

# CoTr: Efficiently Bridging CNN and Transformer for 3D Medical Image Segmentation

Yutong Xie<sup>1,2\*</sup>, Jianpeng Zhang<sup>1,2\*</sup>, Chunhua Shen<sup>2</sup>,  
Yong Xia<sup>✉1</sup>

<sup>1</sup> Northwestern Polytechnical University, China

<sup>2</sup> The University of Adelaide, Australia

**Abstract.** Convolutional neural networks (CNNs) have been the de facto standard for nowadays 3D medical image segmentation. The convolutional operations used in these networks, however, inevitably have limitations in modeling the long-range dependency due to their inductive bias of locality and weight sharing. Although Transformer was born to address this issue, it suffers from extreme computational and spatial complexities in processing high-resolution 3D feature maps. In this paper, we propose a novel framework that efficiently bridges a **C**onvolutional neural network and a **T**ransformer (**CoTr**) for accurate 3D medical image segmentation. Under this framework, the CNN is constructed to extract feature representations and an efficient deformable Transformer (DeTrans) is built to model the long-range dependency on the extracted feature maps. Different from the vanilla Transformer which treats all image positions equally, our DeTrans pays attention only to a small set of key positions by introducing the deformable self-attention mechanism. Thus, the computational and spatial complexities of DeTrans have been greatly reduced, making it possible to process the multi-scale and high-resolution feature maps, which are usually of paramount importance for image segmentation. We conduct an extensive evaluation on the Multi-Atlas Labeling Beyond the Cranial Vault (BCV) dataset that covers 11 major human organs. The results indicate that our CoTr leads to a substantial performance improvement over other CNN-based, transformer-based, and hybrid methods on the 3D multi-organ segmentation task. Code is available at <https://github.com/YtongXie/CoTr>

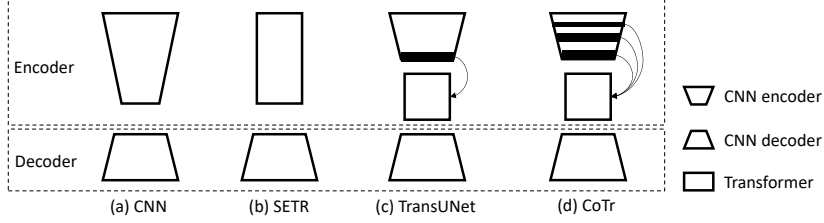
**Keywords:** 3D Medical image segmentation · Deformable self-attention · CNN · Transformer.

## 1 Introduction

Image segmentation is a longstanding challenge in medical image analysis. Since the introduction of U-Net [16], fully convolutional neural networks (CNNs) have become the predominant approach to addressing this task [10, 12, 23–25, 28]. Despite their prevalence, CNNs still suffer from the limited receptive field and fail

---

\* Y. Xie and J. Zhang contributed equally to this work.

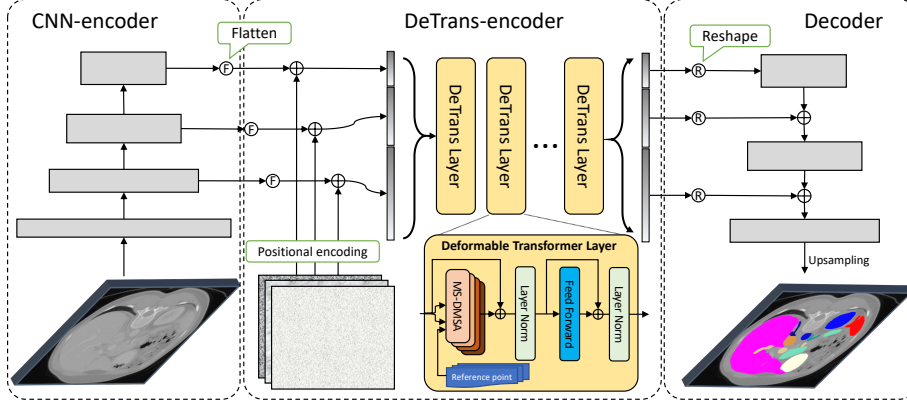


**Fig. 1.** Comparison of different segmentation architectures. All of them have an encoder-decoder structure, but with different encoders. The encoder in CNN (a) is composed of multiple stacked convolutional layers. The encoder in SETR (b) is purely formed from self-attention layers, *i.e.*, Transformer. The encoder in both TransUNet (c) and our proposed CoTr (d) are the hybrid of CNN and Transformer. Differently, TransUNet only processes the low-resolution feature maps from the last stage due to the high computation and spatial complexities. Thanks to the efficient design of Transformer, CoTr is able to process the multi-scale and high-resolution feature maps.

