

Washington University School of Medicine

Digital Commons@Becker

2020-Current year OA Pubs

Open Access Publications

1-1-2022

3D bi-directional transformer U-Net for medical image segmentation

Xiyao Fu

Zhexian Sun

Haoteng Tang

Eric M Zou

Heng Huang

See next page for additional authors

Follow this and additional works at: https://digitalcommons.wustl.edu/oa_4



Part of the [Medicine and Health Sciences Commons](#)

Please let us know how this document benefits you.

Authors

Xiyao Fu, Zhexian Sun, Haoteng Tang, Eric M Zou, Heng Huang, Yong Wang, and Liang Zhan



OPEN ACCESS

EDITED BY

Xiaoqian Wang,
Purdue University, United States

REVIEWED BY

Qi Huang,
The University of Utah, United States
Mengting Liu,
Sun Yat-sen University, China

*CORRESPONDENCE

Haoteng Tang
✉ haoteng.tang@pitt.edu
Liang Zhan
✉ liang.zhan@pitt.edu

SPECIALTY SECTION

This article was submitted to
Medicine and Public Health,
a section of the journal
Frontiers in Big Data

RECEIVED 26 October 2022

ACCEPTED 06 December 2022

PUBLISHED 06 January 2023

CITATION

Fu X, Sun Z, Tang H, Zou EM, Huang H,
Wang Y and Zhan L (2023) 3D
bi-directional transformer U-Net for
medical image segmentation.
Front. Big Data 5:1080715.
doi: 10.3389/fdata.2022.1080715

COPYRIGHT

© 2023 Fu, Sun, Tang, Zou, Huang,
Wang and Zhan. This is an
open-access article distributed under
the terms of the [Creative Commons
Attribution License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use,
distribution or reproduction in other
forums is permitted, provided the
original author(s) and the copyright
owner(s) are credited and that the
original publication in this journal is
cited, in accordance with accepted
academic practice. No use, distribution
or reproduction is permitted which
does not comply with these terms.

3D bi-directional transformer U-Net for medical image segmentation

Xiyao Fu¹, Zhexion Sun², Haoteng Tang^{1*}, Eric M. Zou³,
Heng Huang¹, Yong Wang^{2,4,5,6} and Liang Zhan^{1*}

¹Department of Electrical and Computer Engineering, University of Pittsburgh, Pittsburgh, PA, United States, ²Department of Biomedical Engineering, Washington University in St. Louis, St. Louis, MO, United States, ³Montgomery Blair High School Maryland, 51 University Blvd E, Silver Spring, MD, United States, ⁴Department of Electrical and Systems Engineering, Washington University in St. Louis, St. Louis, MO, United States, ⁵Department of Obstetrics and Gynecology, Washington University in St. Louis, St. Louis, MO, United States, ⁶Department of Radiology, Washington University in St. Louis, St. Louis, MO, United States

As one of the popular deep learning methods, deep convolutional neural networks (DCNNs) have been widely adopted in segmentation tasks and have received positive feedback. However, in segmentation tasks, DCNN-based frameworks are known for their incompetence in dealing with global relations within imaging features. Although several techniques have been proposed to enhance the global reasoning of DCNN, these models are either not able to gain satisfying performances compared with traditional fully-convolutional structures or not capable of utilizing the basic advantages of CNN-based networks (namely the ability of local reasoning). In this study, compared with current attempts to combine FCNs and global reasoning methods, we fully extracted the ability of self-attention by designing a novel attention mechanism for 3D computation and proposed a new segmentation framework (named 3DTU) for three-dimensional medical image segmentation tasks. This new framework processes images in an end-to-end manner and executes 3D computation on both the encoder side (which contains a 3D transformer) and the decoder side (which is based on a 3D DCNN). We tested our framework on two independent datasets that consist of 3D MRI and CT images. Experimental results clearly demonstrate that our method outperforms several state-of-the-art segmentation methods in various metrics.

KEYWORDS

semantic segmentation, COVID, lung, placenta, transformer, 3D UNet, CT, MRI

1. Introduction

In the recent few years, deep convolutional neural networks (DCNNs) (Krizhevsky et al., 2012; Simonyan and Zisserman, 2014; He et al., 2016; Badrinarayanan et al., 2017; Huang et al., 2020; Pan et al., 2020) have achieved considerable progress in medical image segmentation (Long et al., 2015; Noh et al., 2015; Chen L.-C. et al., 2018; Tokunaga et al., 2019; Liu et al., 2022; Zhang et al., 2022). However, limited to the local receptive field of the convolutional filter, DCNN-based frameworks are incapable of capturing long-range dependencies from global features for semantic segmentation. To tackle this, several

strategies can be considered. The first is to use the dilated convolution operation to enlarge the size of the receptive field of the convolutional filter (Yu and Koltun, 2015; Yang et al., 2017; Zhang et al., 2017; Liu et al., 2021). However, this enlarged local receptive field is still limited by the size of dilation. Another solution is to model the feature map as graph structures and investigate the long-range dependencies through the message passing mechanism of different graph learning models (e.g., graph convolution networks) (Li and Gupta, 2018; Chen et al., 2019; Li et al., 2020; Jia et al., 2021). Although these graph learning models have shown great potential in enhancing the global reasoning ability of DCNNs, they have very high requirements for computation and memory due to the constructed large-size graphs.

The attention mechanism (Hochreiter and Schmidhuber, 1997; Vaswani et al., 2017) is a computation scheme that tries to generate representations *via* different types of global features at each step. Since attention can be regarded as the conversion and transformation among the query (q), key (k), and value (v) triplet, attention computation is to generate the q based on the combination of the k–v pair. As it is natural to integrate a cycling computation in recurrent cells, traditional attention mechanisms are integrated within recurrent neural networks (e.g., Hochreiter and Schmidhuber, 1997; Cho et al., 2014), which inevitably impairs the efficiency of recurrent networks compared with linear/residual networks (Vaswani et al., 2017). To cope with this, Vaswani et al. (2017) proposed a transformer, a structure consisting of a series of identical encoder blocks connected with a series of identical decoder blocks, which all have no convolutional layers and are connected in a residual way. The original transformer supported by self-attention works exceptionally well in some tasks like machine translation but not in visual tasks (Chen et al., 2021). This is mainly due to the lack of convolution layers that makes the model struggle to detect local features.

For the aforementioned reasons, convolutional-based frameworks are still preferred for segmentation tasks. Although several other models (Goodfellow et al., 2014; Chen Y. et al., 2018) have been proven feasible, DCNNs remain to be one of the most effective methods. Multiple variants of DCNNs have been proposed to make the segmentation process more effective, one of the most crucial ones is the UNet (Ronneberger et al., 2015), which is a symmetric structure consisting of convolutional blocks with skip connections. These convolutional blocks have descending dimensions on the encoder side and ascending dimensions on the decoder side. However, due to the intrinsic fully convolution structure, UNet is suboptimal to relate local features to global representations with more variant distribution (Chen et al., 2021). To cope with the drawbacks of UNet, many methods have been proposed (Liu et al., 2018; Zhou et al., 2019; Diakogiannis et al., 2020; Huang et al., 2020). However, these methods are either very time-consuming or require

heavy computations, which make it impossible to be applied to 3D objects.

Under such circumstances, the self-attention mechanism seems to be a nearly optimal solution. It is highly modularized and can stretch the number of self-attention cells according to the training environment. It can also train on vast datasets due to the training nature of attention. Therefore, researchers combined the transformer with convolutional layers for medical image segmentation (Li et al., 2022). On the one hand, the transformer encodes tokenized image patches from a CNN feature map as the input sequence for extracting global contexts. On the other hand, the decoder upsamples the encoded features, which are then combined with high-resolution CNN feature maps to enable precise localization.

However, this approach still has some obstacles, especially in the segmentation of 3D objects. This is partially due to transformers (Vaswani et al., 2017) requiring the input features to have temporal information. Since the self-attention does not compute with a clear direction, features have to be preprocessed with temporal info (e.g., cosine function) as input embeddings before training. Although this learning process can be seen as natural (scanning the features linearly and with order), it will restrict the performance of high-dimensional data. For example, many existing transformer approaches (Parmar et al., 2018; Huang et al., 2020; Chen et al., 2021) will cut the 3D object into 2D slice sequences to meet the temporal encoding requirement; however, the segmentation performance is actually worse because the 2D slice cutting will destroy the smoothness of the object in 3D space. Bi-directional transformer (Devlin et al., 2018) is a powerful upgrade version of transformer. It is a structure with no decoder and processes the inputs all at once with masks to create temporal/spatial continuity. However, we will show in the experiment section that bi-directional transformers can serve as a strong encoder but still struggles to get better results on 3D segmentation. To compensate for the loss of feature resolution brought by transformers, we propose 3D transformer UNet (3DTU), which employs a hybrid CNN–transformer architecture to leverage both detailed high-resolution spatial information from CNN features and the global context encoded by our new 3D bi-directional transformer module. We show that such a design allows our framework to preserve the advantages of self-attention mechanisms and also get considerably improved results on 3D image segmentation compared with previous U-Net-based or transformer-based methods. To sum up, our contributions to this article can be summarized as follows:

- We proposed a new 3D bi-directional framework to learn deep 3D features for medical image semantic segmentation.
- We designed a novel attention mechanism specifically suitable for network training and self-attention computation for 3D objects.

