts skip connections to aggregate low-level features from the encoder to the decoder. For the encoder, we design a dual-stream encoder (*ref.* Sec. 3.2), which combines a convolutional neural network (*i.e.*, Res2Net [15]) and a lightweight vision transformer to capture local feature dependencies and long-range feature dependencies between image patches, respectively. Such a combination will not bring many computational overheads. For the decoder, an operator-guided decoder (*ref.* Sec. 3.3) uses a boundary detection operator (*i.e.*, Sobel [25]) to guide the learning process via the generated boundary mask. The whole model is trained in an end-to-end manner.

#### 3.2 Dual-Stream Encoder

**The Mainstream Convolution Stream.** The convolution stream is used to capture local feature dependencies. To this end, we choose the strong yet efficient Res2Net [15] as the backbone, which is composed of one convolution stem and four



**Fig. 1.** Illustration of our CTO, which follows an encoder-decoder paradigm, where the encoder network consists of a mainstream CNNs and an assistant ViT. The decoder network employs a boundary detection operator to guide its learning process.

residual blocks, generating feature maps  $F_c^k$  with the spatial resolution of  $H/4 \times W/4$ ,  $H/8 \times W/8$ ,  $H/16 \times W/16$ , and  $H/32 \times W/32$ , respectively.

**The Assistant Transformer Stream.** The lightweight vision transformer (LightViT) is designed to capture the long-range feature dependencies between image patches in different scales. Specifically, the LightViT consists of multiple parallel lightweight transformer blocks that are fed with feature patches in different scales. All the transformer blocks share a similar structure, which consists of patch embedding layers and transformer encoding layers.

As shown in Figure 1, given the input feature map  $F_1^c \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times C}$ , we first divide it into  $\frac{HW}{16p^2}$  patches with size  $p \times p$ , and then flatten each patch into a vector  $\mathbf{v}_i \in \mathbb{R}^{p^2 \times C}$ . In our paper, we use four parallel transformer blocks, which are fed with feature patches in size of  $p = 4, 8, 16, 32$ . Then, we apply a linear projection to each patch vector to obtain the patch embedding  $\mathbf{e}_i \in \mathbb{R}^C$ . After that, patch embeddings along with the position embeddings are fed into the transformer encoding layers to obtain the output. Following [12], the encoding layers consist of a lightweight multi-head self-attention (MHSA) layer and a feed-forward network. MHSA receives a truncated query  $Q$ , key  $K$ , and value  $V$  as input, and then computes the attention score  $A \in \mathbb{R}^{N \times N}$  as follows:

$$A = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V, \quad (1)$$

where  $N$  is the size of patch number,  $d_k$  is the dimension of the key. The output of the MHSA layer is then fed into a feed-forward layer to obtain output  $F_t$ :

$$F_t = \text{FFN}(A), \quad (2)$$

where FFN is the feed-forward network with two linear layers with ReLU activation function. Then,  $F_t$  is reshaped into the same size as  $F_c^1$  to obtain the output. Outputs of all the transformer blocks are concatenated along the channel dimension and fed into the convolutional layer to obtain the final output.

### 3.3 Boundary-Guided Decoder

The boundary-guided decoder uses a gradient operator module to extract the boundary information of foreground objects. Then, the boundary-enhanced feature  $F_b$  is integrated into multi-level encoder’s features by a boundary optimization module, aiming to simultaneously characterize the intra- and inter-class consistency in the feature space, enriching the feature representative ability.

**Boundary Enhanced Module (BEM).** BEM takes the high-level  $F_c^4$  and low-level features  $F_c^1$  as inputs to extract the boundary information while filtering the trivial boundary irrelevant information. To achieve this goal, we apply Sobel operator [25] at both horizontal  $G_x$  and vertical  $G_y$  directions to obtain the gradient maps. Specifically, we utilize two  $3 \times 3$  parameter-fixed convolutions and apply convolution operation with stride 1. Two convolutions are defined as:

$$K_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}, \quad K_y = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix}. \quad (3)$$

