# Deep 3D Vessel Segmentation based on Cross Transformer Network

**Chengwei Pan**[1], **Baolian Qi**[2,3], **Gangming Zhao**[4], **Jiaheng Liu**[1], **Chaowei Fang**[5],
**Dingwen Zhang**[6], and **Jinpeng Li**[2,3,*]

[1]Institute of Artificial Intelligence, Beihang University, Beijing, China
[2]HwaMei Hospital, University of Chinese Academy of Sciences (UCAS), Ningbo, China
[3]Ningbo Institute of Life and Health Industry, UCAS, Ningbo, China
[4]Department of Computer Science, University of Hong Kong, Hong Kong, China
[5]School of Artificial Intelligence, Xidian University, Xi'an, China
[6]School of Automation, Northwestern Polytechnical University, Xi'an, China
[*]Email: lijinpeng@ucas.ac.cn

*Abstract*—The coronary microvascular disease poses a great threat to human health. Computer-aided analysis/diagnosis systems help physicians intervene in the disease at early stages, where 3D vessel segmentation is a fundamental step. However, there is a lack of carefully annotated dataset to support algorithm development and evaluation. On the other hand, the commonly-used U-Net structures often yield disconnected and inaccurate segmentation results, especially for small vessel structures. In this paper, motivated by the data scarcity, we first construct two large-scale vessel segmentation datasets consisting of 100 and 500 computed tomography (CT) volumes with pixel-level annotations by experienced radiologists. To enhance the U-Net, we further propose the cross transformer network (CTN) for fine-grained vessel segmentation. In CTN, a transformer module is constructed in parallel to a U-Net to learn long-distance dependencies between different anatomical regions; and these dependencies are communicated to the U-Net at multiple stages to endow it with global awareness. Experimental results on the two in-house datasets indicate that this hybrid model alleviates unexpected disconnections by considering topological information across regions. Our codes, together with the trained models are made publicly available at https://github.com/qibaolian/ctn.

*Index Terms*—Coronary microvascular disease, 3D vessel segmentation, Transformer

## I. Introduction

The coronary microvascular disease is a major threat to human health. In the clinic, the key to reducing its impairment is early intervention. Typically, radiologists seek for diagnostic clues for this disease using computed tomography (CT) scans. In recent years, deep learning algorithms have been developed and applied to help radiologists analyze and diagnose the disease in a more efficient manner, where automatic 3D vessel segmentation is a fundamental step.

Although some excellent works have emerged to tackle this task, some problems remain unsolved. First, there is a lack of carefully annotated datasets to facilitate the proper

development and evaluation of algorithms. To our knowledge, the largest open dataset for coronary vessel segmentation was provided by a MICCAI 2020 challenge[1], which only contained 60 CT volumes. Second, the most popular structure in medical image segmentation, i.e., the U-Net [20] family, although being data-efficient, needs to be enhanced on the fine-grained vessel segmentation task. Existing U-Nets often yield disconnected and inaccurate predictions, especially for small vessel structures. One reason is that the stacked convolutions are hard to capture the global information and yield topology-ignorant results. Most recently, the transformer-based models have been expected to learn long-range dependencies between anatomical regions. However, compared with the U-Net family, they require larger datasets to train due to the lack of inductive bias, i.e., structural constraints.

A feasible 3D vessel segmentation model should be lightweight, not too demanding on data, and being able to learn 3D global information. This motivates us to explore possibles ways to overlay the advantages of U-Net (*data-efficiency*) and transformer (*global-awareness*). In this design, an important topic is how to interconnect the features of the two types of models to achieve concise and elegant information interaction.

To implement such design, we propose the cross transformer network (CTN) for fine-grained vessel segmentation. CTN is a hybrid model consisting of a standard U-Net to model the local contextual information and a parallel stacked transformer to learn long-range dependencies between 3D anatomical regions. The dependencies are communicated to the U-Net encoder at multiple stages to endow it with global awareness. Based on this, the model considers both local and global information when segmenting the vessels.

In summary, this paper contributes in two folds:

- We propose the **C**ross **T**ransformer **N**etwork, a novel method to learn 3D vessel features with global/topology awareness. CTN involves a multi-scale feature interaction

---

[1]https://asoca.grand-challenge.org/

between the U-Net and transformer modules, endowing the U-Net with global awareness to confront with disconnections and inaccurate segmentation.

- We evaluate CTN and existing methods on our two large-scale vessel segmentation datasets. Extensive experiments indicate that CTN achieves state-of-the-art 3D vessel segmentation performance. Our codes and pre-trained models are made publicly available to the community.

