

# Dual Cross-Attention for Medical Image Segmentation

Gorkem Can Ates<sup>1\*</sup>, Prasoon Mohan<sup>1</sup>, Emrah Celik<sup>1</sup>

<sup>1</sup>University of Miami

\*Corresponding to: gca45@miami.edu

## Abstract

We propose Dual Cross-Attention (DCA), a simple yet effective attention module that is able to enhance skip-connections in U-Net-based architectures for medical image segmentation. DCA addresses the semantic gap between encoder and decoder features by sequentially capturing channel and spatial dependencies across multi-scale encoder features. First, the Channel Cross-Attention (CCA) extracts global channel-wise dependencies by utilizing cross-attention across channel tokens of multi-scale encoder features. Then, the Spatial Cross-Attention (SCA) module performs cross-attention to capture spatial dependencies across spatial tokens. Finally, these fine-grained encoder features are up-sampled and connected to their corresponding decoder parts to form the skip-connection scheme. Our proposed DCA module can be integrated into any encoder-decoder architecture with skip-connections such as U-Net and its variants. We test our DCA module by integrating it into six U-Net-based architectures such as U-Net, V-Net, R2Unet, ResUnet++, DoubleUnet and MultiResUnet. Our DCA module shows Dice Score improvements up to 2.05% on GlaS, 2.74% on MoNuSeg, 1.37% on CVC-ClinicDB, 1.12% on Kvasir-Seg and 1.44% on Synapse datasets. Our codes are available at: <https://github.com/gorkemcanates/Dual-Cross-Attention>

achieve a better feature representation. Such a scheme helps the model recover the contextual information loss during the down-sampling process in the encoder by simply concatenating encoder features on different scales with their corresponding parts in the decoder. Motivated by the success of U-Net and the skip-connection scheme, several architectural designs have been developed [10, 39, 75, 1, 45, 29, 27]. These U-Net variants have successfully improved vanilla U-Net in various medical image segmentation tasks by adapting sophisticated designs such as residual [20] and recurrent [34] connections into the encoder-decoder framework and/or improving the plain skip-connection scheme through further enhancing the encoder features before connecting to the decoder.

Although U-Net as well as its variants achieved good performance on various medical image segmentation tasks, there still exist performance limitations. The first limitation comes from the locality of convolutions which cannot capture the long-range dependencies across different features. This is mainly caused by the nature of the convolutional operation which gradually obtains the local receptive fields using local kernels rather than extracting global feature interactions at once [22]. The second limitation is the semantic gap caused by skip-connections when simply concatenating encoder and decoder features. Recently, Wang *et al.* [57] showed that the plain skip-connection scheme presented by U-Net is not sufficient enough to model the global multi-scale context and it is in fact essential to effectively fuse the low-level encoder features before connecting them to their corresponding decoder parts. As mentioned before, several U-Net variations tackled such a semantic gap problem by adding a series of convolutional or residual layers to further fuse the low-level encoder and decoder features in order to adequately connect them to their decoder counterparts. One approach is U-Net++ [75], which connects encoder features to decoder features through a series of nested, dense skip pathways. Another approach is MultiResUnet [25], which introduced residual paths by applying a series of convolutional layers with residual connections on encoder features propagating to their decoder counterparts. Despite improving the quality of the skip-connections, both of these meth-

## 1. Introduction

Convolutional Neural Networks (CNNs) have become the de-facto standard for accurate medical image segmentation because of their strong and complex mapping capability [17]. Fully Convolutional Networks (FCNs) [37], particularly convolutional encoder-decoder networks have drawn much attention in the past years due to their tremendous success in various medical image segmentation tasks [56, 23]. U-Net [44], in particular, achieved superior performance due to its skip-connection scheme, which connects low-level features extracted by the encoder to the decoder to

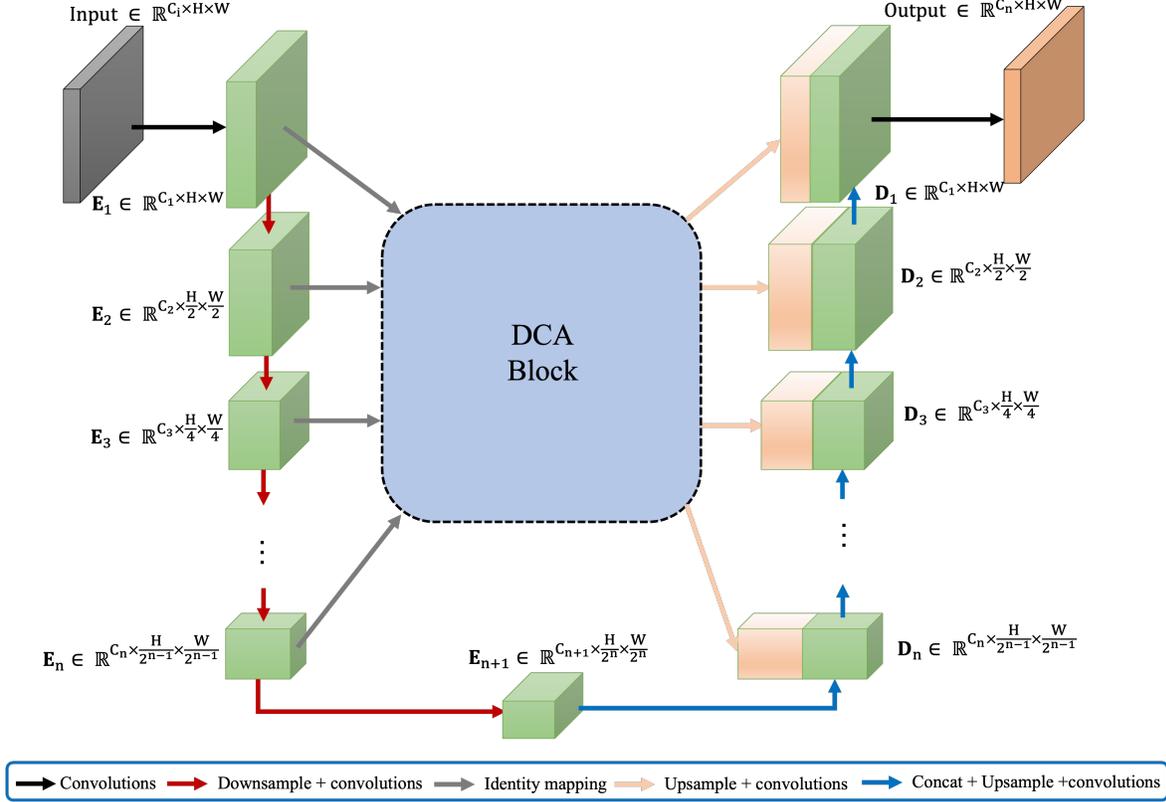


Figure 1. Encoder decoder architecture with our proposed DCA block. DCA block can be integrated into any encoder-decoder architecture with skip connections. It takes multi-scale features from different encoder stages, produces enhanced representations and connects them to their decoder counterparts.

ods still struggle to reduce the semantic gap between encoder and decoder features.