Then, we apply the two convolutions to the input feature map to obtain the gradient maps  $M_x$  and  $M_y$ . After that, the gradient maps are normalized by a sigmoid function and then fused with the input feature map to obtain the edge-enhanced feature map  $F_e$ :

$$F_e = F_c \odot \sigma(M_{xy}), \quad (4)$$

where  $\odot$  denotes the element-wise multiplication,  $\sigma$  is the sigmoid function, and  $M_{xy}$  is the concatenation of  $M_x$  and  $M_y$  along the channel dimension. Then, we fuse the edge-enhanced feature maps of  $F_c^1$  and  $F_c^4$  with a simple stacked convolution layer in the bottleneck. Specifically, we first apply a  $1 \times 1$  convolution with a bilinear upsampling operation to the feature map  $F_c^4$  to obtain the feature map with the same size as  $F_c^1$ . Then, we separately apply  $1 \times 1$  convolution operation to equate the channel size of these two features. Finally, we concatenate these two feature maps along the channel dimension and apply a two-layer convolutions to get the final feature map  $\bar{F}_e$ . The output is supervised by the ground truth boundary map, which in turn eliminates the edge feature inside the objects, producing the boundary-enhanced feature  $F_b$ .

**Boundary Inject Module (BIM).** The obtained boundary-enhanced feature from BEM can be used as a prior to improve the image representation ability of the features produced by the encoder. We propose BIM that introduces a dual path boundary fusion scheme to promote the feature representation in both foreground and

background. Specifically, BIM takes two inputs: the channel-wise concatenation of the boundary-enhanced feature  $F_b$  and the corresponding feature  $F_c$  from the encoder network, and the feature from the previous decoder layer  $F_d^{j-1}$ . Then, these two inputs are fed into BIM, which contains two individual paths aiming to promote the feature representation in the foreground and background, respectively. For the foreground path, we directly concatenate the two inputs along the channel dimension, and then apply a sequential Conv-BN-ReLU (*i.e.*, convolution, batch normalization, ReLU activator) layers to obtain the foreground feature  $F_{fg}$ . For the background path, we design the background attention component to selectively focus on the background information, which is expressed as:

$$F_{bg} = \text{Convs} \left( (1 - \sigma(F_d^{j-1})) \odot F_c \right), \quad (5)$$

where Convs is a three-layer Conv-BN-ReLU layers,  $\sigma$  is the sigmoid function, and  $\odot$  denotes the element-wise multiplication. The term  $(1 - (\sigma F_d^{j-1}))$  is the background attention map, which is computed by first applying the sigmoid function to the feature map from the previous decoder layer, which will generate a foreground attention map. Then, we subtract the foreground attention map from 1 to obtain the background attention map. Finally, we concatenate the foreground feature  $F_{fg}$ , the background feature  $F_{bg}$ , and the previous decoder feature  $F_d^{j-1}$  along the channel dimension to obtain the final output  $F_d^j$ .

### 3.4 Overall Loss Function

Since the proposed CTO is a multi-task model (*i.e.*, interior and boundary segmentation), we define an overall loss function to jointly optimize these two tasks.

**Interior Segmentation Loss.** The interior segmentation loss is the weighted sum of cross-entropy loss  $\mathcal{L}_{\text{CE}}$  and mean intersection-over-union (mIoU) loss  $\mathcal{L}_{\text{mIoU}}$ , which are defined as:

$$\mathcal{L}_{\text{CE}} = -\frac{1}{N} \sum_{i=1}^N (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)), \quad (6)$$

$$\mathcal{L}_{\text{mIoU}} = 1 - \frac{\sum_{i=1}^N (y_i * \hat{y}_i)}{\sum_{i=1}^N (y_i + \hat{y}_i - y_i * \hat{y}_i)}, \quad (7)$$

where  $y_i$  and  $\hat{y}_i$  are the ground truth and the predicted label for the  $i$ -th pixel, respectively, and  $N$  is the total number of pixels of the image.

