3D Medical Axial Transformer: A Lightweight Transformer Model for 3D Brain Tumor Segmentation

Author name(s) withheld

EMAIL(S) WITHHELD

Address withheld

Editors: Under Review for MIDL 2023

Abstract

In recent years, Transformer-based models have gained attention in the field of medical image segmentation, with research exploring ways to integrate them with established architectures such as Unet. However, the high computational demands of these models have led most current approaches to focus on segmenting 2D slices of MRI or CT images, which can limit the ability of the model to learn semantic information in the depth axis and result in output with uneven edges. Additionally, the small size of medical image datasets, particularly those for brain tumor segmentation, poses a challenge for training transformer models. To address these issues, we propose 3D Medical Axial Transformer (MAT), a lightweight, end-to-end model for 3D brain tumor segmentation that employs an axial attention mechanism to reduce computational demands and -distillation to improve performance on small datasets. Results indicate that our approach, which has fewer parameters and a simpler structure than other models, achieves superior performance and produces clearer output boundaries, making it more suitable for clinical applications.

Keywords: Deep learning, 3D brain tumor segmentation, 3D Transformer, axial attention, self-distillation

1. Introduction

Medical image segmentation is a key component in computer-aided diagnosis and a fundamental procedure in medical image processing (Doi, 2007). It helps clinicians make more accurate diagnoses and treatment decisions by segmenting organs or tumors in medical scans. With the development of convolutional neural networks (CNNs), Unet (Ronneberger et al., 2015) emerged as a popular medical image segmentation network with its simple U-shaped structure and innovative skip connections design. Many variations of Unet have been developed, including V-Net (Milletari et al., 2016), Res-Unet (Zhang et al., 2018), H-Dense-Unet (Li et al., 2018) and 3D-Unet (Çiçek et al., 2016) for 3D medical image segmentation. Specifically, the nnUnet (Isensee et al., 2021) model has demonstrated stateof-the-art performance in a wide range of tasks, including brain tumor segmentation. CNNs have achieved success in medical image segmentation, but struggle with long-range dependencies and processing global context (Vaswani et al., 2017). This is especially problematic in brain tumor images, where doctors need to combine information from multiple regions for diagnosis (Valanarasu et al., 2021).

Transformer-based models, utilizing self-attention mechanisms, have become popular in natural language processing and set state-of-the-art benchmarks in recent years (Devlin et al., 2018; Brown et al., 2020). These models are able to efficiently compute dependencies between sequential inputs, even when distant from each other, addressing the problem of

WITHHELD

long-range dependencies which traditional convolutional models struggle with. In computer vision, the Vision Transformer (ViT) (Dosovitskiy et al., 2020) adapted the Transformer for image classification tasks with successful results. Subsequently, Swin-Transformer (Liu et al., 2021) has combined the sliding window concept of CNNs with self-attention mechanisms to form the Transformer using shifted window methods, capable of handling downstream tasks such as classification, detection, and segmentation.

In the field of medical image segmentation, Trans-UNet (Chen et al., 2021) pioneered the use of transformer-based models by integrating them into the traditional Unet architecture. Another notable approach is Medical Transformer (MedT)(Valanarasu et al., 2021), which uses gated axial attention for efficient semantic feature learning. Additionally, Swin-Unet (Cao et al., 2021) leveraged the efficient structure of the Swin-Transformer to achieve superior performance. However, the high computational cost of Transformer models with 3D inputs has limited their application to mainly 2D slices. UNETR(Hatamizadeh et al., 2022b) and Swin-UNETR (Hatamizadeh et al., 2022a) are exception, achieving good results on 3D medical images, but at the cost of high GPU memory consumption. To mitigate this, CoTr (Xie et al., 2021) proposed dividing 3D images into smaller blocks, but this also sacrifices semantic information. AFTer-Unet (Yan et al., 2022) utilized several neighboring axial slices as a 3D input image, reducing resource overhead but limiting global dependencies and requiring extra pre-processing.

In this study, we propose the 3D Medical Axial Transformer (MAT) model for efficient 3D brain tumor segmentation. Building on the success of MedT (Valanarasu et al., 2021), we introduce three one-dimensional gated axial attention mechanisms within the Transformer module to decompose the three-dimensional attentions, reducing GPU resource consumption and memory cost. Unlike AFTer-Unet, we focus on the entire depth axis rather than on neighboring axial slices. Additionally, we use a self-distillation training technique to improve performance on small datasets by using the model's output on the previous mini-batch as a soft target. Our main contributions are: (1) an end-to-end model that eliminates the need for pre-processing and allows for the direct use of 3D images, (2) a 3D self-attention module utilizing axial attention mechanism, which effectively reduces resource consumption, and (3) the use of self-distillation with a warm-up schedule, a novel approach in medical image processing, which helps the transformer module learn information with small-scale datasets.