to capture the long-range dependency, due to the inductive bias of locality and weight sharing [6]. Many efforts have been devoted to enlarge a CNN’s receptive field thus improve its ability to context modeling. Yu *et al.* [22] proposed the atrous convolution with an adjustable dilated rate, which shows superior performance in semantic segmentation [5]. More straightforwardly, Peng *et al.* [15] designed large kernels to capture rich global context information. Zhao *et al.* [26] employed the pyramid pooling at multiple feature scales to aggregate multi-scale global information. Wang *et al.* [20] presented the non-local operations which is usually embedded at the end of encoder to capture the long-range dependency. Although improving the context modeling to some extent, these models still have an inevitably limited receptive field, stranded by the CNN architecture.

Transformer, a sequence-to-sequence prediction framework, has a proven track record in machine translation and nature language processing [8, 19], due to its strong ability to long-range modeling. The self-attention mechanism in Transformer can dynamically adjust the receptive field according to the input content, and hence is superior to convolutional operations in modeling the long-range dependency.

Recently, Transformer has been considered as an alternative architecture, and has achieved competitive performance on many computer vision tasks, like image recognition [9, 17], semantic/instance segmentation [21, 27], object detection [2, 29], low-level vision [3, 14], and image generation [13]. A typical example is the vision Transformer (ViT) [9], which outperforms a ResNet-based CNN on recognition tasks but at a cost of using 300M data for training. Since a huge training dataset is not always available, recent studies attempt to combine a CNN and a Transformer into a hybrid model. Carion *et al.* [2] employed a CNN to extract image features and a Transformer to further process the extracted features. Chen *et al.* [4] designed TransUNet, in which a CNN and a Transformer are combined in a cascade manner to make a strong encoder for 2D medical image segmentation. Although the design of TransUNet is interesting and the



**Fig. 2.** Diagram of CoTr: A CNN-encoder, a DeTrans-encoder, and a decoder. Gray rectangles: CNN blocks. Yellow rectangles: 3D deformable Transformer layers. The CNN-encoder extracts multi-scale feature maps from an input image. The DeTrans-encoder processes the flattened multi-scale feature maps that embedded with the positional encoding in a sequence-to-sequence manner. The features with long-range dependency are generated by the DeTrans-encoder and fed to the decoder for segmentation.

performance is good, it is challenging to optimize this model due to the existence of self-attention [19]. First, it requires extremely long training time to focus the attention, which was initially cast to each pixel uniformly, on salient locations, especially in a 3D scenario. Second, due to its high computational complexity, a vanilla Transformer [19] can hardly process multi-scale and high-resolution feature maps, which play a critical role in image segmentation.

In this paper, we propose a hybrid framework that efficiently bridges **C**onvolutional neural network and **T**ransformer (**CoTr**) for 3D medical image segmentation. CoTr has an encoder-decoder structure. In the encoder, a concise CNN structure is adopted to extract feature maps and a Transformer is used to capture the long-range dependency (see Fig. 1). Inspired by [7, 29], we introduce the deformable self-attention mechanism to the Transformer. This attention mechanism casts attentions only to a small set of key sampling points, and thus dramatically reduces the computational and spatial complexity of Transformer. As a result, it is possible for the Transformer to process the multi-scale feature maps produced by the CNN and keep abundant high resolution information for segmentation. The main contributions of this paper are three-fold: (1) we are the first to explore Transformer for 3D medical image segmentation, particularly in a computationally and spatially efficient way; (2) we introduce the deformable self-attention mechanism to reduce the complexity of vanilla Transformer, and thus enable our CoTr to model the long-range dependency using multi-scale features; (3) our CoTr outperforms the competing CNN-based, Transformer-based, and hybrid methods on the 3D multi-organ segmentation task.

## 2 Materials

The Multi-Atlas Labeling **B**eyond the **C**ranial **V**ault (BCV) dataset<sup>1</sup> was used for this study. It contains 30 labeled CT scans for automated segmentation of 11 abdominal organs, including the spleen (Sp), kidney (Ki), gallbladder (Gb), esophagus (Es), liver (Li), stomach (St), aorta (Ao), inferior vena cava (IVC), portal vein and splenic vein (PSV), pancreas (Pa), and adrenal gland (AG).

## 3 Methods

CoTr aims to learn more effective representations for medical image segmentation via bridging CNN and Transformer. As shown in Fig. 2, it consists of a CNN-encoder for feature extraction, a deformable Transformer-encoder (DeTrans-encoder) for long-range dependency modeling, and a decoder for segmentation. We now delve into the details of each module.