Recently, Transformers [54], originally proposed for Natural Language Processing (NLP), have become the dominant architecture in computer vision with the pioneering work of Vision Transformer (ViT) [14]. Subsequently, its success has been shown in various vision tasks including image classification [36, 52, 70, 58, 24, 13, 32], object detection [5, 15, 63], segmentation [47, 66, 74, 19] and beyond [60, 6, 68, 33, 49, 42, 61]. The self-attention mechanism in transformers certainly plays a key role in their success due to its ability to directly capture long-range dependencies [43]. Motivated by the effectiveness of self-attention, researchers proposed to combine channel self-attention with self-attention to also capture the channel-wise interactions. These dual attention schemes based on self-attention have been shown to improve performance on various vision tasks. For instance, Fu *et al.* [16] proposed DANet for scene segmentation, where a channel and a spatial attention module were integrated at the end of a dilated FCN to model the semantic dependencies in both spatial and channel dimensions. Those two attention modules were applied in a parallel manner to capture channel and

spatial dependencies separately. Then, those two outputs were fused using element-wise summation to further improve the feature representations, which significantly improved the segmentation performance. Mou *et al.* [40] performed a similar fusing strategy by summing the channel and spatial attention in a convolutional encoder-decoder network for the segmentation of curvilinear structures. Liu *et al.* [35], on the other hand, used concatenation to fuse the outputs of channel and spatial attention modules. A recent study [11] proposed a different dual attention approach named DaViT, in which spatial and channel attention modules were utilized sequentially. That is, self-attention is first applied to the spatial tokens, then the outputs of the spatial attention module are fed into the channel attention module which generates the fine-grained features by further extracting long-range feature interactions. Such dual attention mechanisms showed that the channel-self attention mechanism can also provide useful information. In fact, Wang *et al.* [57] showed that even stand-alone channel-wise attention can effectively capture the global context. They introduced channel cross-attention to tackle the semantic gap problem in U-Net by utilizing cross-attention in the channel axis of the multi-scale encoder features to capture the long-

range channel dependencies. Their cross-attention mechanism improved the segmentation performance on different medical image segmentation datasets when utilized with the U-Net structure, showing its promise in reducing the semantic gap between encoder and decoder features.

Motivated by the success of sequential dual attention [11] and the channel cross-attention [57], we propose Dual Cross-Attention (DCA), an attention module that effectively extracts channel and spatial-wise interdependencies across multi-scale encoder features to tackle the semantic gap problem. A major difference between channel and spatial cross-attention mechanisms is that channel-cross attention extracts global channel-wise context by fusing all spatial positions together between any given two channels, whereas spatial-cross attention mechanism extracts global spatial context by capturing spatial interdependencies between any two positions of the multi-scale encoder channels. A delicate incorporation of such channel and spatial cross-attention modules can further extract long-range contextual information by capturing the rich global context in both channel and spatial dimensions. Note that our main purpose here is to develop a well-structured mechanism that acts as a powerful bridge between encoder and decoder at a slight parameter increase rather than building an end-to-end network with an extensive amount of parameters as in [7]. Note that, MetaFormer [71] the general, abstracted architecture in end-to-end transformers consists of a patch embedding operation (e.g., convolutions) [65, 73, 64, 55] and two main building blocks, namely a token mixer (e.g., attention [54]) and a two-layer multi-layer perceptron (MLP) block with a non-linear activation (e.g., GeLU [21]) [72]. All of these blocks typically come with large computational costs, especially when multiple MetaFormer blocks are stacked together to form an end-to-end vision transformer network [14, 36, 50, 51, 52, 71, 69]. We minimize such computational costs by first introducing a new patch embedding operation, where we extract tokens by utilizing 2D average pooling which does not require any additional parameters and project them with  $1 \times 1$  depth-wise convolutions. Then, we replace the linear projections in both channel and spatial attention modules with depth-wise convolutions [9]. Finally, we completely eliminate the MLP layer to further reduce the number of parameters and we up-sample the DCA outputs to connect them to the decoder, thus forming our DCA mechanism. As we will show later, such modifications can still improve the segmentation performance while minimizing the computational overhead that occurs in MetaFormer. Extensive experiments using six state-of-the-art U-Net-based architectures and five benchmark medical image segmentation datasets show that our DCA module can significantly improve segmentation performance with minimal computational overhead.

## 2. Dual Cross-Attention (DCA)

Fig.1 illustrates the integration of our proposed DCA block into a general FCN architecture with skip-connections. The architecture of our DCA block is invariant to the number of encoder stages. That is, given  $n+1$  multi-scale encoder stages, the DCA block takes multi-scale features from the first  $n$  stages as input (outputs of the last convolutional layers in each stage), produces enhanced representations, and connects them to their corresponding  $n$  decoder stages. As shown in Fig.2a, we divide our DCA block into two main stages. The first stage consists of a multi-scale patch embedding module to obtain encoder tokens. In the second stage, we perform our proposed DCA mechanism using channel cross-attention (CCA) and spatial cross-attention (SCA) modules on these encoder tokens to capture long-range dependencies. Finally, we apply layer normalization [2] and GeLU [21] sequence and upsample those tokens to connect them to their decoder counterparts.