**Boundary Loss.** Considering the class imbalance problem between the foreground and background pixels in boundary detection, we employ the Dice Loss:

$$\mathcal{L}_{\text{Dice}} = 1 - \frac{2 \sum_{i=1}^N (y_i * \hat{y}_i)}{\sum_{i=1}^N (y_i + \hat{y}_i)}. \quad (8)$$

**Total Loss.** The total loss is composed of the major segmentation loss  $\mathcal{L}_{\text{seg}}$ , and the boundary loss  $\mathcal{L}_{\text{bnd}}$ . Note that for the boundary detection loss, we only consider the prediction from BEM, which takes encoder’s feature maps from the high-level layer (*i.e.*,  $F_b^4$ ) and low-level layer (*i.e.*,  $F_b^1$ ) as input. As for the major image segmentation loss, we apply the deep supervision strategy to obtain the prediction from the decoder’s feature at different levels. In summary, the total loss can be formulated as:

$$\mathcal{L} = \mathcal{L}_{\text{seg}} + \mathcal{L}_{\text{bnd}} = \sum_i^L (\mathcal{L}_{\text{CE}} + \mathcal{L}_{\text{mIoU}}) + \alpha \mathcal{L}_{\text{Dice}}, \quad (9)$$

where  $L$  is the number of BOMs, which is set to 3 in this work.  $\alpha$  is the weighting factor, which is set to 3 to balance the losses.

## 4 Experiments

### 4.1 Datasets and Evaluation Metrics

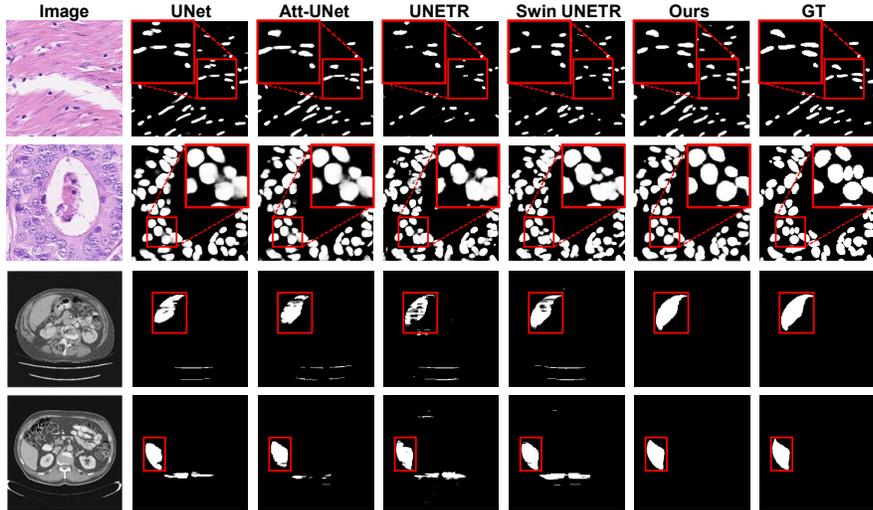
**Datasets.** We evaluate our CTO on six public MISeg datasets, including three datasets for skin lesion segmentation, *i.e.*, ISIC [19,11] and PH2 [31], the Colon Nuclei Identification and Counting (CoNIC) challenge dataset [17], the Liver Tumor Segmentation (LiTS17) Challenge dataset [4], and the Beyond the Cranial Vault (BTCV) challenge dataset [37]. As in [19,11], we perform 5-fold cross-validation on ISIC 2018, and train the model on ISIC 2016 and test it on PH2 [31]. BTCV is divided into 18 cases for training and 12 cases for test [5,8]. CoNIC and LiTS17 are randomly divided into training, validation, and test sets with a ratio of 7:1:2. **Evaluation Metrics.** Following [5,8,26], the commonly used Dice Coefficient (Dice), Intersection over Union (IoU), average Hausdorff Distance (HD) and Panoptic Quality (PQ) are used as the primary accuracy evaluation metrics. Besides, FLOPs and model parameters are used to evaluate the model efficiency.

### 4.2 Implementation Details

We optimize our model using the ADAM optimizer with an initial learning rate 1e-4. The default batch size is set to 32 with the image size of  $256 \times 256$ . The encoder is initialized with the pre-trained weights of Res2Net-50 [15] on ImageNet and then fine-tuned for 90 epochs on a single NVIDIA RTX 3090 GPU. All 3D volumes are inferenced in a sliding-window manner with the stride of 1, and the final segmentation results are obtained by stacking the prediction maps to reconstruct the 3D volume for evaluation. Except for a special statement, all the experimental settings follow the baseline paper [5,8,26,40].