### 3.1 CNN-encoder

The CNN-encoder  $\mathcal{F}^{CNN}(\cdot)$  contains a Conv-IN-ReLU block and three stages of 3D residual blocks. The Conv-IN-ReLU block contains a 3D convolutional layer followed by an instance normalization (IN) [18] and Rectified Linear Unit (ReLU) activation. The numbers of 3D residual blocks in three stages are three, three, and two, respectively.

Given an input image  $\mathbf{x}$  with a height of  $H$ , a width of  $W$ , and a depth (*i.e.*, number of slices) of  $D$ , the feature maps produced by  $\mathcal{F}^{CNN}(\cdot)$  can be formally expressed as

$$\{\mathbf{f}_l\}_{l=1}^L = \mathcal{F}_l^{CNN}(\mathbf{x}; \boldsymbol{\Theta}) \in \mathbb{R}^{C \times \frac{D}{2^l} \times \frac{H}{2^{l+1}} \times \frac{W}{2^{l+1}}}, \quad (1)$$

where  $L$  indicates the number of feature levels,  $\boldsymbol{\Theta}$  denotes the parameters of the CNN-encoder, and  $C$  denotes the number of channels.

### 3.2 DeTrans-encoder

Due to the intrinsic locality of convolution operations, the CNN-encoder cannot capture the long-range dependency of pixels effectively. To this end, we propose the DeTrans-encoder that introduces the multi-scale deformable self-attention (MS-DMSA) mechanism for efficient long-range contextual modeling. The DeTrans-encoder is a composition of an input-to-sequence layer and  $L_D$  stacked deformable Transformer (DeTrans) layers.

**Input-to-sequence Transformation.** Considering that Transformer processes the information in a sequence-to-sequence manner, we first flatten the feature maps produced by the CNN-encoder  $\{\mathbf{f}_l\}_{l=1}^L$  into a 1D sequence. Unfortunately, the operation of flattening the features leads to losing the spatial information that is critical for image segmentation. To address this issue, we supplement the

<sup>1</sup> <https://www.synapse.org/#!Synapse:syn3193805/wiki/217789>

3D positional encoding sequence  $\{\mathbf{p}_l\}_{l=1}^L$  to the flattened  $\{\mathbf{f}_l\}_{l=1}^L$ . For this study, we use sine and cosine functions with different frequencies [19] to compute the positional coordinates of each dimension  $pos$ , shown as follows

$$\begin{cases} PE_{\#}(pos, 2k) = \sin(pos \cdot v) \\ PE_{\#}(pos, 2k + 1) = \cos(pos \cdot v) \end{cases} \quad (2)$$

where  $\# \in \{D, H, W\}$  indicates each of three dimensions,  $v = 1/10000^{2k/\frac{C}{3}}$ . For each feature level  $l$ , we concatenate  $PE_D$ ,  $PE_H$ , and  $PE_W$  as the 3D positional encoding  $\mathbf{p}_l$  and combine it with the flattened  $\mathbf{f}_l$  via element-wise summation to form the input sequence of DeTrans-encoder.

**MS-DMSA Layer.** In the architecture of Transformer, the self-attention layer would look over all possible locations in the feature map. It has the drawback of slow convergence and high computational complexity, and hence can hardly process multi-scale features. To remedy this, we design the MS-DMSA layer that focuses only on a small set of key sampling locations around a reference location, instead of all locations.

Let  $\mathbf{z}_q \in \mathbb{R}^C$  be the feature representation of query  $q$  and  $\hat{\mathbf{p}}_q \in [0, 1]^3$  be the normalized 3D coordinate of the reference point. Given the multi-scale feature maps  $\{\mathbf{f}_l\}_{l=1}^L$  that are extracted in the last  $L$  stages of CNN-encoder, the feature representation of the  $i$ -th attention head can be calculated as

$$\text{head}_i = \sum_l^L \sum_k^K A(\mathbf{z}_q)_{ilqk} \cdot \Psi(\mathbf{f}_l)(\sigma_l(\hat{\mathbf{p}}_q) + \Delta_{\mathbf{p}_{ilqk}}) \quad (3)$$

where  $K$  is the number of sampled key points,  $A(\mathbf{z}_q)_{ilqk} \in [0, 1]$  is the attention weight,  $\Delta_{\mathbf{p}_{ilqk}} \in \mathbb{R}^3$  is the sampling offset of the  $k$ -th sampling point in the  $l$ -th feature level, and  $\sigma_l(\cdot)$  re-scales  $\hat{\mathbf{p}}_q$  to the  $l$ -th level feature. Following [29], both  $A(\mathbf{z}_q)_{ilqk}$  and  $\Delta_{\mathbf{p}_{ilqk}}$  are obtained via linear projection over the query feature  $\mathbf{z}_q$ . Then, the MS-DMSA layer can be formulated as

$$\text{MS-DMSA}(\mathbf{z}_q, \{\mathbf{f}_l\}_{l=1}^L) = \Phi(\text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_H)) \quad (4)$$

where  $H$  is the number of attention heads, and  $\Phi(\cdot)$  is a linear projection layer that weights and aggregates the feature representation of all attention heads.