### 2.1. Patch Embedding from Multi-Scale Encoder Stages

We start by extracting patches from  $n$  multi-scale encoder stages (i.e., skip-connection layers). Given  $n$  encoder stages from different scales  $\mathbf{E}_i \in \mathbb{R}^{C_i \times \frac{H}{2^{i-1}} \times \frac{W}{2^{i-1}}}$  and patch sizes  $P_i^s = \frac{P^s}{2^{i-1}}$  where  $i = 1, 2, \dots, n$ , we extract patches using 2D average pooling where the pool size and stride are  $P_i^s$  and apply projection using  $1 \times 1$  depth-wise convolutions over the flattened 2D patches:

$$\mathbf{T}_i = \text{DConv1D}_{E_i}(\text{Reshape}(\text{AvgPool2D}_{E_i}(\mathbf{E}_i))) \quad (1)$$

where  $\mathbf{T}_i \in \mathbb{R}^{P \times C_i}$ , ( $i = 1, 2, \dots, n$ ) represents the flattened patches for the  $i^{\text{th}}$  encoder stage. Note that  $P$  is the number of patches which is the same for each  $\mathbf{T}_i$  so that we can utilize cross-attention across those tokens.

### 2.2. Channel Cross-Attention (CCA)

As shown in Fig.2b, each token  $\mathbf{T}_i$  is fed into the CCA module. We first perform layer normalization (LN) [2] on each  $\mathbf{T}_i$ . Then, by following the cross-attention strategy in [57], we concatenate tokens  $\mathbf{T}_i$ , ( $i = 1, 2, \dots, n$ ) along the channel dimension to create our keys and values  $\mathbf{T}_c$  while we use  $\mathbf{T}_i$  for the queries. Although linear projection is typically utilized in conventional self-attention, recent studies successfully adapted convolutions into self-attention in order to bring locality as well as to reduce computational complexity [59, 64, 62]. Depth-wise convolutions, in particular, have been utilized in self-attention because of their ability to capture local information with negligible additional computational cost [67, 38, 18, 8, 32]. Motivated by prior work and such advantages, we replace all linear projections with  $1 \times 1$  depth-wise convolutional projections:

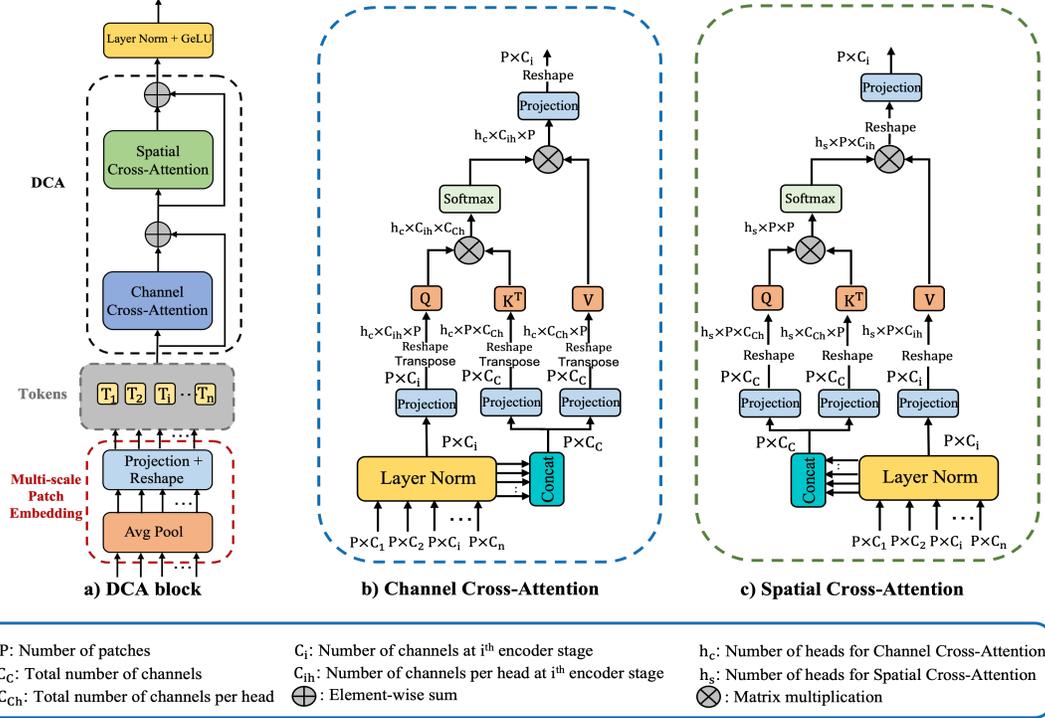


Figure 2. Architecture of our proposed DCA block (a). It consists of b) Channel Cross-Attention and c) Spatial Cross-Attention modules to capture long-range interactions.

$$\mathbf{Q}_i = \text{DConv1D}_{\mathbf{Q}_i}(\mathbf{T}_i) \quad (2)$$

$$\mathbf{K} = \text{DConv1D}_{\mathbf{K}}(\mathbf{T}_c) \quad (3)$$

$$\mathbf{V} = \text{DConv1D}_{\mathbf{V}}(\mathbf{T}_c) \quad (4)$$

where  $\mathbf{Q}_i \in \mathbb{R}^{P \times C_i}$ ,  $\mathbf{K} \in \mathbb{R}^{P \times C_c}$ ,  $\mathbf{V} \in \mathbb{R}^{P \times C_c}$  are the projected queries, keys, and values, respectively. In order to utilize cross-attention along the channel dimension, we take the transpose of queries, keys, and values. Thus, CCA takes the following form:

$$\text{CCA}(\mathbf{Q}_i, \mathbf{K}, \mathbf{V}) = \text{Softmax} \left( \frac{\mathbf{Q}_i^T \mathbf{K}}{\sqrt{C_c}} \right) \mathbf{V}^T \quad (5)$$

where  $\mathbf{Q}_i$ ,  $\mathbf{K}$ , and  $\mathbf{V}$  are matrices representing the queries, keys, and values, respectively, and  $\frac{1}{\sqrt{C_c}}$  is the scaling factor. The output of cross-attention is a weighted sum of the values, where the weights are determined by the similarity between the queries and keys. The output of the softmax function is then used to weight the values. Finally, we apply depth-wise convolutional projections to the outputs of the cross-attention and feed them into the SCA module.