### 4.3 Experimental Results

**Comparisons with State-of-the-Art Methods.** We compare our CTO with the state-of-the-art (SOTA) methods including U-Net [33], ResUNet [33] with ResNet-50 [22] as the backbone, VNet [32], ViT [12], TransUNet [8], and Swin-Unet [5]. On ISIC 2016 [19] & PH2 [31], we compare CTO with five related methods. The results are shown in Table 1. We can observe that CTO achieves 91.89% in Dice and 85.18%



**Fig. 2.** Visualizations on CoNIC [17] (top two rows) and LiTS17 [4] (bottom two rows). The red boxes highlight the main difference of each method.

**Table 1.** Comparisons with other methods on ISIC [19,11] & PH2 [31].

Methods	ISIC 2016 & PH2		Methods	ISIC 2018	
	Dice $\uparrow$	IoU $\uparrow$		Dice $\uparrow$	IoU $\uparrow$
SSLS[1]	78.38	68.16	Deeplabv3 [9]	88.4	80.6
MSCA[2]	81.57	72.33	U-Net++ [42]	87.9	80.5
FCN [30]	89.40	82.15	CE-Net [18]	89.1	81.6
Bi <i>et al</i> [3]	90.66	83.99	MedT [36]	85.9	77.8
Lee <i>et al</i> [26]	<u>91.84</u>	<u>84.30</u>	TransUNet [8]	<u>89.4</u>	<u>82.2</u>
CTO(Ours)	<b>91.89</b>	<b>85.18</b>	Ours	<b>91.2</b>	<b>84.5</b>

in IoU, which outperforms the SOTA methods by 0.05% and 0.88%, respectively. On ISIC 2018 [11], our CTO achieves 91.2% in Dice and 84.5% in IoU by 5-fold cross-validation, which outperforms the SOTA methods by 1.8% and 2.3%, respectively. On CoNIC [17], results are shown in Table 2, we can observe that our CTO achieves 79.77%, 66.42%, and 65.58% in Dice, IoU, and PQ, respectively, consistently outperforming other methods. Qualitative result comparisons are illustrated in Figure 2. We can observe that our CTO delineates more accurate object contours than other methods regarding diverse shapes and sizes of nuclei, especially on some blurred nuclei objects.

We also conduct experiments on 3D MISeg tasks. On LiTS17 [4], as shown in Table 2, our model achieves 91.50% in Dice and 84.59% in IoU, outperforming SOTA methods by 0.26% and 0.45%, respectively. On BTCV [24], as shown in Table 3, our CTO achieves 81.10% in Dice and 18.75% in HD, which outperforms the SOTA methods. In particular, as for Dice, our CTO outperforms the second, third, and fourth best methods by 1.98%, 3.33%, and 3.62%, respectively. Besides, the distinct improvements

**Table 2.** Comparisons with other methods on CoNIC [17] and LiTS17 [4].

Methods	CoNIC			LiTS17		Model Efficiency	
	Dice $\uparrow$	IoU $\uparrow$	PQ $\uparrow$	Dice $\uparrow$	IoU $\uparrow$	Param.(M)	GFLOPs
V-Net [32]	77.46	64.94	63.59	89.20	80.71	11.84	18.54
U-Net [33]	78.42	<u>66.39</u>	64.44	84.66	73.63	7.78	14.59
R50-UNet [33]	77.67	65.34	63.67	<u>91.24</u>	<u>84.14</u>	33.69	20.87
Att-UNet [34]	<u>79.48</u>	66.06	<u>65.25</u>	85.88	75.40	7.88	43.35
R50-AttUNet [34]	78.21	65.86	64.02	89.98	82.13	33.25	49.25
R50-ViT [12]	75.36	62.35	58.03	83.67	72.49	110.62	26.91
UNETR [21]	71.46	57.24	52.26	81.48	69.04	87.51	26.41
Swin-UNETR [20]	70.07	55.56	51.59	84.00	72.76	6.29	4.86
CTO(Ours)	<b>79.77</b>	<b>66.42</b>	<b>65.58</b>	<b>91.50</b>	<b>84.59</b>	59.82	22.72