**DeTrans Layer.** The DeTrans layer is composed of a MS-DMSA layer and a feed forward network, each being followed by the layer normalization [1] (see Fig. 2). The skip connection strategy [11] is employed in each sub-layer to avoid gradient vanishing. The DeTrans-encoder is constructed by repeatedly stacking DeTrans layers.

### 3.3 Decoder

The output sequence of DeTrans-encoder is reshaped into feature maps according to the size at each scale. The decoder, a pure CNN architecture, progressively

upsamples the feature maps to the input resolution (*i.e.*,  $D \times H \times W$ ) using the transpose convolution, and then refines the upsampled feature maps using a 3D residual block. Besides, the skip connections between encoder and decoder are also added to keep more low-level details for better segmentation. We also use the deep supervision strategy by adding auxiliary losses to the decoder outputs with different scales. The loss function of our model is the sum of the Dice loss and cross-entropy loss [12, 24, 28]. More details on the network architecture are in Appendix.

### 3.4 Implementation details

Following [12], we first truncated the HU values of each scan using the range of  $[-958, 327]$  to filter irrelevant regions, and then normalized truncated voxel values by subtracting 82.92 and dividing by 136.97. We randomly split the BCV dataset into two parts: 21 scans for training and 9 scans for test, and randomly selected 6 training scans to form a validation set, which just was used to select the hyper-parameters of CoTr. The final results on the test set are obtained by the model trained on all training scans.

In the training stage, we randomly cropped sub-volumes of size  $48 \times 192 \times 192$  from CT scans as the input. To alleviate the over-fitting of limited training data, we employed the online data augmentation [12], including the random rotation, scaling, flipping, adding white Gaussian noise, Gaussian blurring, adjusting rightness and contrast, simulation of low resolution, and Gamma transformation, to diversify the training set. Due to the benefits of instance normalization [18], we adopted the micro-batch training strategy with a small batch size of 2. To weigh the balance between training time cost and performance reward, CoTr was trained for 1000 epochs and each epoch contains 250 iterations. We adopted the stochastic gradient descent algorithm with a momentum of 0.99 and an initial learning rate of 0.01 as the optimizer. We set the hidden size in MS-DMSA and feed forward network to 384 and 1536, respectively, and empirically set the hyper-parameters  $L_D = 6$ ,  $H = 6$ , and  $K = 4$ . Besides, we formed two variants of CoTr with small CNN-encoders, denoted as CoTr\* and CoTr<sup>†</sup>. In CoTr\*, there is only one 3D residual block in each stage of CNN-encoder. In CoTr<sup>†</sup>, the number of 3D residual blocks in each stage of CNN-encoder is two.

In the test stage, we employed the sliding window strategy, where the window size equals to the training patch size. Besides, Gaussian importance weighting [12] and test time augmentation by flipping along all axes were also utilized to improve the robustness of segmentation. To quantitatively evaluate the segmentation results, we calculated the Dice coefficient scores (Dice) metric that measures the overlapping between a prediction and its ground truth.

## 4 Results

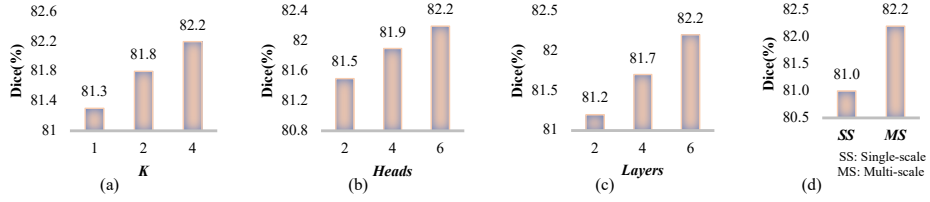
**Comparing to models with only Transformer encoder.** We first evaluated our CoTr against two variants of the state-of-the-art SEGmentation Transformer

**Table 1.** Dice scores of our CoTr and several competing methods on the BCV test set. **CoTr\*** and **CoTr<sup>†</sup>** are two variants of CoTr with small CNN-encoders