## II. RELATED WORK

### A. Coronary Artery Disease

In the last three years, deep learning algorithms have been intensively applied to the automatic diagnosis of coronary artery disease. For example, Denzinger et al. [7] utilized radiomics, deep learning and their combination to assess the coronary artery plaque segments and the results on 345 plague segments indicated that the combination of shape-, intensity- and texture-based radiomic features with 2D CNN yielded the best result. Denzinger et al. [6] further extended the dataset and achieved an accuracy of 0.91 for determining whether the patient suffers from coronary artery disease using a 2D CNN. They used 2D CNN instead of 3D CNN to avoid overfitting. He et al. [11] proposed a hybrid learning algorithm to extract both local and global information, which can guide the automatic extraction of vessel centerlines and confront the discontinuities caused by the segmentation based on local features. Ma et al. [16] sequentially used 3D CNNs and transform modules to capture local and global information respectively. Then a simple classifier is introduced to predict significant stenosis. In all these works, the centerline extraction is a vital procedure. A well-segmented coronary artery is essential to obtain a reliable centerline and can be further used to analyze vessel morphology to determine angiomas, stenosis, occlusions, etc.

### B. CNN-based Segmentation Networks

Since the introduction of U-Net [20], fully convolutional neural networks (CNNs) have been the predominate approach on various 2D and 3D medical image segmentation tasks [8] [29] [31] [35]. A vast number of works are dedicated to retrofit and improve the architecture of U-Net, resulting in many variants such as V-Net [18], 3D U-Net [3], Res-UNet [27], Dense-UNet [13], Y-Net [17] and U-Net++ [34]. These methods achieve impressive performance on many tasks, demonstrating the effectiveness of CNN in learning the discriminate features to segment organs or lesions from medical images. Despite their effectiveness, CNNs still suffer from the limited receptive field and lack the ability to capture the long-range (global) dependencies, due to the inductive bias of locality and weight sharing [4]. Many efforts have been devoted to enlarging the receptive field of CNN with image pyramids [32], atrous convolutions [2] and attention mechanisms [26].

### C. Vision Transformers

Transformer architectures using the self-attention mechanism are the most popular methods in natural language processing (NLP) due to their excellent ability of modeling long-range dependencies [24]. Inspired by this, vision transformers have recently gained traction and achieved competitive performance to CNNs on many computer vision tasks, such as image recognition [12], object detection [5], video recognition [10] and semantic segmentation [33]. For medical image segmentation, Chen et al. [1] designed the TransUNet for multi-organ segmentation by embedding a transformer as an additional layer in the bottleneck of a U-Net network. Gao et al. [9] proposed a hybrid transformer architecture named UTNet, in which self-attention modules are integrated into both encoders and decoders for capturing long-range dependency at multiple scales. Xie et al. [28] introduced a hybrid model that bridges a CNN and a deformable transformer for 3D medical image segmentation. Similarly, Wang et al. [25] proposed a 3d version of TransUNet, in which a transformer is employed in the bottleneck of a 3D U-Net architecture for the task of brain tumor segmentation. It is noteworthy that the vision transformers are based on the attention computation and are not specifically designed for the structure of the input data; therefore, a large amount of data is generally required for vision transformers to learn the inductive bias.

## III. METHOD

### A. Method

The architecture of the proposed method is presented in Figure 1. The CT image volume $X \in \mathbb{R}^{D \times H \times W}$ is fed into two branches at the same time, where $D$, $H$, and $W$ represent the spatial depth, height, and width, respectively. The 3D U-Net is the backbone, which includes an encoder and a decoder network. The 3D Swin Transformer acts as a feature extractor to learn the long-term dependencies. In our design, there are four multi-scale fusions of U-Net feature map $F_u$ and Swin Transformer feature map $F_s$. The first two fusions are designed to integrate the coarse-grained representations of the two branches, and the remaining fusions are designed to integrate the fine-grained representations. Note that the size of the reshaped $F_s$ is the same as $F_u$. The aim of this framework is to utilize the long-range contextual information to improve the performance of the vessel segmentation.