### 2.3. Spatial Cross-Attention (SCA)

SCA module is illustrated in Fig.2c. Given the reshaped outputs  $\bar{\mathbf{T}}_i \in \mathbb{R}^{P \times C_i}$ , ( $i = 1, 2, \dots, n$ ) of the CCA module,

we perform layer normalization and concatenation along the channel dimension. Unlike the CCA module, we utilize concatenated tokens  $\bar{\mathbf{T}}_c$  as queries and keys while we use each token  $\bar{\mathbf{T}}_i$  as values. We utilize  $1 \times 1$  depth-wise projection over the queries, keys, and values:

$$\mathbf{Q} = \text{DConv1D}_{\mathbf{Q}}(\bar{\mathbf{T}}_c), \quad (6)$$

$$\mathbf{K} = \text{DConv1D}_{\mathbf{K}}(\bar{\mathbf{T}}_c), \quad (7)$$

$$\mathbf{V}_i = \text{DConv1D}_{\mathbf{V}_i}(\bar{\mathbf{T}}_i) \quad (8)$$

where  $\mathbf{Q} \in \mathbb{R}^{P \times C_c}$ ,  $\mathbf{K} \in \mathbb{R}^{P \times C_c}$ ,  $\mathbf{V}_i \in \mathbb{R}^{P \times C_i}$  are the projected queries, keys and values, respectively. Then, SCA can be expressed as:

$$\text{SCA}(\mathbf{Q}, \mathbf{K}, \mathbf{V}_i) = \text{Softmax} \left( \frac{\mathbf{Q} \mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}_i \quad (9)$$

Here,  $\mathbf{Q}$ ,  $\mathbf{K}$ , and  $\mathbf{V}_i$  are matrices representing the query, key, and value embeddings, respectively, and  $\frac{1}{\sqrt{d_k}}$  is the scaling factor. For the multi-head case  $d_k = \frac{C_c}{h_c}$  where  $h_c$  is the number of heads. Outputs of the SCA module are then projected using depth-wise convolutions to form DCA outputs. Then, we apply layer normalization and GeLU to those DCA outputs. Finally,  $n$  outputs of the DCA block are connected to their corresponding decoder parts by up-sample layers followed by  $1 \times 1$  convolution, batch normalization [26] and a ReLU [41] sequence. Note that the major

Model	Params	GlaS		MoNuSeg		CVC-ClinicDB		Kvasir-Seg		Synapse	
		DSC (%)	IoU (%)								
U-Net	8.64M	88.87	79.98	77.14	62.79	<b>89.63</b>	<b>81.43</b>	82.99	71.01	78.55	67.37
<b>U-Net (DCA)</b>	8.75M	<b>89.66</b>	<b>81.29</b>	<b>78.13</b>	<b>64.11</b>	89.53	81.28	<b>84.03</b>	<b>72.53</b>	<b>78.98</b>	<b>67.97</b>
ResUnet++	13.1M	85.43	74.62	75.68	60.87	89.46	81.14	<b>82.26</b>	<b>69.93</b>	75.91	64.61
<b>ResUnet++ (DCA)</b>	13.1M	<b>87.35</b>	<b>77.56</b>	<b>77.40</b>	<b>63.13</b>	<b>90.19</b>	<b>82.32</b>	82.07	69.74	<b>77.35</b>	<b>66.43</b>
MultiResUnet	7.24M	<b>88.99</b>	<b>80.18</b>	76.99	62.59	89.52	81.35	81.34	68.66	78.12	67.30
<b>MultiResUnet (DCA)</b>	7.35M	88.86	79.98	<b>78.52</b>	<b>64.63</b>	<b>89.95</b>	<b>81.91</b>	<b>82.32</b>	<b>70.00</b>	<b>79.50</b>	<b>68.65</b>
R2Unet	9.78M	85.16	74.26	78.20	64.20	88.12	78.88	81.07	68.28	75.86	63.94
<b>R2Unet (DCA)</b>	9.89M	<b>87.21</b>	<b>77.37</b>	<b>78.52</b>	<b>64.64</b>	<b>88.39</b>	<b>79.28</b>	<b>82.19</b>	<b>69.89</b>	<b>75.90</b>	<b>64.85</b>
V-Net	35.97M	88.78	79.85	74.79	59.74	88.09	79.02	80.79	68.07	79.27	68.58
<b>V-Net (DCA)</b>	36.08M	<b>89.03</b>	<b>80.27</b>	<b>77.53</b>	<b>63.31</b>	<b>89.46</b>	<b>81.07</b>	<b>81.92</b>	<b>69.53</b>	<b>79.58</b>	<b>69.00</b>
DoubleUnet	29.68M	89.07	80.30	77.16	62.82	90.20	82.35	84.40	73.08	79.76	69.31
<b>DoubleUnet (DCA)</b>	30.68M	<b>89.90</b>	<b>81.68</b>	<b>79.50</b>	<b>65.97</b>	<b>90.86</b>	<b>83.47</b>	<b>85.16</b>	<b>74.34</b>	<b>80.22</b>	<b>69.80</b>

Table 1. Performance comparison for plain and DCA integrated models on different datasets using DSC and IoU metrics. Boldfaced results indicate better results.

difference between cross-attention and self-attention is that cross-attention creates attention maps by fusing multi-scale encoder features together rather than utilizing each stage individually which also allows the cross-attention to capture long-range dependencies between different stages of the encoder.

### 3. Experiments

#### 3.1. Datasets

We conduct our experiments with five benchmark medical image segmentation datasets, including GlaS [46], MoNuSeg [31], CVC-ClinicDB [3], Kvasir-SEG [28] and Synapse. GlaS is a Gland Segmentation dataset consisting of 85 images for training and 80 images for testing. MoNuSeg is a nuclear segmentation dataset for digital microscopic tissue images. It includes 30 images for training and 14 images for testing. CVC-ClinicDB is a colonoscopy image dataset that includes a total of 612 images with their annotations. Kvasir-SEG is a polyp segmentation dataset that has 1000 annotated images. Following [29], we randomly split 80% of CVC-ClinicDB and Kvasir-SEG datasets into training and 20% for testing. Synapse is a multi-organ segmentation dataset, consisting of 30 abdominal CT scans in 8 abdominal organs (aorta, gallbladder, spleen, left kidney, right kidney, liver, pancreas, spleen, stomach). Following [7] and [57], we randomly choose 18 of these scans for training and 12 for testing.