2. Methodology

The proposed 3D Medical Axial Transformer (MAT) model utilizes a Transformer architecture with axial attention mechanisms in the encoder component, in combination with CNNs. The encoder is connected to a decoder, composed of CNNs, through skip connections at various resolutions. It is noteworthy that, due to the constraints imposed by the 3D input on the mini-batch size, group normalization (GN) is employed as a regularization technique in MAT, as opposed to batch normalization. The overall architecture of the model is depicted in Figure 1(a). In this section, we present a thorough description of the encoder and decoder of the proposed MAT model.



Figure 1: (a) Architecture overview of MAT (b) designs of Axial Transformer Block (c) Schematic diagram of the Axial Attention calculation method

2.1. CNN Encoder

MAT utilizes a minimal number of CNN blocks to extract features before passing the input to the Transformer encoder, following the design of MedT (Valanarasu et al., 2021). This approach allows the Transformer encoder to better learn semantic information and reduces the number of model parameters. The CNN encoder consists of one $(7 \times 7 \times 7)$ convolutional block and two $(5 \times 5 \times 5)$ convolutional blocks. The initial convolutional block $(7 \times 7 \times 7)$ with a stride of 2 improves feature extraction, as established in (Simonyan and Zisserman, 2014), with a group normalization layer added between the convolutions. In contrast to MedT, an average pooling layer of size $(2 \times 2 \times 2)$ with stride 2 is implemented between the CNN encoder and the Transformer encoder to extract further features while conserving memory resources.

Suppose the input to the encoder is a 3D image $x \in \mathbb{R}^{C \times D \times H \times W}$, where $(D \times H \times W)$ is the image resolution and C is the channel of the input (e.g. C = 4 channels for MRI). After the CNN encoder, the output should be $x_{CNN} \in \mathbb{R}^{C^{cnn} \times \frac{D}{4} \times \frac{H}{4} \times \frac{W}{4}}$, where C^{cnn} is the number of channels output from the last convolution layer.

2.2. Axial Transformer Encoder

After using CNNs to extract shallow image features, we introduce Transformer encoders to facilitate the learning of deep semantic information. To reduce computational complexity, we implement the Transformer blocks using 3D axial attention, known as Axial Transformer, which allows for self-attention computation in all three dimensions, resulting in a comprehensive 3D modeling of the images, as illustrated in Figure 1(b). This design choice

effectively reduces resource consumption while enabling the model to effectively process and understand the full context of the medical images.

2.2.1. AXIAL ATTENTION

In this work, we propose the use of three axial attention mechanisms in MAT for efficient 3D self-attention computation. This approach decomposes the 3D attention calculation into three 1D calculations in the three dimensions, reducing computational complexity, as illustrated in Figure 1(c). We also incorporate a learnable positional bias term and relative positional encoding in the self-attention module to effectively capture positional information in the 3D medical images. This approach builds upon previous research (Ho et al., 2019; Wang et al., 2020) that have shown the effectiveness of axial attention in capturing semantic information in 3D images. As an example, the output of the height axial attention is illustrated as Eq (1), Where r^q , r^k , and r^v are the relative position encoding for queries, keys, and values respectively.

$$y_{\text{height}_{ijk}} = \sum_{h=1}^{H} \text{Softmax} \left(q_{ijk}^{T} k_{ihk} + q_{ijk}^{T} r_{ihk}^{q} + k_{ihk}^{T} r_{ihk}^{k} \right) \left(v_{ihk} + r_{ihk}^{v} \right) \tag{1}$$

To address the challenge of training the Transformer on small datasets and ensuring adequate position encoding, we incorporate a gating mechanism, as proposed in previous research (Valanarasu et al., 2021). The gated axial attention on the height axis is as follows.

$$y_{\text{height}_{ijk}} = \sum_{h=1}^{H} \text{Softmax} \left(q_{ijk}^{T} k_{ihk} + G_q q_{ijk}^{T} r_{ihk}^{q} + G_k k_{ihk}^{T} r_{ihk}^{k} \right) \left(G_v^1 v_{ihk} + G_v^2 r_{ihk}^v \right)$$
(2)

Where G_q , G_k , G_v^1 , and G_v^2 are gates for queries, keys, and values, respectively. These gates act as learnable parameters that effectively control the final output of attention. In general, the gating mechanism limits the output of poor position encoding and gives higher weights to those that have been learned relatively well (Valanarasu et al., 2021).

2.2.2. Architecture

Our Axial Transformer block is based on the traditional Transformer design and features 3D convolutional layers, group normalization, and axial attention for height, width, and depth, as depicted in Figure 1(b). The output of the last GN layer is connected to the input via a skip connection (Vaswani et al., 2017; He et al., 2016).