### B. 3D Shifted Windows

For the input images $X \in \mathbb{R}^{256 \times 256 \times 256}$, the windows are partitioned into non-overlapped 3D patches $X'$ with the size of $4 \times 4 \times 4$, and the 3D tokens of $\frac{D}{4} \times \frac{W}{4} \times \frac{H}{4}$ are thus obtained. Due to the large size of CT images, the tiny version of 3D Swin Transformer is adopted in our method. The channel numbers of the hidden layers of four stages are $\{48, 96, 192, 384\}$ and the layer numbers of four stages are $\{2, 2, 6, 2\}$. Compared to the 2D models, the 3D shifted window based multi-head self-attention is exploited in place of the multi-head self-attention (MSA) module followed by a feed-forward network (FFN) and the other components keep unchanged to those in the 3D Swin Transformer [15].
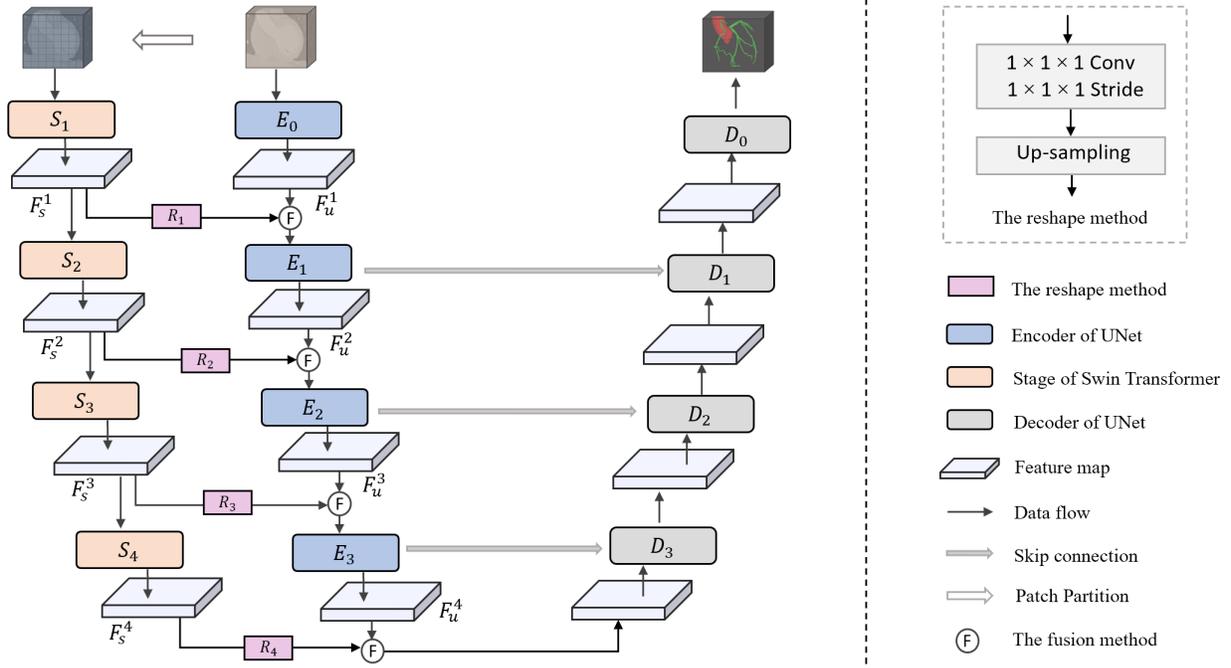
Fig. 1. Flowchart of the proposed CTN. CT volumes are fed into the 3D U-Net and 3D Swin Transformer to learn local and global features simultaneously. To achieve effective feature interactions, there are four fusions of the U-Net feature map $F_u$ and the Swin Transformer feature map $F_s$, resulting in coase-to-fine feature fusions between the two parallel companions.

To reduce computation, the 3D window based MSA (3DW-MSA) module only conducts self-attention within local windows and the computation is implemented as follows.

$$
\begin{aligned}
\hat{z}^l &= 3DW\text{-}MSA(LN(z^{l-1})) + z^{l-1}, \\
z^l &= FFN(LN(\hat{z}^l)) + \hat{z}^l,
\end{aligned}
\tag{1}
$$

where $\hat{z}^l$ and $z^l$ are the outputs of the 3DW-MSA module and the FFN module in the $l^{th}$ layer. LN denotes layer normalization.

The 3DW-MSA module lacks the information interaction between adjacent windows. Thus, the 3D shifted window based MSA (3DSW-MSA) module is introduced followed by the 3DW-MSA to enlarge the range of dependencies modeled, which is beneficial to the learning of global features. The details of computation are defined as follows.

$$
\begin{aligned}
\hat{z}^{l+1} &= 3DSW\text{-}MSA(LN(z^l)) + z^{l-1}, \\
z^{l+1} &= FFN(LN(\hat{z}^{l+1})) + \hat{z}^{l+1},
\end{aligned}
\tag{2}
$$

where $\hat{z}^{l+1}$ and $z^{l+1}$ are the outputs of the 3DSW-MSA module and the FFN module in the $(l+1)^{th}$ layer. With the shifted windows, the 3D Swin Transformer can better extract feature associations with low computational costs.

### C. Fusions in Cross Transformer

The 3D Swin Transformer is used as an auxiliary encoder of the U-Net to extract long-range dependencies and global contextual connections. After obtaining the output features of the two encoders, it is very important to fuse them effectively and efficiently. There are two fusion methods. One is the adding operation of the feature maps $F_s$ and $F_u$, i.e.,

$$
F_{f_+}^{ij} = F_u^i + Reshape(F_s^j),
\tag{3}
$$

where $F_u^i$ is the feature map of the U-Net in the $i^{th}$ encoder and $F_s^j$ is the feature map of the Swin Transformer in the $j^{th}$ stage. $F_{f_+}^{ij}$ is the feature map after fusion. The size of feature map $F_s^j$ is reshaped to the same as the feature map $F_u^i$ using the resizing strategies and the $1 \times 1$ convolution.