#### 3.2. Models

We test our DCA mechanism using six models with skip-connections including U-Net [44], V-Net [39], R2Unet [1], ResUnet++ [29], DoubleUnet [27] and MultiResUnet [25]. For U-Net, V-Net, R2Unet and ResUnet++, we simply place our DCA block between the encoder and decoder stages as shown in Fig.1. For MultiResUnet, we place our DCA

block at the end of the residual paths to improve the residual path outputs. DoubleUnet has three skip-connection schemes since two U-Net architectures are stacked on top of each other. The first scheme connects the first encoder (VGG19) to the first decoder, the second scheme connects the first encoder features to the second decoder and the last scheme connects the second encoder features to the second decoder. In our experiments, we integrated three DCA blocks into all those three skip-connection schemes.

#### 3.3. Implementation Details

We implement our method in Pytorch using two NVIDIA Tesla V100-SXM2-16GB GPUs. We resize images to  $224 \times 224$  for all datasets [7, 57, 4]. We set the patch size to 4 for ResUnet++ and 8 for the remaining models. We set the number of heads in CCA to 1 while we set it to 4 for SCA. As for the data augmentation, we perform random rotations and random vertical and horizontal flips. Following [53] and [57] we set the batch size to 4 for GlaS and MoNuSeg datasets while we set the batch size to 24 for Synapse dataset [7, 57]. For CVC-ClinicDB and Kvasir-SEG datasets, we set the batch size to 16 following [27, 12]. We used Adam optimizer [30] with an initial learning rate of  $10^{-4}$  for all models. We utilized Dice loss [48] as the loss function while we used Dice Similarity Coefficient (DSC) and Intersection over Union (IoU) as the performance metrics. We trained all models for 200 epochs and reported the best results for each model.

#### 3.4. Results

We verify the effectiveness of our proposed DCA block by conducting extensive experiments using six U-Net-based models and five benchmark medical image segmentation datasets. For a fair comparison, we use the same training settings for plain and DCA-integrated models. The overall results are shared in Table 1. As mentioned before,

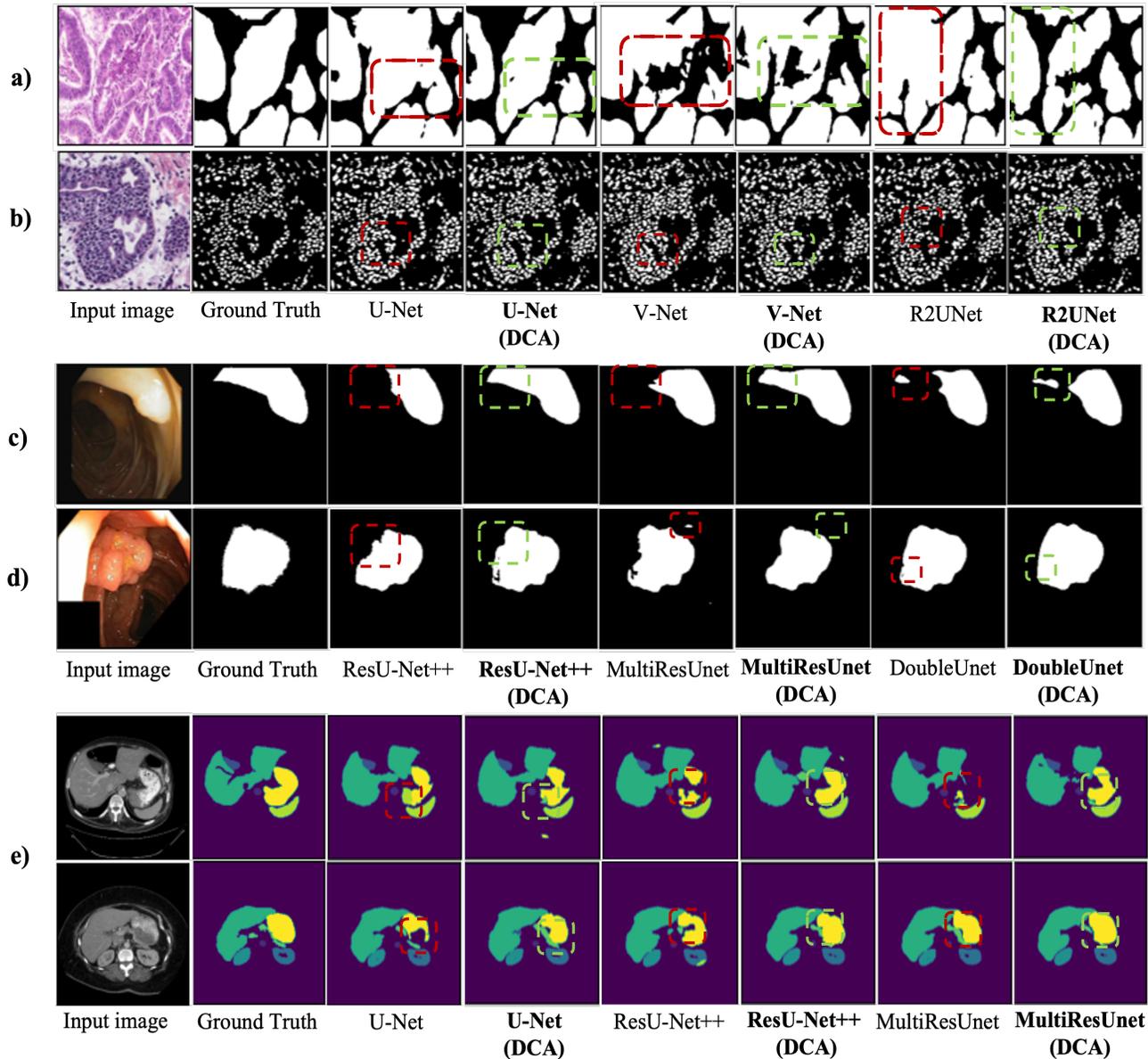


Figure 3. Visual comparison for plain and DCA integrated models. (a) Glas, (b) MoNuSeg, (c) CVC-ClinicDB, (d) Kvasir-Seg, (e) Synapse