Methods	Param (M)	Organs											Ave
		Sp	Ki	Gb	Es	Li	St	Ao	IVC	PSV	Pa	AG	
SETR (ViT-B/16-rand) [27]	100.5	95.2	92.3	55.6	71.3	96.2	80.2	89.7	83.9	68.9	68.7	60.5	78.4
SETR (ViT-B/16-pre) [27]	100.5	94.8	91.7	55.2	70.9	96.2	76.9	89.3	82.4	69.6	70.7	58.7	77.8
CoTr w/o CNN-encoder	21.9	95.2	92.8	59.2	72.2	96.3	81.2	89.9	85.1	71.9	73.3	61.0	79.8
CoTr w/o DeTrans	32.6	96.0	92.6	63.8	77.9	97.0	83.6	90.8	87.8	76.7	81.2	72.6	83.6
APSS [5]	45.5	96.5	93.8	65.6	78.1	97.1	84.0	91.1	87.9	77.0	82.6	73.9	84.3
PP [26]	33.9	96.1	93.1	64.3	77.4	97.0	85.3	90.8	87.4	77.2	81.9	72.8	83.9
Non-local [20]	32.8	96.3	93.7	64.6	77.9	97.1	84.1	90.8	87.7	77.2	82.1	73.3	84.1
TransUnet [4]	43.5	95.9	93.7	63.1	77.8	97.0	86.2	91.0	87.8	77.8	81.6	73.9	84.2
<b>CoTr*</b>	27.9	96.4	94.0	66.2	76.4	97.0	84.2	90.3	87.6	76.3	80.8	72.9	83.8
<b>CoTr<sup>†</sup></b>	36.9	96.2	93.8	66.5	78.6	97.1	86.9	90.8	87.8	77.7	82.8	73.2	84.7
<b>CoTr</b>	41.9	96.3	93.9	66.6	78.0	97.1	88.2	91.2	88.0	78.1	83.1	74.1	<b>85.0</b>

(SETR) [27], which were formed by using randomly initialized and pre-trained ViT-B/16 [9] as the encoder. We also compared to a variant of CoTr that removes the CNN-encoder (CoTr w/o CNN-encoder). To ensure an unprejudiced comparison, all models use the same decoder. The segmentation performance of these models is shown in Table 1, from which three conclusions can be drawn. First, although the Transformer architecture is not limited by the type of input images, the ViT-B/16 pre-trained on 2D natural images does not work well on 3D medical images. The suboptimal performance may be attributed to the domain shift between 2D natural images and 3D medical images. Second, ‘CoTr w/o CNN-encoder’ has about 22M parameters and outperforms the SETR with about 100M parameters. We believe that a lightweight Transformer may be more friendly for medical image segmentation tasks, where there is usually a small training dataset. Third, our CoTr\* with comparable parameters significantly outperforms ‘CoTr w/o CNN-encoder’, improving the average Dice over 11 organs by 4%. It suggests that the hybrid CNN-Transformer encoder has distinct advantages over the pure Transformer encoder in medical image segmentation.

**Comparing to models with only CNN encoder.** Then, we compared CoTr against a variant of CoTr that removes the DeTrans-encoder (CoTr w/o DeTrans) and three CNN-based context modeling methods, *i.e.*, the Atrous Spatial Pyramid Pooling (ASPP) [5] module, pyramid parsing (PP) [26] module, and Non-local [20] module. For a fair comparison, we used the same CNN-encoder and decoder but replaced our DeTrans-encoder with ASPP, PP, and Non-local modules, respectively. The results in Table 1 shows that our CoTr elevates consistently the segmentation performance over ‘CoTr w/o DeTrans’ on all organs and improves the average Dice by 1.4%. It corroborates that our CoTr using a hybrid CNN-Transformer encoder has a stronger ability than using a pure CNN encoder to learn effective representations for medical image segmentation. Moreover, comparing to these context modeling methods, our Transformer architecture contributes to more accurate segmentation.



**Fig. 3.** Average Dice over all organs obtained on the validation set versus (a) the number of sampled key points  $K$ , (b) number of heads  $H$ , and (c) number of DeTrans layers  $L_D$ , and (d) Average Dice obtained by our CoTr using, respectively, single-scale and multi-scale feature maps on the validation set.

**Comparing to models with hybrid CNN-Transformer encoder.** We also compared CoTr to other hybrid CNN-Transformer architectures like TransUNet [4]. To process 3D images directly, we extended the original 2D TransUNet to a 3D version by using 3D CNN-encoder and decoder as done in CoTr. We also set the number of heads and layers of Transformer in 3D TransUNet to be the same as our CoTr. It shows in Table 1 that CoTr steadily beats TransUNet in the segmentation of all organs, particularly for the gallbladder and pancreas segmentation. Even with a smaller CNN-encoder, CoTr<sup>†</sup> still achieves better performance than TransUNet in the segmentation of seven organs. The superior performance owes to the deformable mechanism in CoTr that makes it possible to process high-resolution and multi-scale feature maps due to the reduced computational and spatial complexities.