**Table 3.** Comparisons with other methods on BTCV [24].

Methods	mDice $\uparrow$	HD $\downarrow$	Aorta	Gallb.	Kid(L)	Kid(R)	Liver	Panc.	Spleen	Stom.
V-Net [32]	68.81	-	75.34	51.87	77.10	<u>80.75</u>	87.84	40.05	80.56	56.98
DARR [14]	69.77	-	74.74	53.77	72.31	73.24	94.08	54.18	89.90	45.96
U-Net [33]	76.85	39.70	<u>89.07</u>	<b>69.72</b>	77.77	68.60	93.43	53.98	86.67	75.58
R50-UNet [33]	74.68	36.87	84.18	62.84	79.19	71.29	93.35	48.23	84.41	73.92
Att-UNet [34]	77.77	36.02	<b>89.55</b>	<u>68.88</u>	77.98	71.11	93.57	<u>58.04</u>	87.30	75.75
R50-AttUNet [34]	75.57	36.97	55.92	63.91	79.20	72.71	93.56	49.37	87.19	74.95
R50-ViT [12]	71.29	32.87	73.73	55.13	75.80	72.20	91.51	45.99	81.99	73.95
TransUNet [8]	77.48	31.69	87.23	63.13	81.87	77.02	94.08	55.86	85.08	75.62
SwinUNet [5]	79.12	21.55	85.47	66.53	<u>83.28</u>	79.61	<u>94.29</u>	56.58	<b>90.66</b>	76.60
CTO(Ours)	<b>81.10</b>	<b>18.75</b>	87.72	66.44	<b>84.49</b>	<b>81.77</b>	<b>94.88</b>	<b>62.74</b>	<u>90.60</u>	<b>80.20</b>

can be markedly observed for organs with blurry boundaries, *e.g.*, the “pancreas” and the “stomach”, where our model achieves significant gains over the SOTA methods, *i.e.*, 4.70% and 3.60% in Dice, respectively. As for the model efficiency, we can observe that CTO achieves competitive performance improvements with comparable FLOPs and parameters. **Ablation Study.** We conduct ablation studies to explore the effectiveness of each component in CTO. In Table 4, we compare the performance of CTO variants on ISIC 2018 [11]: 1) CNNs, only the convolution stream; 2) +LightViT, the dual-stream encoder with convolution and Transformer; 3) +CBM, adding the boundary supervision with the same architecture of BEM, except the Sobel layer; 4) +BEM, the boundary-enhanced module; 5) +BIM, the boundary inject module. All the components consistently boost the performance by 0.99%, 1.09%, 1.20%, 2.89% in Dice, respectively. Especially, we observe that the boundary supervision (*i.e.*, BIM) is crucial for MISeg.

**Table 4.** Ablation study results on ISIC 2018 [11]. \* means the component achieves significant performance improvement with  $p < 0.05$  via paired t-test.

CNNs	LightViT	CBM	BEM	BIM	Dice $\uparrow$	IoU $\uparrow$
✓					88.32	81.51
✓	✓				89.31* <sub>+0.99</sub>	82.47* <sub>+0.96</sub>
✓	✓	✓			89.41* <sub>+1.09</sub>	82.51 <sub>+1.00</sub>
✓	✓	✓	✓		89.52 <sub>+1.20</sub>	82.81 <sub>+1.30</sub>
✓	✓	✓	✓	✓	91.21* <sub>+2.89</sub>	84.45* <sub>+2.94</sub>

## 5 Conclusion

In this study, a new network architecture named CTO is proposed for MISeg. Compared to advanced MISeg architectures, CTO achieves a better balance between recognition accuracy and computational efficiency. The contribution of this paper is the utilization of intermediate feature maps to synthesize a high-quality boundary supervision mask without requiring additional information. Results from experiments conducted on six publicly available datasets demonstrate the superiority of CTO over state-of-the-art methods, and the effectiveness of each of its components. Future work includes the extension of the concept of a couple-stream encoder to various advanced backbone architectures, and the potential adaptation of CTO to a 3D manner.