we compare each model with their DCA integrated forms using DSC and IoU metrics where bold-faced results denote better performance. As can be seen in Table 1, DCA can significantly improve the segmentation performance at a slight parameter increase for each model. Note that the parameter increase in each model with DCA depends on the total number of skip-connection layers. ResUnet++ consists of three skip-connection layers which bring less than 0.7% parameter increase with DCA while for U-Net, MultiResUnet, R2Unet and V-Net all of which include four skip-connection layers, the parameter increase varies between 0.3% and 1.5%, depending on the model capac-

ity. Since DoubleUnet consists of three skip-connection schemes, each having four skip-connections, the total parameter increase is 3.4% which is still a small increase. Considering such small additional parameters, DCA can provide DSC improvements up to 2.05% on GlaS, 2.74% on MoNuSeg, 1.37% on CVC-ClinicDB, 1.12% on Kvasir-Seg and 1.44% on Synapse datasets which show the effectiveness of our DCA mechanism.

We also visually compare the model predictions to further verify our DCA mechanism. Fig.3 depicts some of the model segmentation predictions. Red dashed rectangles show the regions where plain models struggled to provide

Model	Params	GlaS		MoNuSeg	
		DSC (%)	IoU (%)	DSC (%)	IoU (%)
U-Net	8.64M	88.87	79.98	77.14	62.79
U-Net (CCA)	8.74M	89.07	80.31	77.78	63.64
U-Net (SCA)	8.74M	89.48	80.96	77.36	63.07
U-Net (SCA-CCA)	8.75M	89.03	80.24	77.90	63.80
U-Net (CCA-SCA)	8.75M	<b>89.66</b>	<b>81.29</b>	<b>78.13</b>	<b>64.11</b>

Table 2. Quantitative comparison of different DCA layouts.

Model	Params	GlaS		MoNuSeg	
		DSC (%)	IoU (%)	DSC (%)	IoU (%)
U-Net	8.75M	88.87	79.98	77.14	62.79
U-Net (CCA+SCA)	8.75M	89.09	80.38	77.30	63.00
U-Net (CCA  SCA)	8.75M	89.19	80.52	78.03	63.97
U-Net (CCA-SCA)	8.75M	<b>89.66</b>	<b>81.29</b>	<b>78.13</b>	<b>64.11</b>

Table 3. Quantitative comparison of different fusion strategies for CCA and SCA.

accurate predictions while green dashed rectangles show the improvement of the DCA block in those same regions. As can be seen in Fig.3, models with DCA mechanism outperform the plain models by providing more consistent boundaries and preserving accurate shape information. Besides, models with DCA can better distinguish discrete parts by eliminating false positive predictions.

## 4. Ablation Study

### 4.1. DCA Layout

We start our ablation study by searching for the best layout for the proposed DCA mechanism. We first perform CCA and SCA modules individually. As shown in Table. 2, both CCA and SCA modules outperform U-Net on GlaS dataset by 0.2% and 0.61% DSC improvements, respectively. When we incorporate CCA and SCA modules by CCA-SCA sequence (i.e., CCA first), the performance further improves by 0.79% whereas the SCA-CCA sequence only improves the performance by 0.16%. As for the MoNuSeg dataset, CCA and SCA modules improve U-Net by 0.64% and 0.22%, respectively while the SCA-CCA sequence performs slightly better than both individual CCA and SCA modules. Nonetheless, the dual attention scheme with the CCA-SCA sequence provides the best performance on both datasets, showing that channel and spatial-wise cross-attention mechanisms complement each other.

### 4.2. Fusion of CCA and SCA

As mentioned before, several fusion strategies for dual attention schemes have been proposed and shown to be effective. We conduct experiments using three fusion strategies: (i) fusion by summation (performing CCA and SCA

Model	Params	GlaS		MoNuSeg	
		DSC (%)	IoU (%)	DSC (%)	IoU (%)
U-Net	8.64M	88.87	79.98	77.14	62.79
U-Net (DCA-Conv)	9.01M	89.52	81.07	77.62	63.43
U-Net (DCA-AP)	8.75M	<b>89.66</b>	<b>81.29</b>	<b>78.13</b>	<b>64.11</b>

Table 4. Quantitative comparison of average pooling and convolution for patch embedding.

in a parallel manner and taking the sum of their outputs) [16, 40]; (ii) fusion by concatenation (performing CCA and SCA in a parallel manner and concatenating their outputs) [35]; and (iii) sequential fusion (performing CCA and SCA in a sequential manner) [11]. Table. 3 shows the comparison of those three fusion strategies where +, || and − denote for fusion by summation (i), fusion by concatenation (ii) and sequential fusion (iii), respectively. Although both summation and concatenation fusion strategies improve U-Net on both datasets, the sequential fusion scheme performs the best results.

### 4.3. Average Pooling for Patch Embedding

As for the last part of our ablation study, we compare simple 2D average pooling with convolutional patch embedding. As shown in Table. 4, DCA with convolutional patch embedding strategy, although still significantly improving U-Net, performs slightly worse than DCA with 2D average pooling. Besides, for multi-scale encoder features  $\mathbf{E}_i \in \mathbb{R}^{C_i \times \frac{H}{2^{i-1}} \times \frac{W}{2^{i-1}}}$ , convolutional patch embedding requires kernel sizes ( $P_i^s, P_i^s$ ) where  $P_i^s = \frac{P}{2^{i-1}}$ , ( $i = 1, 2, \dots, n$ ) which brings additional parameters ( $\approx 260K$ ) while 2D average pooling operation is parameter-free and performs better when combined with  $1 \times 1$  depth-wise convolutional projections.