- We verified our new framework on multiple datasets, consisting of different imaging modalities (MRI and CT images) and different organs (placenta and lungs infected with COVID) and obtained state-of-the-art (SOTA) results. Our method beat baselines in performances on multiple metrics.

2. Related work

2.1. Fully convolutional network in medical image segmentation

Many studies have attempted to adopt convolutional networks to medical image segmentation. For example, [Liu et al. \(2018\)](#) presented a hybrid network consisting of both 3D CNN and 2D CNN in the brain image segmentation for Alzheimer's disease (AD) studies. [Ronneberger et al. \(2015\)](#) presented UNet, one of the most iconic encoder-decoder-based methods for medical image segmentation. Their method consists of convolutional blocks that have a U-shaped dimension variation. Specifically, from the input layer of the encoder to the input layer of the decoder, each block's dimension is descending. And the decoder has an ascending dimension that is matched to the encoder blocks. Such a design makes sure that the learning ability of the framework is powerful enough to find the abstract of the locality and output a global representation map. Several adjustments (e.g., [Zhou et al., 2019](#); [Huang et al., 2020](#)) have been made to the original UNet model. For example, U-Net3+ ([Huang et al., 2020](#)) and its variations, although proved effective, still suffer from the locality-heavy learning scheme. Some researchers tried to boost the local reasoning of convolutional layers through the residual structure. For example, ResUNet ([Diakogiannis et al., 2020](#)) proposed a residual block between every two convolutional blocks on both the encoder side and decoder side as well as skip-connection between residual blocks with the same dimension between the encoder and decoder. [Isensee et al. \(2021\)](#) argued that the understanding of the datasets needed for training is more important than the network itself since most UNet-based moderations have achieved little progress. The authors proposed nnUNet, a robust network, that is designed based on the combination of 2D and 3D UNet. The authors also made different training configurations (normalization tricks, cropping, activation functions, etc.) based on the datasets.

2.2. Transformers

Transformers ([Vaswani et al., 2017](#)) were initially proposed for general NLP tasks and quickly gain widespread attention by beating previous most state-of-the-art results by a large margin. [Devlin et al. \(2018\)](#) converted the original transformer

model into BERT, and introduced the called bi-directional transformers, which are proven effective again. Naturally, multiple efforts have been made to adjust the learning ability of transformers in the computer vision domain. Several variants of transformers have emerged recently. [Parmar et al. \(2018\)](#) presented one of the early works to adjust vanilla transformers by incorporating visual information. This model pre-processes each pixel of one image through a 1×1 convolution layer. Then, the embeddings are computed with positional embeddings before feeding into transformers for super-resolution tasks. In another attempt at visual tasks, [Dosovitskiy et al. \(2020\)](#) proposed Vision transformer (ViT), which presented a novel way of input embedding on visual information. It achieved state-of-the-art on ImageNet classification by directly applying transformers with global self-attention to full-sized images. Specifically, ViT flattens an image to fixed-sized pixels, which then be linearly added to positional embeddings before feeding to transformer encoders. [Valanarasu et al. \(2021\)](#) presented gated axial attention that creates a gated scheme to improve learning ability on the local scale.

2.3. Combination of UNet and transformer in medical image segmentation

Multiple attempts have been made to combine the UNet with transformer in both framework structure and inner encoder/decoder computation. TransUNet ([Chen et al., 2021](#)) consists of a series of transformer units as the encoder and the right half of the UNet as the decoder to generate predictions in medical image segmentation. Both the encoder and the decoder ([Chen et al., 2021](#)) are computed in a 2D scenario. [Yun et al. \(2021\)](#) introduced SpecTr, a framework that takes spectral normalization into the computation between convolution and attention blocks. Their methods achieved better results than the baseline when training on hyperspectral medical images. [Wang et al. \(2021\)](#) presented TransBTS that utilizes 3D CNN to extract input representations. UNet transformer, presented by [Petit et al. \(2021\)](#), replaces self-attention modules in transformer encoder/decoder cells by convolutional blocks and batch normalization computations. Another attempt is Swin-UNet ([Cao et al., 2021](#)), which instead replaces convolution blocks in the UNet-Structure network with self-attention modules. Several works follow similar methods including UNETR ([Hatamizadeh et al., 2022b](#)), SWIN UNETR ([Hatamizadeh et al., 2022a](#)), CoTr ([Xie et al., 2021](#)), nnFormer ([Zhou et al., 2021](#)), DS-TransUNet ([Lin et al., 2022](#)), UTNet ([Gao et al., 2021](#)), and PNS-Net ([Ji et al., 2021](#)). In UNETR, the authors presented a novel 3D transformer encoder and a voxel-wise loss for model training. For the positional embedding, they adopted a strategy from the Visual

transformer, which divides the 3D images into 3D patches. The decoder in their work consists of several convolutional blocks in different dimensions and skip connections to the encoder. The SWIN UNETR is proposed for 3D multi-modal MRI brain image studies, which is different from the SWIN UNET that is proposed for 2D images. The CoTr utilized a DeTrans-encoder with a novel attention mechanism and a CNN-based decoder. The nnFormer utilizes CNN as part of an encoder, which leverages the ability of local feature extraction of CNN structures. Moreover, it utilizes transformer structures as its decoder and the second part of its encoder. There are two differences between our 3DTU and the nnFormer. First, we utilize a CNN-based structure (i.e., the right part of 3DUNet) as our decoder. Then, we design an attention mechanism that computes the attention scores from different directions.

The aforementioned methods adjust the transformers in visual tasks by introducing their own positional embedding rules. Although these rules are to an extent useful, their performance all suffers from the slicing of 3D data to adjust the positional embeddings. In this study, positional embeddings are not needed technically, even for 3D data. We modify the multi-head attention from its original form to a refined computation scheme that fully utilizes the potentials of transformer and UNet. More importantly, our encoder is a refined bi-directional transformer, which learns the feature from three (i.e., along x, y, and z) directions simultaneously.¹

3. Methods

We propose a 3D UNet-based framework with bi-directional transformers (named 3DTU) in this work. The self-attention mechanism in the proposed bi-directional transformers can improve the ability of generalization of the framework encoder. We will delve into the technical details in this section.

As shown in Figure 1, our proposed 3DTU is an encoder-decoder framework, where the encoder consists of two modules including a feature extraction module (see Part I in Figure 1) and a bi-directional transformer module (see Part II in Figure 1). Given a 3D image $I \in \mathcal{R}^{h \times w \times d \times c}$, where h , w , and d are the shapes of the image and c is the image channel number, the feature extraction module projects the 3D image I as a latent representation X via basic convolutional neural networks (CNNs). Then, the 3D bi-directional transformer cells take the latent representation X as input and yield the masked latent representation X_M by using Masked-LM (MLM) (Devlin et al., 2018) step by step. Finally, the decoder part utilizes the

masked latent representations to reconstruct the segmentation predictions for loss computation.

3.1. Encoder with 3D bi-directional transformer

As aforementioned, the encoder of the 3DTU consists of two parts. The first part of the encoder is a CNN-based feature extraction module. We aim to convert the original 3D image (I) into an iso-dimensional latent cube representation ($X \in \mathcal{R}^{1 \times p \times p \times p}$) via this module as assistance to capture the image locality for transformer modules, since the transformer module may not have enough ability to capture the image local features. We will show this point in the ablation studies. Particularly, the feature extraction module includes two convolutional layers followed by a fully-connected (FC) layer and a max-pooling layer in between the two convolutional layers. The FC layer is used to adapt the feature dimension.