The other method is concatenating the feature maps $F_s$ and $F_u$ and then performing the convolution operation, i.e.,

$$
\begin{aligned}
F^{ij} &= Concat(F_u^i, Reshape(F_s^j)) \\
F_{f_{cat}}^{ij} &= Conv(F^{ij}).
\end{aligned}
\tag{4}
$$

## IV. EXPERIMENTS

### A. Dataset

Automated Segmentation of Aorta and Coronary Artery (ASACA) is a large in-house dataset for evaluating the performance of vessel segmentation. Figure 2 shows some examples. Each CT image is annotated by one radiologist and verified by another. All the data described in this paper is derived from studies that have received appropriate approvals from institutional ethics committees. The ASACA contains two datasets named ASACA100 and ASACA500, where the main difference is the number of CT images. The ASACA100 consists of a total of 100 coronary computed tomography angiography (CCTA) images, including a training set of 80 CCTA images and a test set of 20 CCTA images. The
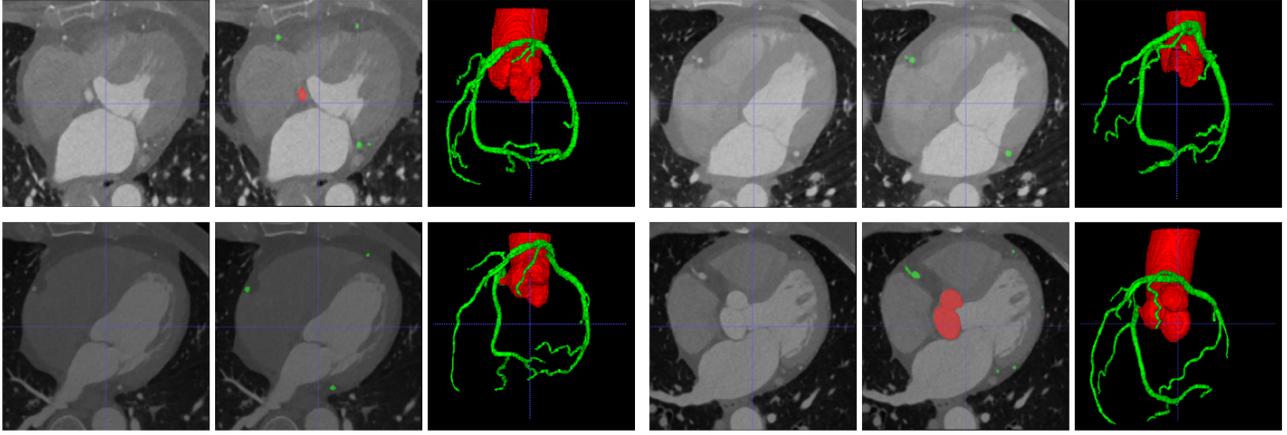
Fig. 2. Visualization of the ASACA dataset. In each instance, the first column shows the original CT axial images; the second and third columns show images covered with vessel labels and the 3D vessels, respectively. The red and green area mean the aorta and coronary artery, respectively. The aorta is relatively easy to segment, whereas the coronary artery is difficult to segment due to its high diversity and complexity in shape and appearance.

ASACA500 consists of a total of 500 CCTA images, including a training set of 400 CCTA images, a validation set of 50 CCTA images, and a test set of 50 CCTA images. Please note that all CT images in ASACA are resized to $256 \times 256 \times 256$.

### B. Experimental Setting

The model is trained on ASACA500 and ASACA100 using the AdamW optimizer. Within a total of 200 epochs, the learning rate starts from 0.0001 and decreases by 10 times after every 50 epochs. The weight decay and the momentum are set to 0.0001 and 0.9, respectively. The mini-batch size is set to 1 because of the large volume size of the whole 3D CT and the limited GPU memory of a single NVIDIA A100 GPU. The proposed method and all the comparing methods are implemented with PyTorch [19].

### C. Evaluation Metrics

The dice coefficient (DICE) and average symmetric surface distance (ASSD) are used to evaluate the performance of models, which is commonly used in medical image segmentation [30]. Note that the $DICE_A$ denotes the Dice coefficient of aorta and the $DICE_C$ denotes the dice coefficient of coronary artery.

$$DICE = \frac{2|S \cap G|}{|S| + |G|}, \quad (5)$$

where $S$ denotes the segmentation result and $G$ denotes the ground truth. ASSD measures the average symmetrical distance between the segmentation result and the ground truth:

$$
\begin{aligned}
d(G, S) &= \min_{s \in T(S)} ||g - s||, d(S, G) = \min_{g \in T(G)} ||s - g||, \\
N &= |T(S)| + |T(G)|, \\
ASSD &= \frac{1}{N}\left( \sum_{g \in T(G)} d(G, S) + \sum_{s \in T(S)} d(S, G) \right),
\end{aligned}
\quad (6)
$$

where $T(S)$ and $T(G)$ denote the set of surface voxels of $S$ and $G$, respectively. $s$ and $g$ are a surface voxel of $T(S)$ and

$T(G)$, respectively. $S$ denotes the segmentation result and $G$ denotes the ground truth.

Additionally, the metric of skeleton recall (SR) and skeleton precision (SP) are used to evaluate the tubular structure of coronary artery, which preserve more accurate connectivity of coronary arteries.

$$SR(S, G) = \frac{|S \cap Q(G)|}{|Q(G)|}, SP(S, G) = \frac{|G \cap Q(S)|}{|Q(S)|}, \quad (7)$$

where $S$ and $G$ are the segmentation result and the ground truth, respectively. Q(*) is the skeletonization function [23], which has been used to preserve the original vascular topology and connectivity.