## 5. Conclusion

In this paper, we introduced Dual Cross-Attention (DCA) to strengthen skip-connections in U-Net-based architectures for medical image segmentation. DCA consists of Channel Cross-Attention (CCA) and Spatial Cross-Attention (SCA) modules which sequentially capture long-range dependencies in channel and spatial dimensions, respectively. Besides, DCA utilizes cross-attention in order to effectively fuse low-level multi-scale encoder features and extract fine-grained representations to narrow the semantic gap. Our DCA mechanism is formed with lightweight operations such as 2D average pooling for patch embedding and depth-wise convolutions for projection layers to minimize the computational overhead. Comprehensive experiments using six state-of-the-art U-Net-based architectures and five benchmark medical image segmentation datasets show that DCA can significantly improve the segmentation

performance for models with skip-connections.

## Acknowledgment

We would like to acknowledge the Engineering Cancer Cures funding and IDSC Early Career Researcher Grant at the University of Miami for supporting this project. We also thank the University of Miami Institute for Data Science and Computing for allocating two NVIDIA Tesla V100-SXM2-16GB GPUs for our experiments.

## References

- [1] Md Zahangir Alom, Mahmudul Hasan, Chris Yakopcic, Tarek M Taha, and Vijayan K Asari. Recurrent residual convolutional neural network based on u-net (r2u-net) for medical image segmentation. *arXiv preprint arXiv:1802.06955*, 2018. [1](#), [5](#)
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. [3](#)
- [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. [5](#)
- [4] Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning Wang. Swin-unet: Unet-like pure transformer for medical image segmentation. In *Computer Vision–ECCV 2022 Workshops: Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part III*, pages 205–218. Springer, 2023. [5](#)
- [5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 213–229. Springer, 2020. [2](#)
- [6] 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. [2](#)
- [7] 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. [3](#), [5](#)
- [8] Qiang Chen, Qiman Wu, Jian Wang, Qinghao Hu, Tao Hu, Errui Ding, Jian Cheng, and Jingdong Wang. Mixformer: Mixing features across windows and dimensions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5249–5259, 2022. [3](#)
- [9] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017. [3](#)
- [10] Özgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger. 3d u-net: learning dense volumetric segmentation from sparse annotation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2016: 19th International Conference, Athens, Greece, October 17–21, 2016, Proceedings, Part II 19*, pages 424–432. Springer, 2016. [1](#)
- [11] Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jingdong Wang, and Lu Yuan. Davit: Dual attention vision transformers. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXIV*, pages 74–92. Springer, 2022. [2](#), [3](#), [7](#)
- [12] 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. [5](#)
- [13] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and Baining Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12124–12134, 2022. [2](#)
- [14] 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. [2](#), [3](#)
- [15] Yuxin Fang, Bencheng Liao, Xinggang Wang, Jiemin Fang, Jiyang Qi, Rui Wu, Jianwei Niu, and Wenyu Liu. You only look at one sequence: Rethinking transformer in vision through object detection. *Advances in Neural Information Processing Systems*, 34:26183–26197, 2021. [2](#)
- [16] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3146–3154, 2019. [2](#), [7](#)
- [17] Jiuxiang Gu, Zhenhua Wang, Jason Kuen, Lianyang Ma, Amir Shahroudy, Bing Shuai, Ting Liu, Xingxing Wang, Gang Wang, Jianfei Cai, et al. Recent advances in convolutional neural networks. *Pattern recognition*, 77:354–377, 2018. [1](#)
- [18] Jianyuan Guo, Kai Han, Han Wu, Yehui Tang, Xinghao Chen, Yunhe Wang, and Chang Xu. Cmt: Convolutional neural networks meet vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12175–12185, 2022. [3](#)
- [19] Ruohao Guo, Dantong Niu, Liao Qu, and Zhenbo Li. Sotr: Segmenting objects with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7157–7166, 2021. [2](#)
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [1](#)
- [21] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016. [3](#)

- [22] Han Hu, Zheng Zhang, Zhenda Xie, and Stephen Lin. Local relation networks for image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3464–3473, 2019. 1
- [23] Peijun Hu, Fa Wu, Jialin Peng, Yuanyuan Bao, Feng Chen, and Dexing Kong. Automatic abdominal multi-organ segmentation using deep convolutional neural network and time-implicit level sets. *International journal of computer assisted radiology and surgery*, 12:399–411, 2017. 1
- [24] Zilong Huang, Youcheng Ben, Guozhong Luo, Pei Cheng, Gang Yu, and Bin Fu. Shuffle transformer: Rethinking spatial shuffle for vision transformer. *arXiv preprint arXiv:2106.03650*, 2021. 2
- [25] Nabil Ibtihaz and M Sohel Rahman. Multiresunet: Rethinking the u-net architecture for multimodal biomedical image segmentation. *Neural networks*, 121:74–87, 2020. 1, 5
- [26] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr, 2015. 4
- [27] Debesh Jha, Michael A Riegler, Dag Johansen, Pål Halvorsen, and Håvard D Johansen. Doubleu-net: A deep convolutional neural network for medical image segmentation. In *2020 IEEE 33rd International Symposium on Computer-Based Medical Systems (CBMS)*, pages 558–564. IEEE, 2020. 1, 5
- [28] 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 *MultiMedia Modeling: 26th International Conference, MMM 2020, Daejeon, South Korea, January 5–8, 2020, Proceedings, Part II* 26, pages 451–462. Springer, 2020. 5
- [29] 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. 1, 5
- [30] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [31] Neeraj Kumar, Ruchika Verma, Sanuj Sharma, Surabhi Bhargava, Abhishek Vahadane, and Amit Sethi. A dataset and a technique for generalized nuclear segmentation for computational pathology. *IEEE transactions on medical imaging*, 36(7):1550–1560, 2017. 5
- [32] Youngwan Lee, Jonghee Kim, Jeffrey Willette, and Sung Ju Hwang. Mpvit: Multi-path vision transformer for dense prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7287–7296, 2022. 2, 3
- [33] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019. 2
- [34] Ming Liang and Xiaolin Hu. Recurrent convolutional neural network for object recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3367–3375, 2015. 1
- [35] Xin Liu, Guobao Xiao, Luanyuan Dai, Kun Zeng, Changcai Yang, and Riqing Chen. Scsa-net: Presentation of two-view reliable correspondence learning via spatial-channel self-attention. *Neurocomputing*, 431:137–147, 2021. 2, 7
- [36] 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. 2, 3
- [37] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 1
- [38] Jiasen Lu, Roozbeh Mottaghi, Aniruddha Kembhavi, et al. Container: Context aggregation networks. *Advances in neural information processing systems*, 34:19160–19171, 2021. 3
- [39] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, pages 565–571. Ieee, 2016. 1, 5
- [40] Lei Mou, Yitian Zhao, Huazhu Fu, Yonghuai Liu, Jun Cheng, Yalin Zheng, Pan Su, Jianlong Yang, Li Chen, Alejandro F Frangi, et al. Cs2-net: Deep learning segmentation of curvilinear structures in medical imaging. *Medical image analysis*, 67:101874, 2021. 2, 7
- [41] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 807–814, 2010. 4
- [42] Daniel Neimark, Omri Bar, Maya Zohar, and Dotan Asselmann. Video transformer network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3163–3172, 2021. 2
- [43] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jon Shlens. Stand-alone self-attention in vision models. *Advances in neural information processing systems*, 32, 2019. 2
- [44] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III* 18, pages 234–241. Springer, 2015. 1, 5
- [45] Jo Schlemper, Ozan Oktay, Michiel Schaap, Mattias Heinrich, Bernhard Kainz, Ben Glocker, and Daniel Rueckert. Attention gated networks: Learning to leverage salient regions in medical images. *Medical image analysis*, 53:197–207, 2019. 1
- [46] Korsuk Sirinukunwattana, Josien PW Pluim, Hao Chen, Xiaojuan Qi, Pheng-Ann Heng, Yun Bo Guo, Li Yang Wang, Bogdan J Matuszewski, Elia Bruni, Urko Sanchez, et al. Gland segmentation in colon histology images: The glas challenge contest. *Medical image analysis*, 35:489–502, 2017. 5
- [47] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmenta-