The bi-directional transformer module takes the latent cube representation X as input and computes multi-head attentions with the MLM strategy (Devlin et al., 2018). Details of the bi-directional transformer module are shown in Figure 2. In general, each cell in the bi-directional transformer module generates the latent feature map X_1 by the following steps:

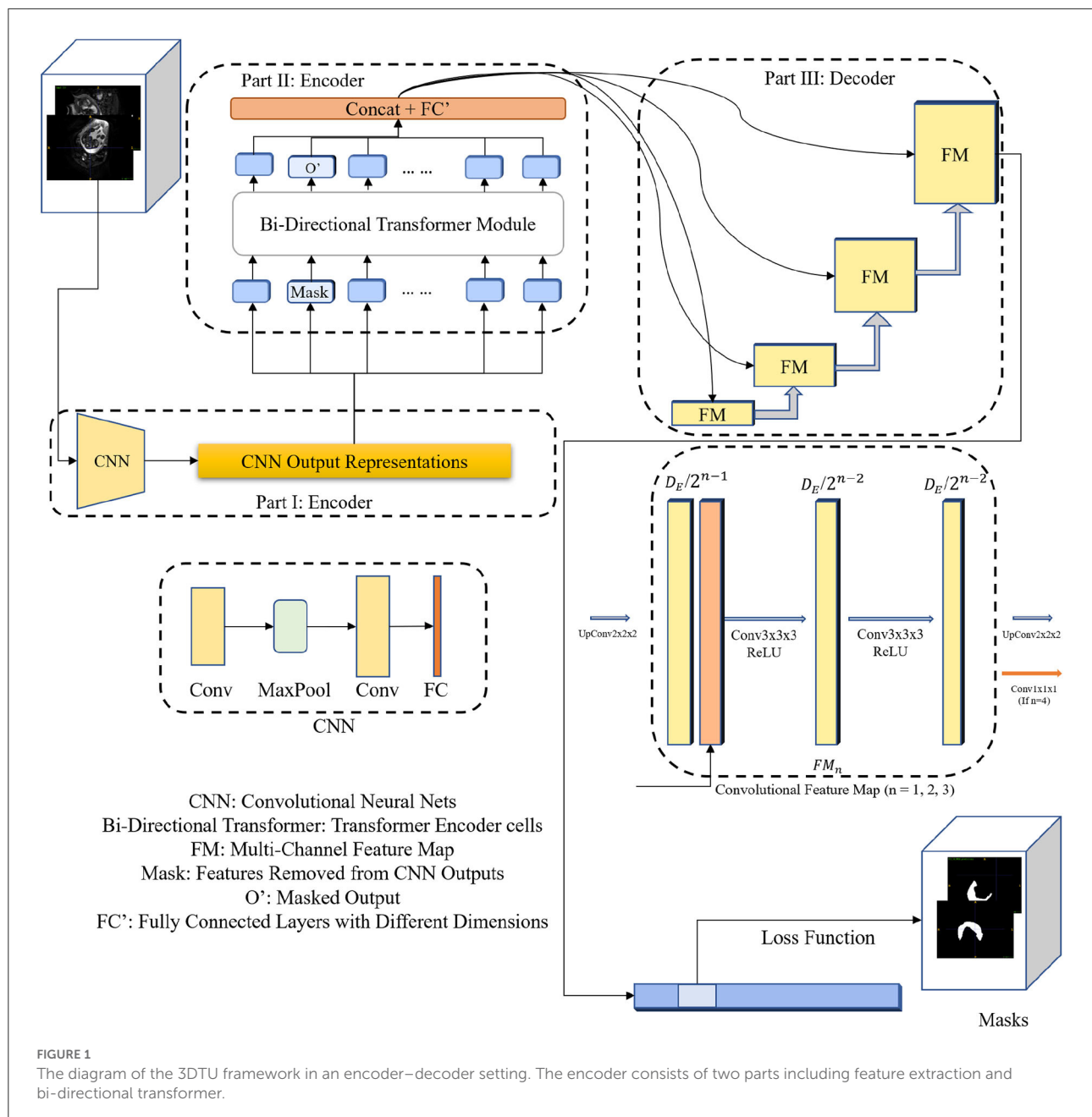
$$\begin{aligned} X' &= \text{Att}(\text{Norm}(X)) + X, \\ X'' &= \text{FF}(\text{Norm}(X')), \\ X_1 &= X' + X'', \end{aligned} \quad (1)$$

where $\text{Att}(\cdot)$ is the multi-head self-attention operation, $\text{Norm}(\cdot)$ is a 3D normalization operation, and $\text{FF}(\cdot)$ is the feed forward layer (i.e., FC layer). $+$ denotes a pixel-wise add operation. Particularly, the multi-head attention is computed as follows:

$$\begin{aligned} \text{Att_head}_i^{x,y,z} &= \text{SDP}(Q, K, V) \times W, \\ \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_i^x, \text{head}_i^y, \text{head}_i^z), \end{aligned} \quad (2)$$

where $\text{SDP}(\cdot)$ is the Scaled Dot-Product Attention, W is the trainable parameters for linear projections (i.e., L_q, L_k , and L_v in Figure 2) and $\text{Concat}(\cdot)$ denotes a concatenation operation. Q, K , and V are the query-key-value triplets defined by the transformer cell. Note that our proposed attention mechanism can yield the attention score by scanning the query-key-value triplets in three different directions (i.e., along x, y, and z axes, respectively), which gain plentiful discriminative and anisotropic semantic information for the 3D image segmentation.

¹ We use the term “bi-directional” by following previous studies. However, our 3DTU learns the features from three directions instead.



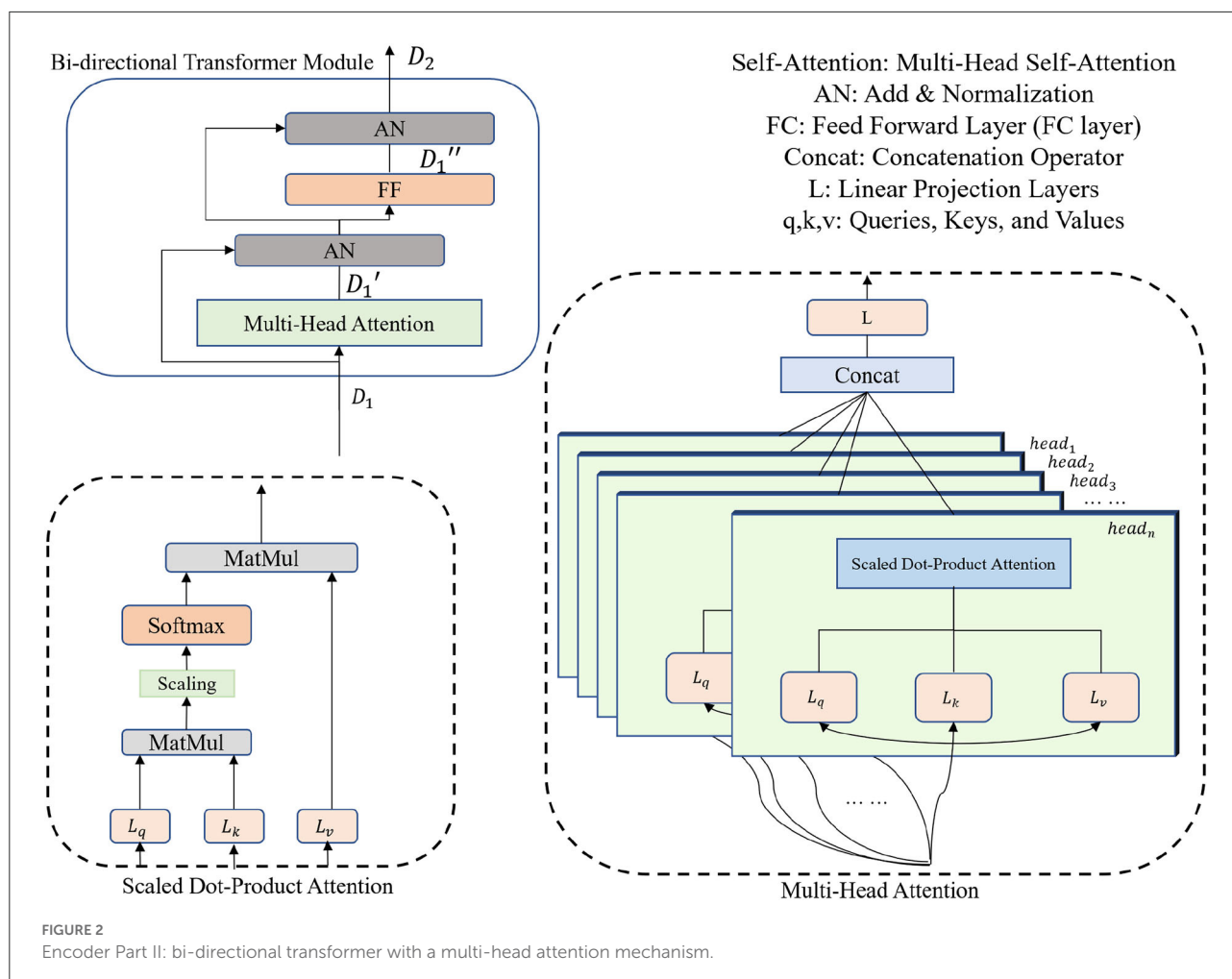
3.2. UNet-based decoder

As shown in Figure 1, we utilize convolutional blocks with ascensional dimensions in the decoder part. A residual connection is adopted between the encoder side and the decoder side. Particularly, a cascaded of multi-channel feature map (FM) blocks are integrated into the decoder part, each of which contains two $3 \times 3 \times 3$ convolutional layers and an upsampling layer. The channel number of feature maps reduces by half after each FM block. In the last FM block, instead of upsampling

layer, a $1 \times 1 \times 1$ convolutional layer is used to generate final segmentation predictions.