### D. Comparison with the Reference Model

To evaluate the effectiveness of the proposed model in a fair manner, we reproduce and compare six representative methods on ASACA500 dataset and ASACA100 dataset under the same data partitioning and experimental setup. The six compared methods consist of the 3D U-Net [3], 3D Swin Transformer [15], clDice [22], TransUNet [1], UTNet [9] and CoTr [28]. The 3D U-Net [3] is the 3D version of U-Net [21], where the 2D operations of which are all replaced by their 3D counterparts for volumetric segmentation. 3D Swin Transformer [15] is the 3D version of Swin Transformer [14], which extends the scope of local attention computation from only the spatial domain to the spatiotemporal domain. The clDice [22] is the state-of-the-art method for tubular structure segmentation, which guarantees topology preservation up to homotopy equivalence for binary 2D and 3D segmentation. The TransUNet [1] embeds a transformer as an additional layer in the bottleneck of a U-Net network for multi-organ segmentation. UTNet [9] is a hybrid transformer architecture, in which self-attention modules are integrated in both encoders and decoders for capturing long-range dependency at multiple scales. CoTr [28] is a hybrid model that bridges a CNN and a deformable transformer for 3D medical image segmentation.

TABLE I
THE PERFORMANCE COMPARISON OF VESSEL SEGMENTATION AMONG THE MODELS USING ASACA500 AND ASACA100 DATASETS. NOTE THAT MM DENOTES MILLIMETERS AND $_{(128)}$ DENOTES THE INPUT RESOLUTION OF $(128 \times 128 \times 128)$.

| Dataset | Models | DICE(%) | $DICE_A$(%) | $DICE_C$(%) | ASSD(mm) | SP(%) | SR(%) |
|---|---|---|---|---|---|---|---|
| ASACA500 | UNet [3] | 91.57 | 98.03 | 85.15 | 0.479 | 95.43 | 91.82 |
| | Swin Transformer [15] | 89.61 | 97.41 | 82.14 | 0.633 | 94.44 | 87.88 |
| | clDice [22] | 91.71 | 98.04 | 85.46 | 0.466 | 95.60 | 92.40 |
| | TransUNet [1] | 90.89 | 97.86 | 83.96 | 0.546 | 96.49 | 88.39 |
| | UTNet$_{(128)}$ [9] | 86.50 | 98.01 | 85.37 | 0.497 | 94.67 | 92.25 |
| | CoTr$_{(128)}$ [28] | 84.87 | 97.68 | 83.72 | 0.605 | 95.73 | 88.43 |
| | Ours | **91.89** | **98.07** | **85.79** | **0.429** | **95.75** | **93.20** |
| ASACA100 | UNet [3] | 89.93 | 97.13 | 83.21 | 0.616 | 94.00 | 90.35 |
| | Swin Transformer [15] | 86.43 | 94.72 | 78.15 | 1.299 | 94.23 | 78.38 |
| | clDice [22] | 90.27 | 97.26 | 83.36 | 0.656 | **94.97** | 89.30 |
| | TransUNet [1] | 88.48 | 96.50 | 80.75 | 0.994 | 93.93 | 84.28 |
| | UTNet$_{(128)}$ [9] | 83.48 | 97.17 | 84.38 | 0.663 | 94.00 | 91.24 |
| | CoTr$_{(128)}$ [28] | 77.85 | 96.32 | 80.03 | 1.016 | 91.10 | 83.89 |
| | Ours | **91.01** | **97.33** | **84.93** | **0.499** | 94.77 | **93.14** |

TABLE II
THE PERFORMANCE COMPARISON OF DIFFERENT FUSION METHODS USING ASACA500 AND ASACA100 DATASET. NOTE THAT MM DENOTES MILLIMETERS.

| Dataset | Models | DICE(%) | $DICE_A$(%) | $DICE_C$(%) | ASSD(mm) | SP(%) | SR(%) |
|---|---|---|---|---|---|---|---|
| ASACA500 | $F_{f_{cat}}$ | 91.87 | 98.04 | **85.85** | 0.433 | **95.77** | **93.35** |
| | $F_{f_+}$ | **91.89** | **98.07** | 85.79 | **0.429** | 95.75 | 93.20 |
| ASACA100 | $F_{f_{cat}}$ | 90.54 | 97.29 | 84.08 | 0.597 | 94.32 | 92.07 |
| | $F_{f_+}$ | **91.01** | **97.33** | **84.93** | **0.499** | **94.77** | **93.14** |

Especially, UTNet and CoTr have a low input resolution $(128 \times 128 \times 128)$ because of the GPU memory constraint, and the other methods have a default input resolution of $256 \times 256 \times 256$ as described in the Section 3.1.