**Computational Complexity.** The proposed CoTr was trained using a workstation with a NVIDIA GTX 2080Ti GPU and the Pytorch software packages. It took about 2 days for training, and less than 30ms to segment a volume of size  $48 \times 192 \times 192$ .

## 5 Discussion on Hyper-parameter Settings

In the DeTrans-encoder, there are three hyper-parameters, *i.e.*,  $K$ ,  $H$ , and  $L_D$ , which represent the number of sampled key points, heads, and stacked DeTrans layers, respectively. To investigate the impact of their settings on the segmentation, we set  $K$  to 1, 2, and 4, set  $H$  to 2, 4, and 6, and set  $L_D$  to 2, 4, and 6. In Fig. 3 (a-c), we plotted the average Dice over all organs obtained on the validation set versus the values of  $K$ ,  $H$ , and  $L_D$ . It shows that increasing the number of  $K$ ,  $H$ , or  $L_D$  can improve the segmentation performance. To demonstrate the performance gain resulted from the multi-scale strategy, we also attempted to train CoTr with single-scale feature maps from the last stage. The results in Fig. 3 (d) show that using multi-scale feature maps instead of single-scale feature maps can effectively improve the average Dice by 1.2%.



## 6 Conclusion

In this paper, we propose a hybrid model of CNN Transformer, namely CoTr, for 3D medical image segmentation. In this model, we design the deformable Transformer (DeTrans) that employs the deformable self-attention mechanism to reduce the computational and spatial complexities of modelling the long-range dependency on multi-scale and high-resolution feature maps. Comparative experiments were conducted on the BCV dataset. The superior performance of our CoTr over both CNN-based and vanilla Transformer-based models suggests that, via combining the advantages of CNN and Transformer, the proposed CoTr achieves the balance in keeping the details of low-level features and modeling the long-range dependency. As a stronger baseline, our CoTr can be extended to deal with other structures (*e.g.*, brain structure or tumor segmentation) in the future.

## References

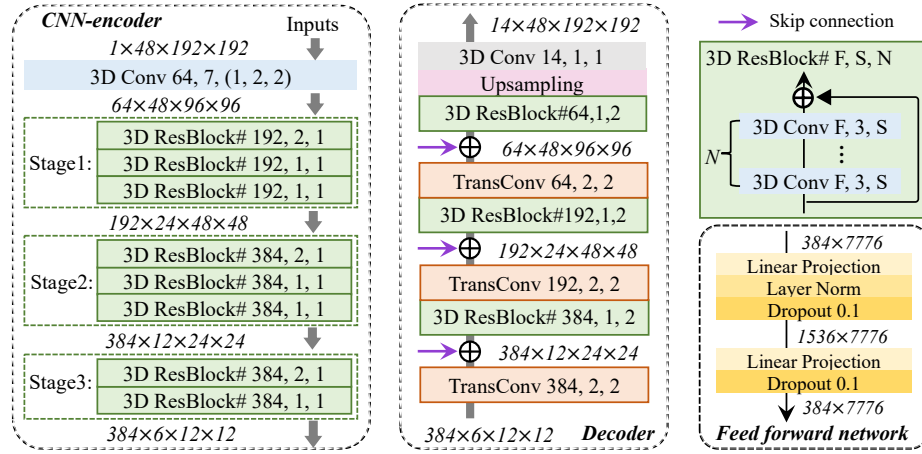
1. Ba, J.L., Kiros, J.R., Hinton, G.E.: Layer normalization. arXiv preprint arXiv:1607.06450 (2016)
2. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: European Conference on Computer Vision. pp. 213–229. Springer (2020)
3. Chen, H., Wang, Y., Guo, T., Xu, C., Deng, Y., Liu, Z., Ma, S., Xu, C., Xu, C., Gao, W.: Pre-trained image processing transformer. arXiv preprint arXiv:2012.00364 (2020)
4. 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 preprint arXiv:2102.04306 (2021)
5. Chen, L.C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-decoder with atrous separable convolution for semantic image segmentation. In: Proceedings of the European conference on computer vision (ECCV). pp. 801–818 (2018)
6. Cohen, N., Shashua, A.: Inductive bias of deep convolutional networks through pooling geometry. arXiv preprint arXiv:1605.06743 (2016)
7. Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., Wei, Y.: Deformable convolutional networks. In: Proceedings of the IEEE international conference on computer vision. pp. 764–773 (2017)
8. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
9. 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. arXiv preprint arXiv:2010.11929 (2020)
10. Fang, X., Yan, P.: Multi-organ segmentation over partially labeled datasets with multi-scale feature abstraction. IEEE Transactions on Medical Imaging **39**(11), 3619–3629 (2020)
11. 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. pp. 770–778 (2016)