## References

- Ahn, E., Bi, L., Jung, Y.H., Kim, J., Li, C., Fulham, M., Feng, D.D.: Automated saliency-based lesion segmentation in dermoscopic images. In: International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC) (2015)
- Bi, L., Kim, J., Ahn, E., Feng, D., Fulham, M.: Automated skin lesion segmentation via image-wise supervised learning and multi-scale superpixel based cellular automata. In: IEEE International Symposium on Biomedical Imaging (ISBI) (2016)
- Bi, L., Kim, J., Ahn, E., Kumar, A., Fulham, M., Feng, D.: Dermoscopic image segmentation via multistage fully convolutional networks. *IEEE Transactions on Biomedical Engineering* **64**(9), 2065–2074 (2017)
- Bilic, P., Christ, P., Li, H.B., Vorontsov, E., Ben-Cohen, A., Kaissis, G., Szeskin, A., Jacobs, C., Mamani, G.E.H., Chartrand, G., et al.: The liver tumor segmentation benchmark (lits). *Medical Image Analysis* **84**, 102680 (2023)
- Cao, H., Wang, Y., Chen, J., Jiang, D., Zhang, X., Tian, Q., Wang, M.: Swin-Unet: Unet-like pure transformer for medical image segmentation. *arXiv* (2021)
- Chen, H., Dou, Q., Yu, L., Qin, J., Heng, P.A.: Voxresnet: Deep voxelwise residual networks for brain segmentation from 3d mr images. *NeuroImage* **170**, 446–455 (2018)
- Chen, H., Qi, X., Yu, L., Heng, P.A.: Dcan: deep contour-aware networks for accurate gland segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
- Chen, J., Lu, Y., Yu, Q., Luo, X., Adeli, E., Wang, Y., Lu, L., Yuille, A.L., Zhou, Y.: TransUnet: Transformers make strong encoders for medical image segmentation. *arXiv* (2021)

9. Chen, L.C., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous convolution for semantic image segmentation. arXiv (2017)
10. Chen, X., Dong, C., Ji, J., Cao, J., Li, X.: Image manipulation detection by multi-view multi-scale supervision. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2021)
11. Codella, N., Rotemberg, V., Tschandl, P., Celebi, M.E., Dusza, S., Gutman, D., Helba, B., Kalloo, A., Liopyris, K., Marchetti, M., et al.: Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (ISIC). arXiv (2019)
12. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. In: International Conference on Learning Representations (ICLR) (2020)
13. Fan, D.P., Ji, G.P., Sun, G., Cheng, M.M., Shen, J., Shao, L.: Camouflaged object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2020)
14. Fu, S., Lu, Y., Wang, Y., Zhou, Y., Shen, W., Fishman, E., Yuille, A.: Domain adaptive relational reasoning for 3D multi-organ segmentation. In: International Conference on Medical Image Computing and Computer Assisted Intervention (MICCAI). Springer (2020)
15. Gao, S.H., Cheng, M.M., Zhao, K., Zhang, X.Y., Yang, M.H., Torr, P.: Res2net: A new multi-scale backbone architecture. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **43**(2), 652–662 (2019)
16. Gao, X., Jin, Y., Zhao, Z., Dou, Q., Heng, P.A.: Future frame prediction for robot-assisted surgery. In: IPMI (2021)
17. Graham, S., Jahanifar, M., Vu, Q.D., Hadjigeorgiou, G., Leech, T., Snead, D., Raza, S.E.A., Minhas, F., Rajpoot, N.: CoNIC: Colon nuclei identification and counting challenge 2022. arXiv (2021)
18. Gu, Z., Cheng, J., Fu, H., Zhou, K., Hao, H., Zhao, Y., Zhang, T., Gao, S., Liu, J.: CE-NET: Context encoder network for 2D medical image segmentation. *IEEE Transactions on Medical Imaging* **38**(10), 2281–2292 (2019)
19. Gutman, D., Codella, N.C.F., Celebi, E., Helba, B., Marchetti, M., Mishra, N., Halpern, A.: Skin lesion analysis toward melanoma detection: A challenge at the international symposium on biomedical imaging (ISBI) 2016, hosted by the international skin imaging collaboration (ISIC). In: arXiv (2016)
20. Hatamizadeh, A., Nath, V., Tang, Y., Yang, D., Roth, H.R., Xu, D.: Swin UNETR: Swin transformers for semantic segmentation of brain tumors in MRI images. In: Inter. MICCAI Brain Les. Workshop (2022)
21. Hatamizadeh, A., Tang, Y., Nath, V., Yang, D., Myronenko, A., Landman, B., Roth, H.R., Xu, D.: UNETR: Transformers for 3D medical image segmentation. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) (2022)
22. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
23. Huang, X., Deng, Z., Li, D., Yuan, X.: Missformer: An effective medical image segmentation transformer. arXiv (2021)
24. Irshad, S., Gomes, D.P., Kim, S.T.: Improved abdominal multi-organ segmentation via 3D boundary-constrained deep neural networks. arXiv (2022)
25. Kanopoulos, N., Vasanthavada, N., Baker, R.L.: Design of an image edge detection filter using the Sobel operator. *IEEE Jour. of Sol.-Stat. Cir.* **23**(2), 358–367 (1988)