The experimental results in Table I show that our model achieves state-of-the-art performance in six evaluation metrics on the two datasets. For example, our model achieves the DICE of 91.89%, the ASSD of 0.429mm, the SP of 95.75% and the SR of 93.20% on ASACA500 dataset, and the DICE outperforms U-Net, Swin transformer, clDice, TransUNet, UTNet, and CoTr by 0.32%, 2.28%, 0.18%, 1.00%, 5.39%, and 7.02%, respectively. It is noteworthy that UTNet and CoTr achieve poor performance because of low input resolution compared with the other methods. Additionally, since aorta is relatively easy to segment and contributes high performance, the Dice of aorta ($DICE_A$) and the DICE of coronary artery ($DICE_C$) are reported separately. Overall, the experimental results demonstrate that the proposed CTN can efficiently represent sparse and anisotropic vessel structures and has a good performance for vessel segmentation.

## V. ABLATION STUDY

### A. Ablation of Fusion Method

We investigate the effectiveness of the two fusion methods, including the addition operation $F_{f_+}$ and the concatenation operation $F_{f_{cat}}$. As is shown in Table II, the performance of the fusion method $F_{f_+}$ is better than that of the fusion method $F_{f_{cat}}$ on ASACA500 and ASACA100 dataset. Especially, the fusion method $F_{f_+}$ has more obvious advantages on ASACA100 dataset. For example, $F_{f_+}$ outperforms $F_{f_{cat}}$ by 0.02% (DICE) on ASACA500 dataset, but $F_{f_+}$ outperforms $F_{f_{cat}}$ by 0.47% (DICE) on ASACA100 dataset. The reason may be that the concatenation operation requires more data to learn a good weight of feature vectors, while the addition operation does not. Therefore, the fusion method $F_{f_+}$ is used to fuse features in our experiment.

### B. Ablation of Fusion Stages

We also investigate the effectiveness of different fusion stages. We discard different fusion stages and evaluate the effect of each fusion stage. According to Table III, the models with different fusion stages demonstrate more advantages over the baseline model and the model with all fusion stages achieves the best performance. It can be seen that all fusion stages play important roles to improve segmentation performance. Moreover, removing $F_f^{33}$ leads to 0.15% (DICE) performance drop and removing $F_f^{44}$ leads to 0.18% (DICE) performance drop on ASACA100 dataset. Additionally, if the fusion of $F_f^{33}$ and $F_f^{44}$ is discarded, the model suffers from more performance drop. Therefore, $F_f^{22}$, $F_f^{33}$, and $F_f^{44}$ are the most important communication in the model. It demonstrates that the extraction of global information is necessary for vessel

TABLE III

THE PERFORMANCE COMPARISON OF DIFFERENT FUSION STAGES USING ASACA500 AND ASACA100 DATASET. NOTE THAT MM DENOTES MILLIMETERS.

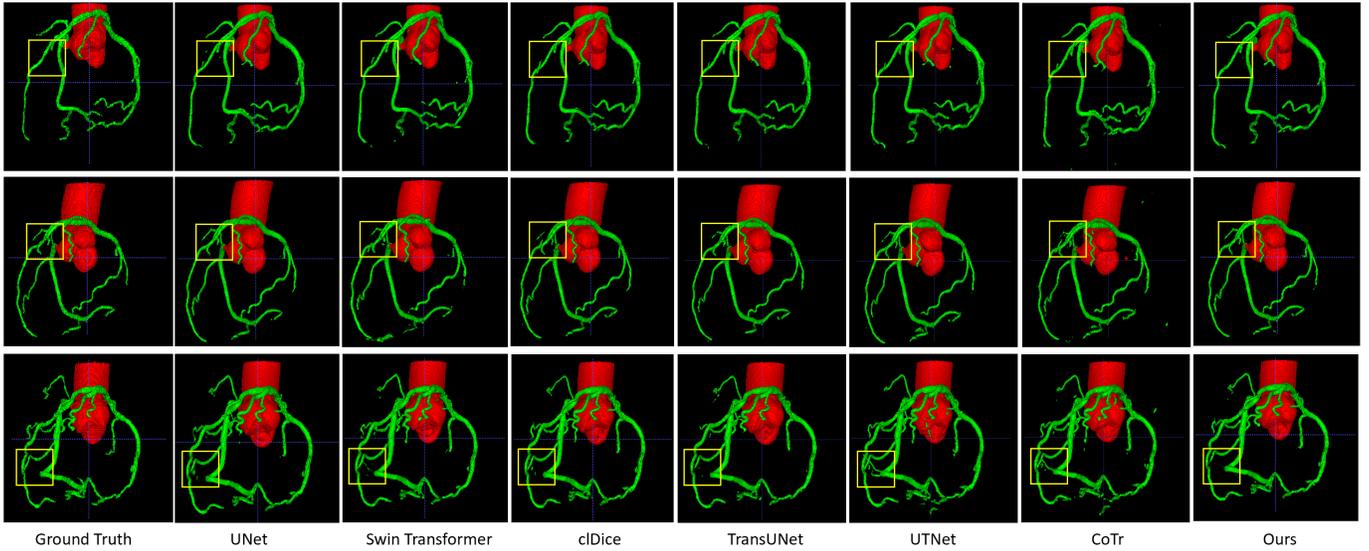| Dataset | $F_f^{11}$ | $F_f^{22}$ | $F_f^{33}$ | $F_f^{44}$ | DICE(%) | DICE$_A$(%) | DICE$_C$(%) | ASSD(mm) | SP(%) | SR(%) |
|---|---|---|---|---|---|---|---|---|---|---|
| | ⊠ | ⊠ | ⊠ | ⊠ | 91.57 | 98.03 | 85.15 | 0.479 | 95.43 | 91.82 |
| | ✓ | ✓ | ✓ | ⊠ | 91.79 | 98.04 | 85.63 | 0.430 | 95.73 | 93.44 |
| ASACA500 | ✓ | ✓ | ⊠ | ✓ | 91.80 | 98.04 | 85.76 | 0.426 | 94.83 | 94.60 |
| | ✓ | ⊠ | ✓ | ✓ | 91.77 | 98.04 | 85.59 | 0.444 | 96.35 | 92.15 |
| | ⊠ | ✓ | ✓ | ✓ | 91.83 | 98.07 | 85.65 | 0.440 | 96.59 | 92.42 |
| | ✓ | ✓ | ✓ | ✓ | 91.89 | 98.07 | 85.79 | 0.429 | 95.75 | 93.20 |
| | ⊠ | ⊠ | ⊠ | ⊠ | 89.93 | 97.13 | 83.21 | 0.616 | 94.00 | 90.35 |
| | ✓ | ✓ | ✓ | ⊠ | 90.83 | 97.30 | 84.74 | 0.535 | 93.29 | 94.22 |
| ASACA100 | ✓ | ✓ | ⊠ | ✓ | 90.85 | 97.20 | 84.62 | 0.530 | 94.72 | 92.11 |
| | ✓ | ⊠ | ✓ | ✓ | 90.87 | 97.18 | 84.62 | 0.538 | 95.69 | 91.33 |
| | ⊠ | ✓ | ✓ | ✓ | 90.87 | 97.31 | 84.56 | 0.535 | 94.27 | 92.81 |
| | ✓ | ✓ | ✓ | ✓ | 91.01 | 97.33 | 84.93 | 0.499 | 94.77 | 93.14 |