3.3. Loss function and supervision manner

Since the MLM strategy is used in the encoder part, where a portion of image features are masked (i.e., set to 0 values)



and the other portions remain the same. Hence, our goal is to use the uncovered portions to predict the masked portions (Devlin et al., 2018), in which the loss is only estimated based on the masked regions. Particularly, the loss function can be formulated as:

$$\mathcal{L} = \alpha \times \ell_{\text{dice}}(\hat{y}_{\text{mask}}, y_{\text{mask}}) + (1 - \alpha) \times \ell_{\text{BCE}}(\hat{y}_{\text{mask}}, y_{\text{mask}}), \quad (3)$$

where \hat{y}_{mask} and y_{mask} are the masked regions of segmentation prediction and ground truth, respectively. $\alpha \in [0, 1]$ is the loss weight.

4. Experiments

4.1. Datasets

We used three datasets obtained from different modalities for this study, including Placenta MRI (Placenta) dataset, COVID-19 CT lung and infection segmentation (Covid20) dataset, and Multi-Atlas

Labeling Beyond the Cranial Vault (Synapse) dataset. Details of data description and preprocessing are shown below.

Placenta MRI dataset was collected from the Washington University in Saint Louis (WUSTL) (Sun et al., 2022), where all data were de-identified before processing. The data collection and related studies were approved by the Institutional Review Board at the WUSTL. A total of 81 MRI scans were collected from 46 pregnant patients (mean age = 23.91 ± 3.02 yo, mean BMI = 25 ± 3.66 at recruitment) with normal singleton pregnancy who underwent MRI during the third trimester, by a Siemens 3T VIDA scanner. Of the 46 patients, 21 patients had the single scan and 25 patients had multiple longitudinal scans. The average gestational ages (GA) during MRI scans were 34.12 ± 1.07 weeks (Min GA 28 weeks 3 days, max GA 38 weeks 6 days). T2-weighted MRI of the entire uterus was acquired with a 2D EPI sequence in the left lateral position. The MRI data has a fixed acquisition matrix of $128 \times 128 \times 115$, and variable voxel sizes from $3 \times 3 \times 3$ mm to $3.5 \times 3.5 \times 3.5$ mm, up to the patient's size. Manual segmentation of the placenta

TABLE 1 Quantitative segmentation results of different methods on two datasets, where mIOU and DICE are in %.

	Placenta dataset			Covid20 dataset			Synapse dataset		
	mIOU	DICE	HD95	mIOU	DICE	HD95	mIOU	DICE	HD95
2D UNet	67.6	72.3	12.0	73.6	78.3	112.5	56.3	60.6	45.7
3D UNet	72.5	78.6	10.7	78.1	84.0	97.6	59.4	62.2	42.2
UNet++	74.5	77.1	8.2	80.3	84.6	63.0	67.1	73.7	34.0
TransUNet	73.6	80.0	7.4	83.1	89.2	45.8	70.2	77.5	31.7
ViT	72.9	79.7	8.5	84.2	89.0	70.3	65.3	67.9	36.1
nnFormer	78.3	82.1	10.2	81.0	89.9	66.2	81.8	86.6	10.6
nnUNet	78.9	83.6	8.7	90.3	91.6	59.9	84.2	89.8	16.6
3DTU (Ours)	79.8	84.0	7.2	90.5	92.0	59.4	85.0	87.3	18.4

The best results are shown in red and the second-best results are shown in blue.

regions was conducted by experienced radiologists for all MRI images.

COVID19-CT-Seg20 dataset (Covid20) contains 20 COVID-19 3D CT images, where lungs and infections were annotated by two radiologists and verified by an experienced radiologist² (Jun et al., 2021). We only focused on the segmentation of the COVID-19 infections in this study, since it is more challenging and important.

Multi-atlas labeling beyond the cranial vault (Synapse) dataset.³ We use the 30 abdominal CT scans from the MICCAI 2015 Multi-Atlas Abdomen Labeling Challenge. These scans were captured during the portal venous contrast phase with variable volume sizes ($512 \times 512 \times 85$ – $512 \times 512 \times 198$) and field of views (approximately $280 \times 280 \times 280$ mm³– $500 \times 500 \times 650$ mm³). The in-plane resolution varies from 0.54×0.54 mm² to 0.98×0.98 mm², while the slice thickness ranges from 2.5 to 5.0 mm. We report the average experimental results on eight abdominal organs (aorta, gallbladder, spleen, left kidney, right kidney, liver, pancreas, spleen, and stomach) with 5-fold validation.

4.2. Implementation details

In the pre-processing step, we simply normalized the intensities of each 3D image to zero mean and unit variance. In the training phase, we applied data augmentation techniques to reduce potential overfitting, including random rotation of the image by 90° along three dimensions and adjusting the brightness of the top 3% pixels. The training iterations were set to 10⁵. We trained the model using the Adam optimizer with a batch size of 1 and synchronized batch normalization. The initial learning rate was set to $1e - 2$ and was decayed

by $(1 - \frac{\text{current_epoch}}{\text{max_epoch}})^{0.9}$. We also regularized the training with dropout in the transformer cells. All experiments are conducted using a 5-fold cross-validation, based on Pytorch 1.7.1 on a workstation with 2 NVIDIA TITAN RTX GPUs. The data division on the public Covid20 dataset is adopted by following the division strategy given by Qiu et al. (2021).

As aforementioned, our encoder consists of two parts. In the feature extraction module, we used a CNN network with two convolutional layers, one max-pooling layer, and one 1-D fully-connected layer with the direction of $x - y$ plane to z coordinate to convert the representations with the original dimension to a cube. The first convolutional layer, with a kernel size of $3 \times 3 \times 3$, embeds the input 3-D image into local representation maps, while the second convolutional layer project the local representation maps for the second part of the encoder via a linear transformation. The output dimension of the feature extraction module is converted (i.e., reshape) to $X \in \mathcal{R}^{1 \times 256 \times 256 \times 256}$. In the bi-directional transformer module, we utilize multiple transformer cells with the bi-directional self-attention mechanism. Specifically, the input embedding strategy that we adopted is Masked LM (MLM) (Devlin et al., 2018). The Masked LM has been proven to be useful within the previous BERT paper (Vaswani et al., 2017), where the image portion masked in the encoder is matched to that in the loss computation stage. Moreover, since we do not embed the data with the positional encoding in our framework, we require a way to learn the 3D representations through a certain sequence. MLM can well meet this requirement. We set the number of transformer cells as 12, 6, and 6 for Placenta, Covid20, and Synapse datasets, respectively. The number of heads within each transformer cell is 15, where each direction (i.e., $x - y$, $x - z$, and $y - z$ plane) contains five heads to compute self-attention scores. The length of each mask is set to 16, 32, and 32 for the Placenta, Covid20, and Synapse datasets, respectively. Each cube representation is divided into 16 parts in the training phase.