12. Isensee, F., Jäger, P.F., Kohl, S.A., Petersen, J., Maier-Hein, K.H.: Automated design of deep learning methods for biomedical image segmentation. arXiv preprint arXiv:1904.08128 (2019)
13. Jiang, Y., Chang, S., Wang, Z.: Transgan: Two transformers can make one strong gan. arXiv preprint arXiv:2102.07074 (2021)
14. Parmar, N., Vaswani, A., Uszkoreit, J., Kaiser, L., Shazeer, N., Ku, A., Tran, D.: Image transformer. In: International Conference on Machine Learning. pp. 4055–4064. PMLR (2018)
15. Peng, C., Zhang, X., Yu, G., Luo, G., Sun, J.: Large kernel matters—improve semantic segmentation by global convolutional network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 4353–4361 (2017)
16. 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. pp. 234–241. Springer (2015)
17. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-efficient image transformers & distillation through attention. arXiv preprint arXiv:2012.12877 (2020)
18. Ulyanov, D., Vedaldi, A., Lempitsky, V.: Instance normalization: The missing ingredient for fast stylization. arXiv preprint arXiv:1607.08022 (2016)
19. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Proceedings of the 31st International Conference on Neural Information Processing Systems. pp. 6000–6010 (2017)
20. Wang, X., Girshick, R., Gupta, A., He, K.: Non-local neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7794–7803 (2018)
21. Wang, Y., Xu, Z., Wang, X., Shen, C., Cheng, B., Shen, H., Xia, H.: End-to-end video instance segmentation with transformers. arXiv preprint arXiv:2011.14503 (2020)
22. Yu, F., Koltun, V.: Multi-scale context aggregation by dilated convolutions. In: International Conference on Learning Representations (ICLR) (2016)
23. Zhang, J., Xie, Y., Wang, Y., Xia, Y.: Inter-slice context residual learning for 3d medical image segmentation. IEEE Transactions on Medical Imaging (2020)
24. Zhang, J., Xie, Y., Zhang, P., Chen, H., Xia, Y., Shen, C.: Light-weight hybrid convolutional network for liver tumor segmentation. In: IJCAI. pp. 4271–4277 (2019)
25. Zhang, L., Zhang, J., Shen, P., Zhu, G., Li, P., Lu, X., Zhang, H., Shah, S.A., Bennamoun, M.: Block level skip connections across cascaded v-net for multi-organ segmentation. IEEE Transactions on Medical Imaging **39**(9), 2782–2793 (2020)
26. Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid scene parsing network. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2881–2890 (2017)
27. Zheng, S., Lu, J., Zhao, H., Zhu, X., Luo, Z., Wang, Y., Fu, Y., Feng, J., Xiang, T., Torr, P.H., et al.: Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. arXiv preprint arXiv:2012.15840 (2020)
28. Zhou, Z., Siddiquee, M.M.R., Tajbakhsh, N., Liang, J.: Unet++: Redesigning skip connections to exploit multiscale features in image segmentation. IEEE Transactions on Medical Imaging **39**(6), 1856–1867 (2019)
29. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable detr: Deformable transformers for end-to-end object detection. arXiv preprint arXiv:2010.04159 (2020)

## 7 Appendix

### 7.1 Detailed network architecture

Fig. 4 shows the architecture of CNN-encoder, decoder and feed forward network in Detrans-encoder. It consists of a Conv-In-Relu and three stages of 3D residual blocks. The numbers of 3D residual blocks are three, three, and two in three stages, respectively. The decoder contains four upsampling modules. Each of first three modules has a TransConv layer followed by a residual block, and a pixel-wise summation with the corresponding feature maps from the encoder and the TransConv layer. The last module comprises of an Upsampling layer followed by a  $1 \times 1$  Conv layer that maps the 64-channel feature maps to the desired number of classes. The feed forward network in Detrans-encoder has two linear projection layers. The first layer is followed by a layer normalization layer and a Dropout layer. The second layer is followed by a Dropout layer.



**Fig. 4.** Detailed architecture of CNN-encoder, decoder and feed forward network. Blue ‘Conv’: Conv-In-Relu block that contains a 3D convolutional layer followed by an instance normalization (IN) layer [18] and ReLU activation; Gray ‘Conv’: a 3D convolutional layer; Orange ‘TransConv’: 3D transposed convolutional layer. Note that the numbers in each Conv block / layer indicate the number of filters, kernel size, and stride, respectively. The numbers in each residual block indicate the number of filters, stride, and Conv blocks, respectively.