Fig. 3. Visualization of the vessel segmentation of different models, including the ground truth, the results of U-Net, Swin Transformer, clDice, TransUNet, UTNet, CoTr, and our model (from left to right). The red and green areas mean the aorta and the coronary vessels. The proposed CTN makes more accurate and continuous predictions, which is in consistent with the quantitative results.

segmentation, which is consistent with our motivation to take advantage of the ability of transformer modules to capture long-range dependencies. These observations are inconsistent with our motivation.

*C. Visualization*

We visualize some typical predictions of our model and the compared methods in Figure 3. The proposed CTN predicts more accurate segmentation results, and other models achieve relatively low performance due to over-segmentation or under-segmentation in high-density regions. It demonstrates that the method of the multi-scale feature interaction between the U-Net and transformer modules can further boost the segmentation performance.

## VI. CONCLUSION

We exploit and explore the CTN to better segment coronary arteries by addressing the insufficiency of existing 3D U-Nets, i.e., the neglect of global structural information. In practice, this can partially lead to discontinuous and inaccurate segmentation results. Comparison experiments and ablation experiments on two largest datasets to date demonstrate that our hybrid model is superior to state-of-the-art. The CTN is designed to better take into account global topological information to combat disconnections and inaccurate segmentation results. Qualitative and quantitative results indicate that the global features learned by transformers are able to compensate for the deficiency of the U-Net structure, i.e., the topology-ignorance. The long distances encourage the continuity of tiny vessel structures. Our implementation will serve as a strong baseline in the CT-based vessel segmentation. In future

work, we will explore alternative ways of interacting global-local information between parallel models and investigate lightweight model implementations.

REFERENCES

[1] J. Chen, Y. Lu, Q. Yu, X. Luo, E. Adeli, Y. Wang, L. Lu, A. L. Yuille, and Y. Zhou, "Transunet: Transformers make strong encoders for medical image segmentation," *arXiv preprint arXiv:2102.04306*, 2021.

[2] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Semantic image segmentation with deep convolutional nets and fully connected crfs," *arXiv preprint arXiv:1412.7062*, 2014.

[3] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, "3d u-net: learning dense volumetric segmentation from sparse annotation," in *International conference on medical image computing and computer-assisted intervention*. Springer, 2016, pp. 424–432.

[4] N. Cohen and A. Shashua, "Inductive bias of deep convolutional networks through pooling geometry," *arXiv preprint arXiv:1605.06743*, 2016.

[5] X. Dai, Y. Chen, J. Yang, P. Zhang, L. Yuan, and L. Zhang, "Dynamic detr: End-to-end object detection with dynamic attention," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 2988–2997.

[6] F. Denzinger, M. Wels, K. Breininger, M. A. Gülsün, M. Schöbinger, F. André, S. Buß, J. Görich, M. Sühling, and A. Maier, "Automatic cadrads scoring using deep learning," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2020, pp. 45–54.

[7] F. Denzinger, M. Wels, N. Ravikumar, K. Breininger, A. Reidelshöfer, J. Eckert, M. Sühling, A. Schmermund, and A. Maier, "Coronary artery plaque characterization from ccta scans using deep learning and radiomics," *medical image computing and computer assisted intervention*, 2019.