² <https://zenodo.org/record/3757476#.Y1NGmy1h1B1>

³ <https://www.synapse.org/ISynapse:syn3193805/wiki/217789>

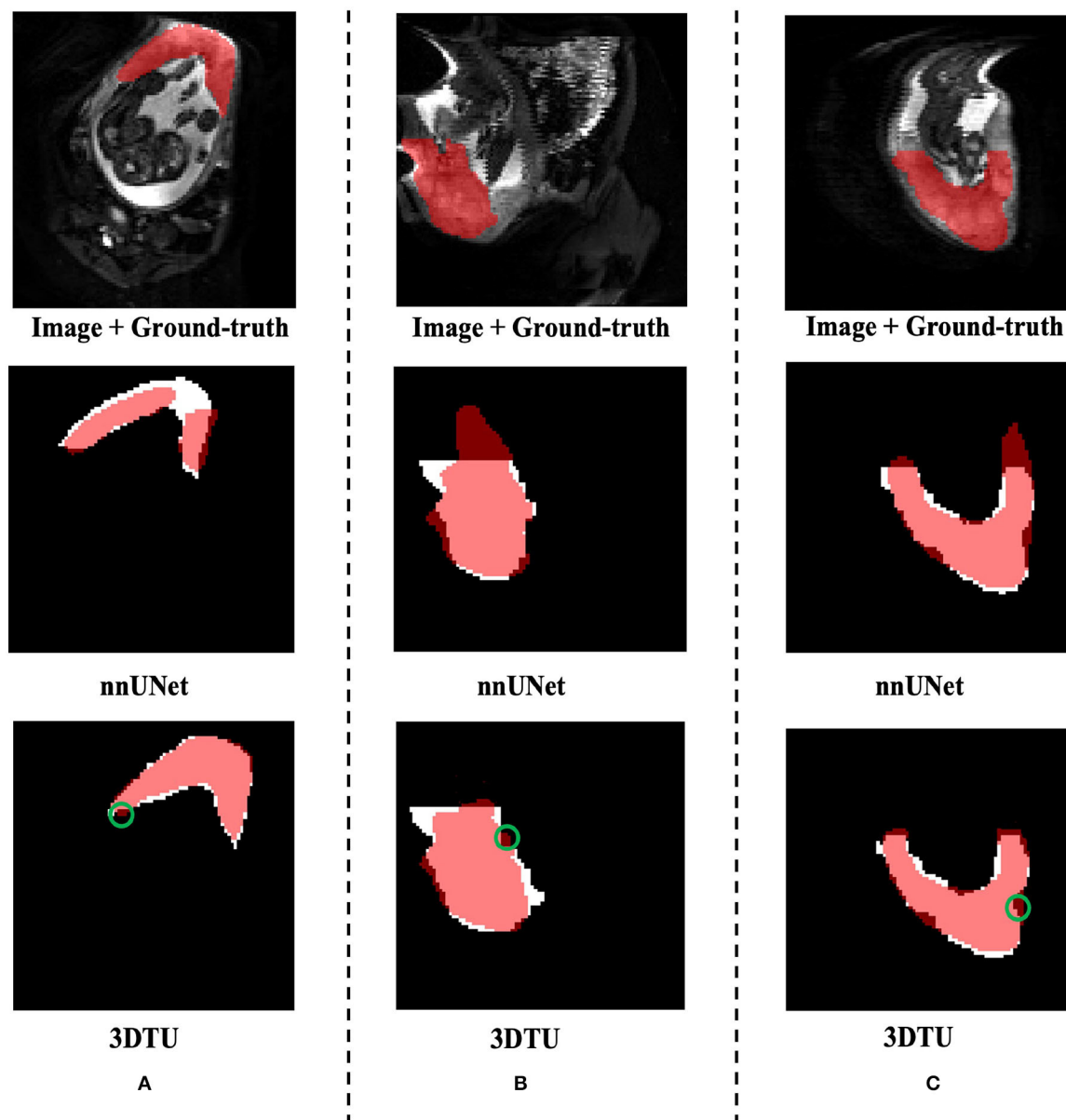


FIGURE 3

Visualization of the segmentation results on the Placenta dataset produced by our 3DTU and nnUNet. Columns (A–C) show the x – y plane, y – z plane, and x – z plane of 3D segmentation predictions, respectively. The true-positive regions are highlighted in pink. The false-negative regions are highlighted in red (e.g., the green circle regions in the last row). Better view with colors and zooming in.

4.3. Baseline settings and evaluation metrics

To evaluate our 3DTU's performance, we choose the following frameworks as baselines: 2DU-Net (Ronneberger et al., 2015), 3D U-Net (Çiçek et al., 2016), UNet++ (Zhou et al., 2019), TransUNet (Chen et al., 2021), ViT (visual transformer) (Dosovitskiy et al., 2020), nnFormer (Zhou et al., 2021), and nnUNet (Isensee et al., 2021). Both 2D

and 3D UNet are FCN-based encoder–decoder structures with convolutional blocks and skip-connections between the encoder and decoder. The UNet++ is a nested-connected encoder–decoder structure, where each convolutional block is connected to all other blocks. The TransUNet is an encoder–decoder network, where the encoder of UNet is replaced by a 2D transformer including a positional embedding scheme followed by a visual transformer (ViT). The nnFormer is a 3D UNet-type framework that

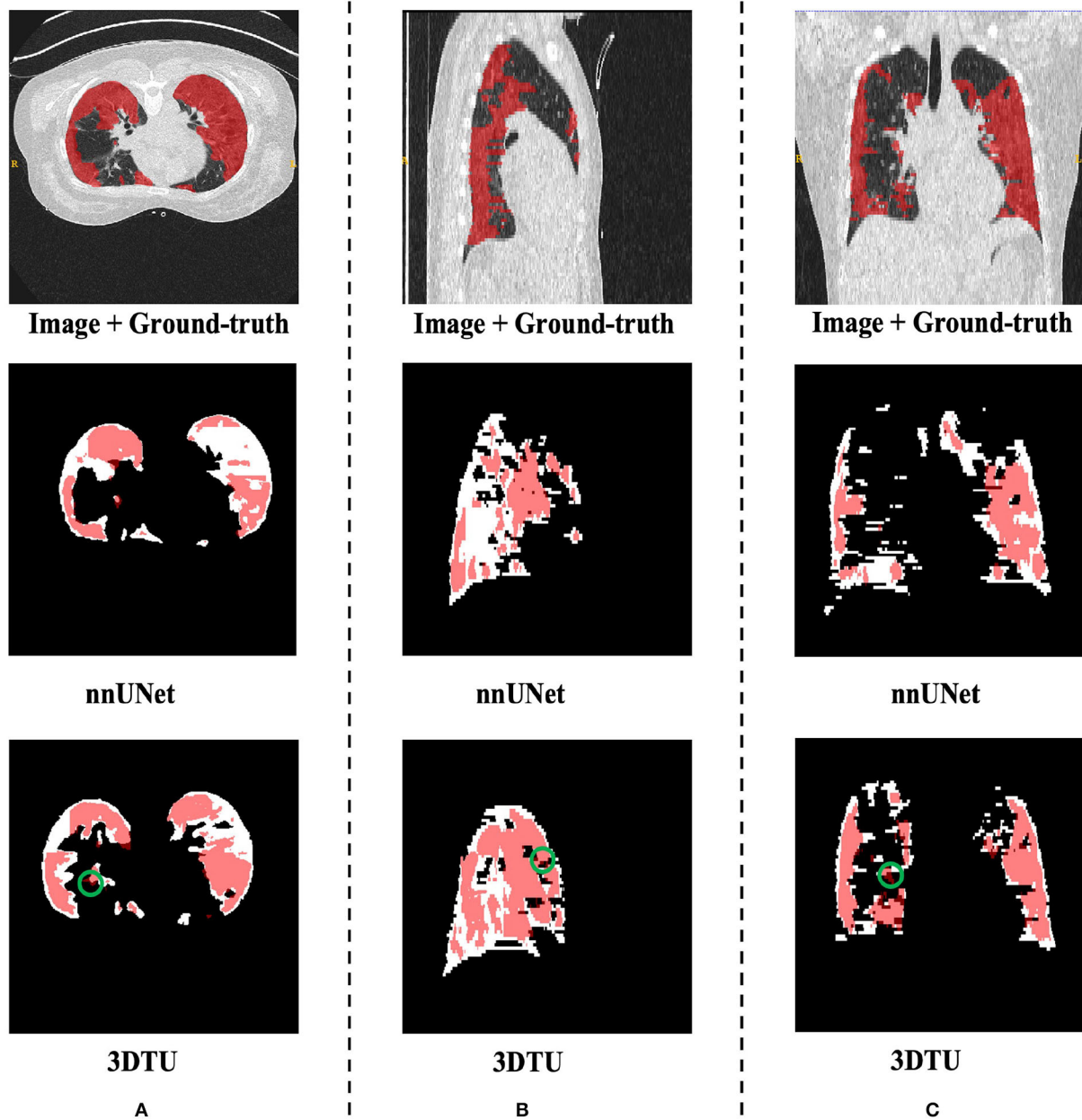


FIGURE 4
Visualization of the infection segmentation results on the Covid20 dataset produced by our 3DTU and nnUNet. Columns (A–C) show the x–y plane, y–z plane, and x–z plane of 3D segmentation predictions, respectively. The true-positive regions are highlighted in pink. The false-negative regions are highlighted in red (e.g., the green circle regions in the last row). Better view with colors and zooming in.

replaces the convolutional blocks with three different novel attention mechanisms.