### 7.2 Loss function

We jointly use the Dice loss and cross-entropy loss for optimization, which is popular in many medical image segmentation applications and has achieved

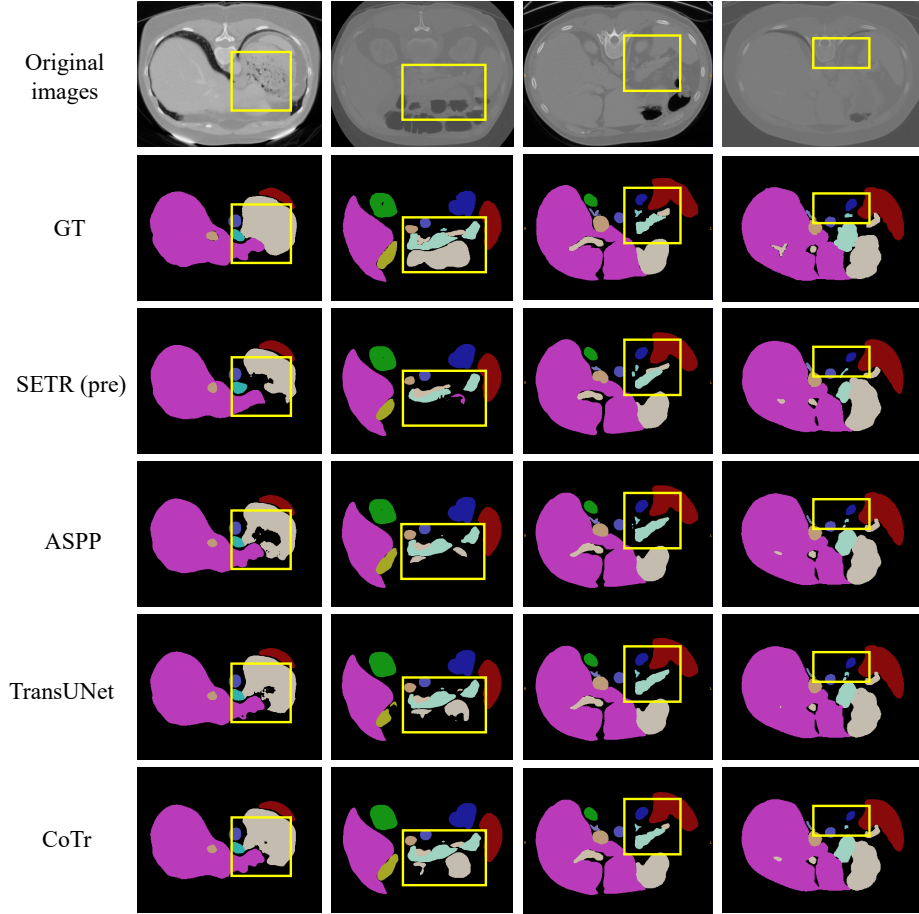
prominent success [16, 24, 28]. The loss function is formulated as

$$\mathcal{L} = \frac{1}{c} \sum_{n=1}^c \left\{ -\frac{2 \sum \tilde{\mathbf{y}}^n \mathbf{y}^n + \varepsilon}{\sum (\tilde{\mathbf{y}}^n + \mathbf{y}^n) + \varepsilon} - \mathbb{E}[\mathbf{y}^n \log \tilde{\mathbf{y}}^n] \right\} \quad (5)$$

where the first item is the soft Dice loss, the second item is the cross-entropy loss, the prediction and ground truth are denoted by  $\tilde{\mathbf{y}}$  and  $\mathbf{y}$ , respectively,  $\mathbb{E}$  is the expectation operation,  $\varepsilon$  is a smoothing factor, and  $c$  is the number of categories. To speed up convergence and alleviate the vanishing gradient problem, we also use the deep supervision strategy that adds auxiliary losses to the decoder outputs with different resolutions. The total loss function is the sum of the losses at all resolutions.

### 7.3 Visualization

The segmentation results produced by (1) SETR with pre-trained ViT-B/16, (2) replacing DeTrans-encoder with ASPP module, (3) 3D TransUNet, and (4) our CoTr, were visually compared in Fig. 5. We can see that: 1) comparing to the pure Transformer encoder method (SETR) and pure CNN encoder method (ASPP), our CoTr with the hybrid CNN-Transformer encoder is able to produce the segmentation results that are more similar to the ground truth, and 2) our CoTr are more likely to produce less false positives compared to TransUNet, which confirms the superiority of our 3D deformable Transformer over vanilla Transformer.



**Fig. 5.** Visualization of segmentation results of four cases. The regions in yellow rectangles indicate our superiority. Each type of organs are denoted by a unique color.