- tion. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7262–7272, 2021. 2
- [48] Carole H Sudre, Wenqi Li, Tom Vercauteren, Sebastien Ourselin, and M Jorge Cardoso. Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: Third International Workshop, DLMIA 2017, and 7th International Workshop, ML-CDS 2017, Held in Conjunction with MICCAI 2017, Québec City, QC, Canada, September 14, Proceedings 3*, pages 240–248. Springer, 2017. 5
- [49] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7464–7473, 2019. 2
- [50] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer: An all-mlp architecture for vision. *Advances in neural information processing systems*, 34:24261–24272, 2021. 3
- [51] Hugo Touvron, Piotr Bojanowski, Mathilde Caron, Matthieu Cord, Alaaeldin El-Nouby, Edouard Grave, Gautier Izacard, Armand Joulin, Gabriel Synnaeve, Jakob Verbeek, et al. Resmlp: Feedforward networks for image classification with data-efficient training. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 3
- [52] 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. 2, 3
- [53] Jeya Maria Jose Valanarasu, Poojan Oza, Ilker Hacihaliloglu, and Vishal M Patel. Medical transformer: Gated axial-attention for medical image segmentation. In *Medical Image Computing and Computer Assisted Intervention—MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part I 24*, pages 36–46. Springer, 2021. 5
- [54] 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. 2, 3
- [55] Cong Wang, Hongmin Xu, Xiong Zhang, Li Wang, Zhitong Zheng, and Haifeng Liu. Convolutional embedding makes hierarchical vision transformer stronger. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XX*, pages 739–756. Springer, 2022. 3
- [56] Changhan Wang, Xinchun Yan, Max Smith, Kanika Kochhar, Marcie Rubin, Stephen M Warren, James Wrobel, and Honglak Lee. A unified framework for automatic wound segmentation and analysis with deep convolutional neural networks. In *2015 37th annual international conference of the IEEE engineering in medicine and biology society (EMBC)*, pages 2415–2418. IEEE, 2015. 1
- [57] 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. 1, 2, 3, 5
- [58] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 568–578, 2021. 2
- [59] 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. 3
- [60] Xinpeng Wang, Chandan Yeshwanth, and Matthias Nießner. Sceneformer: Indoor scene generation with transformers. In *2021 International Conference on 3D Vision (3DV)*, pages 106–115. IEEE, 2021. 2
- [61] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8741–8750, 2021. 2
- [62] Yujing Wang, Yaming Yang, Jiangang Bai, Mingliang Zhang, Jing Bai, Jing Yu, Ce Zhang, Gao Huang, and Yunhai Tong. Evolving attention with residual convolutions. In *International Conference on Machine Learning*, pages 10971–10980. PMLR, 2021. 3
- [63] Yingming Wang, Xiangyu Zhang, Tong Yang, and Jian Sun. Anchor detr: Query design for transformer-based detector. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 2567–2575, 2022. 2
- [64] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22–31, 2021. 3
- [65] Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollár, and Ross Girshick. Early convolutions help transformers see better. *Advances in Neural Information Processing Systems*, 34:30392–30400, 2021. 3
- [66] 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*, 34:12077–12090, 2021. 2
- [67] Weijian Xu, Yifan Xu, Tyler Chang, and Zhuowen Tu. Co-scale conv-attentional image transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9981–9990, 2021. 3
- [68] Fuzhi Yang, Huan Yang, Jianlong Fu, Hongtao Lu, and Bain-ing Guo. Learning texture transformer network for image super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5791–5800, 2020. 2

- [69] Jianwei Yang, Chunyuan Li, and Jianfeng Gao. Focal modulation networks. *arXiv preprint arXiv:2203.11926*, 2022. 3
- [70] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Xiyang Dai, Bin Xiao, Lu Yuan, and Jianfeng Gao. Focal attention for long-range interactions in vision transformers. *Advances in Neural Information Processing Systems*, 34:30008–30022, 2021. 2
- [71] Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10819–10829, 2022. 3
- [72] Weihao Yu, Chenyang Si, Pan Zhou, Mi Luo, Yichen Zhou, Jiashi Feng, Shuicheng Yan, and Xinchao Wang. Metaformer baselines for vision. *arXiv preprint arXiv:2210.13452*, 2022. 3
- [73] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 558–567, 2021. 3
- [74] 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. 2
- [75] 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: 4th International Workshop, DLMIA 2018, and 8th International Workshop, ML-CDS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, Proceedings 4*, pages 3–11. Springer, 2018. 1