The metrics we used to evaluate our 3DTU include mIoU, DICE score, and Hausdorff Distance (HD). IoU is the area of overlap between the predicted segmentation and the ground truth divided by the area of union between them. For binary (two classes) or multi-class segmentation, the mean IoU (mIoU)

of the image is calculated by taking the IoU of each class and averaging them. DICE score is the harmonic mean of precision and recall of the segmentation results. mIoU and DICE scores are two overlap-based metrics measuring the similarity between the ground truths and segmentation predictions. The range of mIoU and DICE scores is from 0 to 1 and the larger value indicates better segmentation performance. The directed

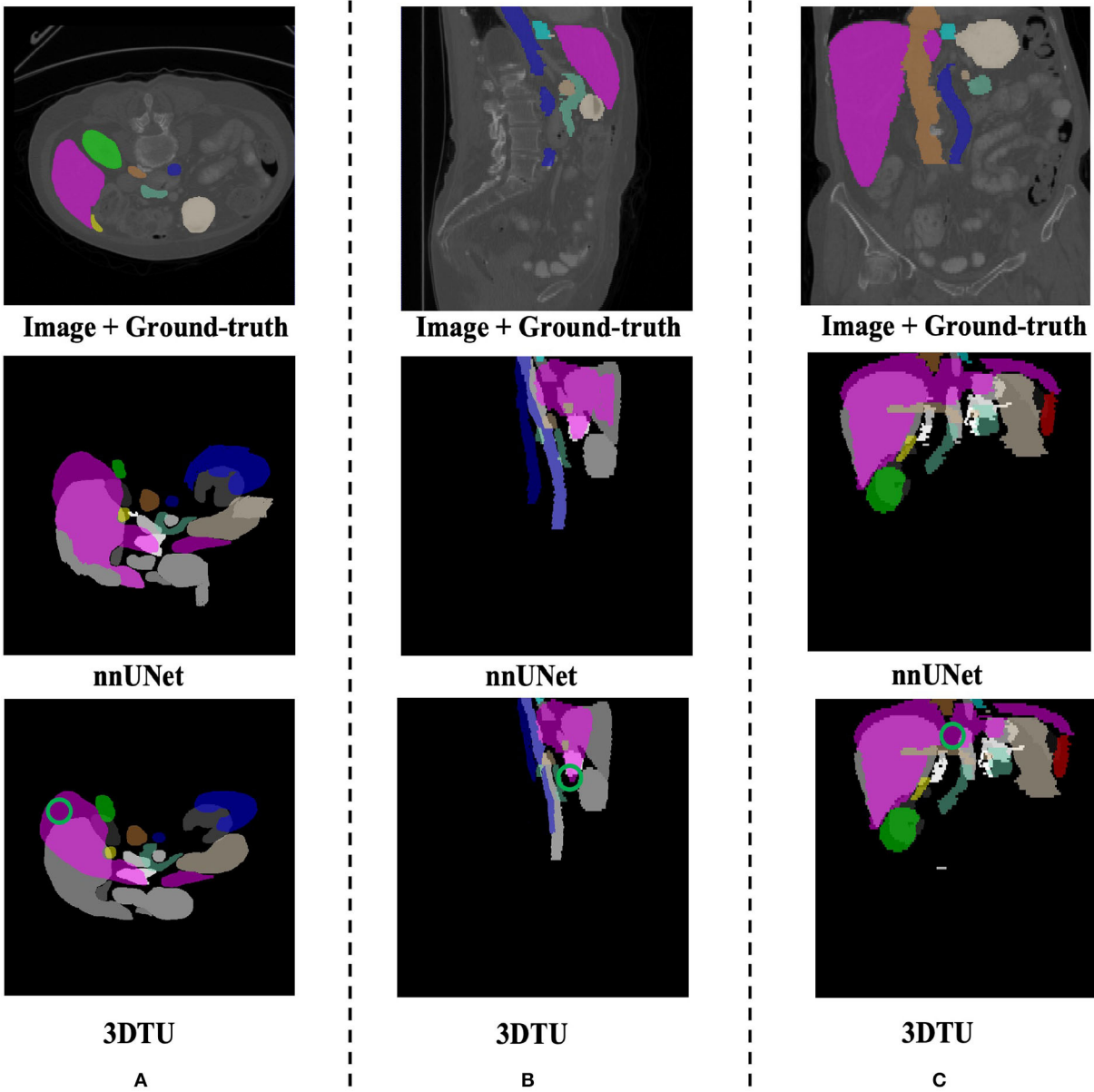


FIGURE 5
Visualization of the segmentation results on the Synapse dataset produced by our 3DTU and nnUNet. Columns (A–C) show the x–y plane, y–z plane, and x–z plane of 3D segmentation predictions, respectively. The green circle indicates part of false-negative regions. Better view with colors and zooming in.

average Hausdorff distance (HD) from point set X to Y is computed by the sum of all minimum distances from all points from point set X to Y divided by the number of points in X . HD is a shape distance-based metric, which measures the dissimilarity between the surfaces of the segmentation results and the related ground truths. A lower value of HD indicates better performance.

4.4. Comparative experiments

Table 1 provides the performance of our proposed 3DTU and the six competing baselines, including 2D UNet (Ronneberger et al., 2015), 3D UNet (Ronneberger et al., 2015), UNet++ (Zhou et al., 2019), TransUNet (Chen et al., 2021), visual transformer (ViT) (Dosovitskiy et al., 2020), and

TABLE 2 Dice scores (in %) of our 3DTU on three datasets.

DICE score	Placenta dataset	Covid20 dataset	Synapse dataset
CNN + UNet decoder	68.6	74.3	59.5
BiT + UNet decoder	66.9	72.8	70.2
CNN + BiT	80.0	89.2	65.1
3DTU	84.0	92.0	87.3

The best results are shown in bold.

TABLE 3 Dice scores (in %) of our 3DTU running on data that has been preprocessed with/without positional encoding.

	Placenta dataset	Covid20 dataset	Synapse dataset
3DTU w/o Positional encoding	84.0	92.0	87.3
3DTU with Positional encoding	82.7	92.1	86.8

nnFormer (Zhou et al., 2021) on the Placenta and Covid20 datasets. It shows that our 3DTU outperforms all competing baseline methods consistently in terms of mIOU and DICE scores on both datasets, while beating most of the methods in the baseline in the Synapse dataset, indicating that the segmentation results of our models match well with the ground truth. For example, our proposed 3DTU outperforms baselines with at least 0.48% and 0.44% increases in DICE scores on Placenta and Covid20 datasets, respectively. This may attribute to the attention mechanism proposed in the 3DTU, which can compute the attention scores from three different directions to yield discriminative and anisotropic semantic features for 3D images. In general, the transformer-based methods (e.g., TransUNet, ViT, etc.) perform better than the other baseline methods. In addition, we visualized the segmentation results of our 3DTU and the best baseline method (i.e., nnUNet) on three datasets in Figures 3–5, respectively.

4.5. Ablation study

We conducted an ablation study on both datasets (i.e., Placenta and Covid20) to evaluate the effectiveness of each part in our 3DTU framework. Our 3DTU is an encoder-decoder-based framework, where the encoder consists of a CNN networks part as well as a bi-directional transformer (BiT) part, where the decoder is in the UNet decoder setting. Hence, we designed the following four experiments in our ablation study.

- We removed the CNN networks in the encoder and directly fed the input images to the BiT part.
- We removed the BiT part in the encoder and directly connected the CNN networks to the UNet decoder.