[8] Q. Dou, H. Chen, Y. Jin, L. Yu, J. Qin, and P.-A. Heng, "3d deeply supervised network for automatic liver segmentation from ct volumes," in *International conference on medical image computing and computer-assisted intervention*. Springer, 2016, pp. 149–157.

[9] Y. Gao, M. Zhou, and D. N. Metaxas, "Utnet: a hybrid transformer architecture for medical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021, pp. 61–71.

[10] R. Girdhar, J. Carreira, C. Doersch, and A. Zisserman, "Video action transformer network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 244–253.

[11] J. He, C. Pan, C. Yang, M. Zhang, Y. Wang, X. Zhou, and Y. Yu, "Learning hybrid representations for automatic 3d vessel centerline extraction," *medical image computing and computer-assisted intervention*, 2020.

[12] A. Kolesnikov, A. Dosovitskiy, D. Weissenborn, G. Heigold, J. Uszkoreit, L. Beyer, M. Minderer, M. Dehghani, N. Houlsby, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," 2021.

[13] X. Li, H. Chen, X. Qi, Q. Dou, C.-W. Fu, and P.-A. Heng, "H-denseunet: hybrid densely connected unet for liver and tumor segmentation from ct volumes," *IEEE transactions on medical imaging*, vol. 37, no. 12, pp. 2663–2674, 2018.

[14] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows." *arXiv: Computer Vision and Pattern Recognition*, 2021.

[15] Z. Liu, J. Ning, Y. Cao, Y. Wei, Z. Zhang, S. Lin, and H. Hu, "Video swin transformer," 2021.

[16] Y. Ma, Y. Hua, H. Deng, T. Song, H. Wang, Z. Xue, H. Cao, R. Ma, and H. Guan, "Self-supervised vessel segmentation via adversarial learning," *international conference on computer vision*, 2021.

[17] S. Mehta, E. Mercan, J. Bartlett, D. Weaver, J. G. Elmore, and L. Shapiro, "Y-net: joint segmentation and classification for diagnosis of breast biopsy images," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2018, pp. 893–901.

[18] F. Milletari, N. Navab, and S.-A. Ahmadi, "V-net: Fully convolutional neural networks for volumetric medical image segmentation," in *2016 fourth international conference on 3D vision (3DV)*. IEEE, 2016, pp. 565–571.

[19] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," 2017.

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

[21] ——, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention*, 2015.

[22] S. Shit, J. C. Paetzold, A. Sekuboyina, I. Ezhov, A. Unger, A. Zhylka, J. P. Pluim, U. Bauer, and B. H. Menze, "cldice-a novel topology-preserving loss function for tubular structure segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 16 560–16 569.

[23] S. Van der Walt, J. L. Schönberger, J. Nunez-Iglesias, F. Boulogne, J. D. Warner, N. Yager, E. Gouillart, and T. Yu, "scikit-image: image processing in python," *PeerJ*, vol. 2, p. e453, 2014.

[24] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[25] W. Wang, C. Chen, M. Ding, H. Yu, S. Zha, and J. Li, "Transbts: Multimodal brain tumor segmentation using transformer," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021, pp. 109–119.

[26] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7794–7803.

[27] X. Xiao, S. Lian, Z. Luo, and S. Li, "Weighted res-unet for high-quality retina vessel segmentation," in *2018 9th international conference on information technology in medicine and education (ITME)*. IEEE, 2018, pp. 327–331.

[28] Y. Xie, J. Zhang, C. Shen, and Y. Xia, "Cotr: Efficiently bridging cnn and transformer for 3d medical image segmentation," in *International conference on medical image computing and computer-assisted intervention*. Springer, 2021, pp. 171–180.

[29] L. Yu, X. Yang, H. Chen, J. Qin, and P. A. Heng, "Volumetric convnets with mixed residual connections for automated prostate segmentation from 3d mr images," in *Thirty-first AAAI conference on artificial intelligence*, 2017.

[30] Q. Yue, X. Luo, Q. Ye, L. Xu, and X. Zhuang, "Cardiac segmentation from lge mri using deep neural network incorporating shape and spatial priors," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2019, pp. 559–567.

[31] L. Zhang, J. Zhang, P. Shen, G. Zhu, P. Li, X. Lu, H. Zhang, S. A. Shah, and M. Bennamoun, "Block level skip connections across cascaded v-net for multi-organ segmentation," *IEEE Transactions on Medical Imaging*, vol. 39, no. 9, pp. 2782–2793, 2020.

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

[33] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. 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*, 2021, pp. 6881–6890.

[34] Z. Zhou, M. M. Rahman Siddiquee, N. Tajbakhsh, and J. 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*. Springer, 2018, pp. 3–11.

[35] Q. Zhu, B. Du, B. Turkbey, P. L. Choyke, and P. Yan, "Deeply-supervised cnn for prostate segmentation," in *2017 international joint conference on neural networks (IJCNN)*. IEEE, 2017, pp. 178–184.