26. Lee, H.J., Kim, J.U., Lee, S., Kim, H.G., Ro, Y.M.: Structure boundary preserving segmentation for medical image with ambiguous boundary. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2020)
27. Lin, Y., Liu, L., Ma, K., Zheng, Y.: Seg4reg+: Consistency learning between spine segmentation and cobb angle regression. In: MICCAI (2021)
28. Lin, Y., Qu, Z., Chen, H., Gao, Z., Li, Y., Xia, L., Ma, K., Zheng, Y., Cheng, K.T.: Label propagation for annotation-efficient nuclei segmentation from pathology images. arXiv preprint arXiv:2202.08195 (2022)
29. Lin, Y., Su, J., Wang, X., Li, X., Liu, J., Cheng, K.T., Yang, X.: Automated pulmonary embolism detection from CTPA images using an end-to-end convolutional neural network. In: International Conference on Medical Image Computing and Computer Assisted Intervention (MICCAI) (2019)
30. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2015)
31. Mendonça, T., Ferreira, P.M., Marques, J.S., Marcal, A.R., Rozeira, J.: Ph 2-a dermoscopic image database for research and benchmarking. In: EMBC (2013)
32. Milletari, F., Navab, N., Ahmadi, S.A.: V-Net: Fully convolutional neural networks for volumetric medical image segmentation. In: 3DV. IEEE (2016)
33. Ronneberger, O., Fischer, P., Brox, T.: U-Net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical Image Computing and Computer Assisted Intervention (MICCAI) (2015)
34. Schlemper, J., Oktay, O., Schaap, M., Heinrich, M., Kainz, B., Glocker, B., Rueckert, D.: Attention gated networks: Learning to leverage salient regions in medical images. *Medical Image Analysis* **53**, 197–207 (2019)
35. Shamshad, F., Khan, S., Zamir, S.W., Khan, M.H., Hayat, M., Khan, F.S., Fu, H.: Transformers in medical imaging: A survey. arXiv (2022)
36. Valanarasu, J.M.J., Oza, P., Hacihaliloglu, I., Patel, V.M.: Medical transformer: Gated axial-attention for medical image segmentation. arXiv (2021)
37. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)* (2017)
38. Wijeratne, P.A., Alexander, D.C., Initiative, A.D.N., et al.: Learning transition times in event sequences: The temporal event-based model of disease progression. In: IPMI (2021)
39. Wu, J., Fang, H., Shang, F., Yang, D., Wang, Z., Gao, J., Yang, Y., Xu, Y.: SeATrans: Learning segmentation-assisted diagnosis model via transformer. In: International Conference on Medical Image Computing and Computer Assisted Intervention (MICCAI) (2022)
40. Zhang, D., Lin, Y., Chen, H., Tian, Z., Yang, X., Tang, J., Cheng, K.T.: Deep learning for medical image segmentation: tricks, challenges and future directions. arXiv (2022)
41. Zhang, D., Tang, J., Cheng, K.T.: Graph reasoning transformer for image parsing. In: ACM MM (2022)
42. Zhou, Z., Siddiquee, M.M.R., Tajbakhsh, N., Liang, J.: Unet++: A nested U-Net architecture for medical image segmentation. In: DLMI (2018)