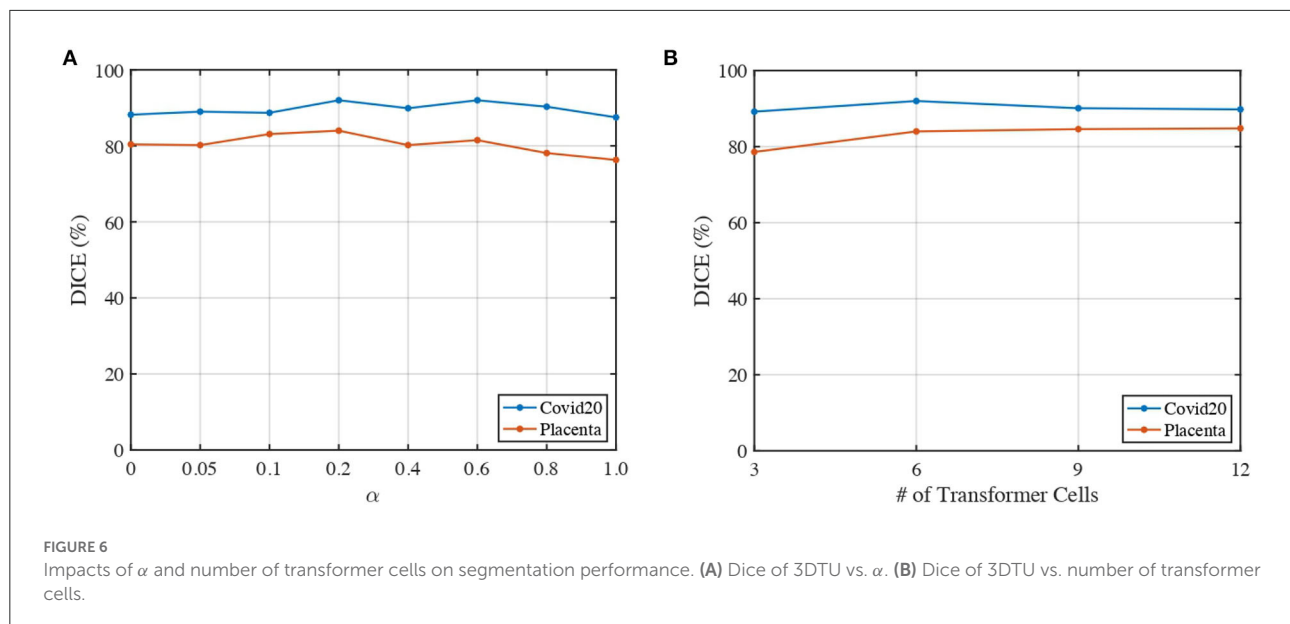
- We removed the UNet decoder part and considered the BiT as both (part of) encoder and decoder.⁴
- We designed a comparative experiment where we trained 3DTU with positional encoded representations. We encoded the representations at the input of the transformer encoder.

The results in Table 2 show the effectiveness and necessity of all the sub-parts in our 3DTU. The results in Table 3 indicate that positional encoding is not necessary in our framework since our attention mechanism can process the 3D data as a whole. Compared with the 3DTU w/o positional encoding, the segmentation dice scores yielded by 3DTU with positional encoding are not changed or even decreased. When we removed the CNN networks and only utilized BiT as the encoder (see results of BiT+Unet decoder in Table 2), the segmentation performance decreased on both datasets (e.g., DICE decrease from 84.0 to 66.9% and from 92.0 to 72.8% on Placenta and COVID datasets, respectively). This indicates an essential role of CNN-based convolutional layers in the encoder, without which the self-attention transformer layers may not localize the raw image pixels precisely. Meanwhile, the segmentation performance increase when we use BiT instead of UNet as a decoder (see results of CNN + UNet Decoder and CNN + BiT). This manifests that, compared with UNet-based methods, the (bi-directional) transformers are more powerful in boosting the segmentation results.

4.6. Parameter analysis

We analyze the impact of two parameters, including the loss weights α and the number of transformer cells, on the segmentation performance of our proposed 3DTU across two datasets in Figure 6. In general, Figure 6 indicates that the segmentation results performed by our 3DTU are consistent. Figure 6A shows that the dice results increase and then decrease with the increase of α from 0 to 1. The best dice scores are achieved when $\alpha = 0.2$ on both Placenta and Covid20 datasets. Figure 6B shows that the segmentation performance improves when increasing the number of transformer cells from 3 to 6. However, the performance will keep stable (on the Placenta dataset) or even slightly decrease (on the Covid20 dataset) when the framework goes deeper. The reason for the slight decrease in the performance of the Covid20 dataset may result from the small size of the dataset. Only 20 3D images are included in the Covid20 dataset, which may not facilitate the training process when the network goes deep. Moreover, our 3DTU has a total of 70M parameters (when training on the Covid20 dataset and the Synapse dataset), which is more than 2D UNet (7M) and

⁴ It shows in Devlin et al. (2018) that the bi-directional transformer can serve as both encoder and decoder.



3D UNet (17M) but beats the other transformer-based or hybrid framework in the baseline (the TransUNet has 80M parameters, and nnFormer has 158M parameters).

5. Conclusion

In this article, we propose a novel 3D transformer UNet (3DTU) framework to capture global contextual information for 3D medical image segmentation. A new attention mechanism is proposed with our 3DTU framework, which is especially suitable for computing self-attentions for 3D objects. The experimental results on two 3D medical image datasets demonstrate that our method can outperform several state-of-the-art segmentation baselines. In the future, we plan to explore how to reduce the computation loads in transformer layers, which may improve the efficiency of most current transformer-based methods.

Data availability statement

The Covid20 dataset is from the community of Coronavirus Disease Research-COVID-19 (Jun et al., 2021) and is available via <https://zenodo.org/record/3757476#.Y1NGmy1h1B1>. The Placenta dataset is available upon request.

Author contributions

XF took charge of conception, design, method implementation, statistical analysis, and manuscript writing. ZS and YW took charge of data collection and preprocessing. ZS, EZ, HH, and YW took charge of experimental design, results discussion, and manuscript proofreading. HT and LZ

took charge of project design, analysis, interpretation, and manuscript writing/revising. All authors contributed to the article and approved the submitted version.

Funding

This project was partially supported by NSF IIS 2045848 and NIH/NICHD (R01HD094381 and R01HD104822), as well as by the Burroughs Wellcome Fund Preterm Birth Initiative (NGP10119) and the Bill & Melinda Gates Foundation (INV-005417, INV-035476, and INV-037302).

Acknowledgments

We thank the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National Science Foundation (NSF) grant number ACI-1548562 and NSF award number ACI-1445606, which provide the computation resources based on the Pittsburgh Supercomputing Center (PSC) for part of our work. We would like to appreciate the efforts devoted by the community of Coronavirus Disease Research-COVID-19 and Zenodo to collect and share the COVID-19 CT image dataset. Meanwhile, we appreciate the Washington University in Saint Louis for collecting and sharing the data Placenta MRI dataset for our segmentation algorithm evaluations.

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated

organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

References

- Badrinarayanan, V., Kendall, A., and Cipolla, R. (2017). Segnet: a deep convolutional encoder-decoder architecture for image segmentation. *IEEE Trans. Pattern. Anal. Mach. Intell.* 39, 2481–2495. doi: 10.1109/TPAMI.2016.2644615
- Cao, H., Wang, Y., Chen, J., Jiang, D., Zhang, X., Tian, Q., et al. (2021). Swin-Unet: unet-like pure transformer for medical image segmentation. *arXiv preprint arXiv:2105.05537*. doi: 10.48550/arXiv.2105.05537
- Chen, J., Lu, Y., Yu, Q., Luo, X., Adeli, E., Wang, Y., et al. (2021). TransUnet: transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306*. doi: 10.48550/arXiv.2102.04306
- Chen, L.-C., Zhu, Y., Papandreou, G., Schroff, F., and Adam, H. (2018). "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proceedings of the European Conference on Computer Vision (ECCV)* (Cham: Springer), 801–818.
- Chen, Y., Rohrbach, M., Yan, Z., Shuicheng, Y., Feng, J., and Kalantidis, Y. (2019). "Graph-based global reasoning networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (Long Beach, CA: IEEE), 433–442.
- Chen, Y., Wang, Z., Peng, Y., Zhang, Z., Yu, G., and Sun, J. (2018). "Cascaded pyramid network for multi-person pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (Salt Lake City, UT: IEEE), 7103–7112.
- Cho, K., Van Merriënboer, B., Bahdanau, D., and Bengio, Y. (2014). On the properties of neural machine translation: encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*. doi: 10.3115/v1/W14-4012
- Çiçek, Ö., Abdulkadir, A., Lienkamp, S. S., Brox, T., and Ronneberger, O. (2016). "3D U-Net: learning dense volumetric segmentation from sparse annotation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention* (Athens: Springer), 424–432.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*. doi: 10.48550/arXiv.1810.04805
- Diakogiannis, F. I., Waldner, F., Caccetta, P., and Wu, C. (2020). Resunet-a: a deep learning framework for semantic segmentation of remotely sensed data. *ISPRS J. Photogram. Remote Sens.* 162, 94–114. doi: 10.1016/j.isprsjrs.2020.01.013
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., et al. (2020). An image is worth 16x16 words: transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*. doi: 10.48550/arXiv.2010.11929
- Gao, Y., Zhou, M., and Metaxas, D. N. (2021). "Unet: a hybrid transformer architecture for medical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention* (Strasbourg: Springer), 61–71.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., et al. (2014). "Generative adversarial nets," *Advances in Neural Information Processing Systems* 27. Montreal, QC.
- Hatamizadeh, A., Nath, V., Tang, Y., Yang, D., Roth, H., and Xu, D. (2022a). Swin unet: Swin transformers for semantic segmentation of brain tumors in mri images. *arXiv preprint arXiv:2201.01266*. doi: 10.1007/978-3-031-08999-2_22
- Hatamizadeh, A., Tang, Y., Nath, V., Yang, D., Myronenko, A., Landman, B., et al. (2022b). "UNETR: transformers for 3d medical image segmentation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision* (Waikoloa, HI: IEEE), 574–584.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (Las Vegas, NV: IEEE), 770–778.
- Hochreiter, S., and Schmidhuber, J. (1997). Long short-term memory. *Neural Comput.* 9, 1735–1780. doi: 10.1162/neco.1997.9.8.1735
- Huang, H., Lin, L., Tong, R., Hu, H., Zhang, Q., Iwamoto, Y., et al. (2020). "Unet 3+: a full-scale connected unet for medical image segmentation," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (Barcelona: IEEE), 1055–1059.
- Isensee, F., Jaeger, P. F., Kohl, S. A., Petersen, J., and Maier-Hein, K. H. (2021). NNU-Net: a self-configuring method for deep learning-based biomedical image segmentation. *Nat. Methods* 18, 203–211. doi: 10.1038/s41592-020-01008-z
- Ji, G.-P., Chou, Y.-C., Fan, D.-P., Chen, G., Fu, H., Jha, D., et al. (2021). "Progressively normalized self-attention network for video polyp segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention* (Strasbourg: Springer), 142–152.
- Jia, H., Tang, H., Ma, G., Cai, W., Huang, H., Zhan, L., et al. (2021). PSGR: pixel-wise sparse graph reasoning for covid-19 pneumonia segmentation in ct images. *arXiv preprint arXiv:2108.03809*. doi: 10.48550/arXiv.2108.03809
- Jun, M., Yixin, W., Xingle, A., Cheng, G., Ziqi, Y., Jianan, C., et al. (2021). Toward data-efficient learning: A benchmark for COVID-19 CT lung and infection segmentation. *Med. Phys.* 48, 1197–1210.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems* 25. Lake Tahoe.
- Li, J., Chen, J., Tang, Y., Landman, B. A., and Zhou, S. K. (2022). Transforming medical imaging with transformers? a comparative review of key properties, current progresses, and future perspectives. *arXiv preprint arXiv:2206.01136*. doi: 10.48550/arXiv.2206.01136
- Li, X., Yang, Y., Zhao, Q., Shen, T., Lin, Z., and Liu, H. (2020). "Spatial pyramid based graph reasoning for semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (Seattle: IEEE), 8950–8959.
- Li, Y., and Gupta, A. (2018). "Beyond grids: learning graph representations for visual recognition," in *Advances in Neural Information Processing Systems* 31. Montréal.
- Lin, A., Chen, B., Xu, J., Zhang, Z., Lu, G., and Zhang, D. (2022). Ds-transunet: dual swin transformer u-net for medical image segmentation. *IEEE Trans. Instrum. Meas.* 71, 4005615. doi: 10.1109/TIM.2022.3178991
- Liu, M., Cheng, D., Wang, K., and Wang, Y. (2018). Multi-modality cascaded convolutional neural networks for alzheimer's disease diagnosis. *Neuroinformatics* 16, 295–308. doi: 10.1007/s12021-018-9370-4
- Liu, M., Maiti, P., Thomopoulos, S., Zhu, A., Chai, Y., Kim, H., et al. (2021). "Style transfer using generative adversarial networks for multi-site mri harmonization," in *International Conference on Medical Image Computing and Computer-Assisted Intervention* (Strasbourg: Springer), 313–322.
- Liu, M., Zhu, A., Maiti, P., Thomopoulos, S. I., Gadewar, S., Chai, Y., et al. (2022). Style transfer generative adversarial networks to harmonize multi-site mri to a single reference image to avoid over-correction. *bioRxiv*. doi: 10.1101/2022.09.12.506445
- Long, J., Shelhamer, E., and Darrell, T. (2015). "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (Boston, MA: IEEE), 3431–3440.
- Noh, H., Hong, S., and Han, B. (2015). "Learning deconvolution network for semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision* (Santiago: IEEE), 1520–1528.
- Pan, X., Zhao, Y., Chen, H., Wei, D., Zhao, C., and Wei, Z. (2020). Fully automated bone age assessment on large-scale hand x-ray dataset. *Int. J. Biomed. Imaging* 2020, 8460493. doi: 10.1155/2020/8460493
- Parmar, N., Vaswani, A., Uszkoreit, J., Kaiser, L., Shazeer, N., Ku, A., et al. (2018). "Image transformer," in *International Conference on Machine Learning*. (Stockholmsmässan: PMLR), 4055–4064.
- Petit, O., Thome, N., Rambour, C., Themyr, L., Collins, T., and Soler, L. (2021). "U-net transformer: self and cross attention for medical image segmentation," in *International Workshop on Machine Learning in Medical Imaging*. (Strasbourg: Springer), 267–276.

- Qiu, Y., Liu, Y., Li, S., and Xu, J. (2021). Miniseg: an extremely minimum network for efficient COVID-19 segmentation. *Proc. AAAI Conf. Artif. Intell.* 35, 4846–4854. doi: 10.1609/aaai.v35i6.16617
- Ronneberger, O., Fischer, P., and Brox, T. (2015). “U-Net: convolutional networks for biomedical image segmentation,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention* (Munich: Springer), 234–241.
- Simonyan, K., and Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*. doi: 10.48550/arXiv.1409.1556
- Sun, Z., Wu, W., Zhao, P., Wang, Q., Woodard, P., Nelson, D., et al. (2022). Dual-contrast mri reveals intraplacental oxygenation patterns, detects placental abnormalities and fetal brain oxygenation. *Ultrasound Obstetr. Gynecol.* doi: 10.1002/uog.24959
- Tokunaga, H., Teramoto, Y., Yoshizawa, A., and Bise, R. (2019). “Adaptive weighting multi-field-of-view cnn for semantic segmentation in pathology,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (Long Beach, CA: IEEE), 12597–12606.
- Valanarasu, J. M. J., Oza, P., Hacihaliloglu, I., and Patel, V. M. (2021). “Medical transformer: gated axial-attention for medical image segmentation,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention* (Strasbourg: Springer), 36–46.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., et al. (2017). “Attention is all you need,” in *Advances in Neural Information Processing Systems 30* (Long Beach Convention & Entertainment Center).
- Wang, W., Chen, C., Ding, M., Yu, H., Zha, S., and Li, J. (2021). “Transbts: multimodal brain tumor segmentation using transformer,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention* (Strasbourg: Springer), 109–119.
- Xie, Y., Zhang, J., Shen, C., and Xia, Y. (2021). “COTR: efficiently bridging cnn and transformer for 3D medical image segmentation,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention* (Strasbourg: Springer), 171–180.
- Yang, Z., Hu, Z., Salakhutdinov, R., and Berg-Kirkpatrick, T. (2017). “Improved variational autoencoders for text modeling using dilated convolutions,” in *International Conference on Machine Learning* (Long Beach Convention & Entertainment Center, PMLR), 3881–3890.
- Yu, F., and Koltun, V. (2015). Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*. doi: 10.48550/arXiv.1511.07122
- Yun, B., Wang, Y., Chen, J., Wang, H., Shen, W., and Li, Q. (2021). Spectr: spectral transformer for hyperspectral pathology image segmentation. *arXiv preprint arXiv:2103.03604*. doi: 10.48550/arXiv.2103.03604
- Zhang, J., Zhou, L., Wang, L., Liu, M., and Shen, D. (2022). Diffusion kernel attention network for brain disorder classification. *IEEE Trans. Med. Imaging* 41, 2814–2827. doi: 10.1109/TMI.2022.3170701
- Zhang, X., Zou, Y., and Shi, W. (2017). “Dilated convolution neural network with leakyrelu for environmental sound classification,” in *2017 22nd International Conference on Digital Signal Processing (DSP)* (London: IEEE), 1–5.
- Zhou, H.-Y., Guo, J., Zhang, Y., Yu, L., Wang, L., and Yu, Y. (2021). nnFormer: interleaved transformer for volumetric segmentation. *arXiv preprint arXiv:2109.03201*. doi: 10.48550/arXiv.2109.03201
- Zhou, Z., Siddiquee, M. M. R., Tajbakhsh, N., and Liang, J. (2019). Unet++: redesigning skip connections to exploit multiscale features in image segmentation. *IEEE Trans. Med. Imaging* 39, 1856–1867. doi: 10.1109/TMI.2019.2959609