Medical images have been found to require higher accuracy rather than a complex method for processing, based on research of various datasets (Isensee et al., 2021). To address this, we have chosen to use a lightweight model with fewer blocks to avoid issues such as overfitting or difficulty in practical application. We use 3 modules for axial attention computation and divide it into three stages: S1, S2, and S3, with an average pooling layer (strides = 2) behind each of the modules to reduce feature map resolution, except for T1 to prevent loss of important semantic information at an early stage. The input is assumed to be $x_{CNN} \in \mathbb{R}^{C^{cnn} \times \frac{D}{4} \times \frac{H}{4} \times \frac{W}{4}}$ (the output of the CNN encoder).

The input is assumed to be $x_{CNN} \in \mathbb{R}^{C^{cnn} \times \frac{D}{4} \times \frac{H}{4} \times \frac{W}{4}}$ (the output of the CNN encoder). Each stage expands the number of channels by a factor of two through multi-headed attention, so the output of each stage turns to be $(2C^{cnn} \times \frac{D}{4} \times \frac{H}{4} \times \frac{W}{4}), (4C^{cnn} \times \frac{D}{8} \times \frac{H}{8} \times \frac{W}{8})$ and $(8C^{cnn} \times \frac{D}{16} \times \frac{H}{16} \times \frac{W}{16}).$

2.3. CNN Decoder

Our CNN decoder is designed using the Unet architecture (Ronneberger et al., 2015), with a structure largely symmetric to the Axial Transformer encoder. It comprises five CNNinterpolation blocks, where the first block has, a stride of 1, while the remaining blocks have a stride of 2. The ReLU activation function is employed, and the image resolution is gradually increased by a factor of two with each block on the last four blocks. A skip connection is utilized to connect the encoder and decoder to restore details lost during downsampling. Instead of the traditional four skip connections in the up-sampling blocks, we used only three, discarding the topmost one as suggested by (Guo et al., 2022) to avoid negative impact from noisy shallow layer semantic information.

Let the final output of Axial Transformer encoder be $X_{encoder}$ with a resolution of $\left(\frac{D}{16} \times \frac{H}{16} \times \frac{W}{16}\right)$, then the output after five CNN-interpolation blocks will be $\left(\frac{D}{16} \times \frac{H}{16} \times \frac{W}{16}\right)$, $\left(\frac{D}{8} \times \frac{H}{8} \times \frac{W}{8}\right)$, $\left(\frac{D}{4} \times \frac{H}{4} \times \frac{W}{4}\right)$, $\left(\frac{D}{2} \times \frac{H}{2} \times \frac{W}{2}\right)$ and $(D \times H \times W)$, which means that the resolution reduction is completed. At the last, a convolutional layer is used for classification, thus completing the downstream task of image segmentation.

2.4. Self-distillation for Regularization with warm-upped DLB

To address the challenge of training Transformer-based models on small datasets, we propose the use of self-distillation, a technique that utilizes the model's own output as a soft target for learning (Zhang et al., 2019). This approach allows for the optimization of the model directly through the training schedule, without the need for extensive modification of the architecture or the use of a large teacher model like traditional knowledge distillation (Hinton et al., 2015). (Bhat et al., 2021; Gani et al., 2022) has shown that self-distillation can be effective in improving performance on small datasets, and may also act as a regularization-like effect to aid in model training(He et al., 2022). Recently, (Shen et al., 2022) proposed a new method of implementing self-distillation, named DLB, where the model from the previous iteration is used as the teacher for the current iteration. The loss function for this method is formulated as:

$$\mathcal{L}_{LB} = \frac{1}{n} \sum_{i=1}^{n} \mathcal{T}^2 \cdot D_{KL} \left(p_i^{\mathcal{T},t-1} \| p_i^{\mathcal{T},t} \right)$$
(3)

Where n is the number of samples in a batch, \mathcal{T} is the temperature of the distillation, and p_i is the predicted distribution of each sample (Shen et al., 2022).

To date, there is a lack of literature on the application of self-distillation to the task of medical image segmentation. We propose to integrate the DLB loss as an additional term to the primary loss function in our model architecture, which is a combination of weighted cross-entropy loss and Dice loss, commonly used in medical image segmentation tasks. The overall loss function for our model is presented in Eq. (4).

$$\mathcal{L}_{MAT} = 0.4 \cdot \mathcal{L}_{CE} + 0.6 \cdot \mathcal{L}_{\text{Dice}} + \alpha \cdot \mathcal{L}_{LB} \tag{4}$$

The implementation of our DLB follows the approach of the original DLB, with a constraint of maintaining half of the mini-batch consistency between consecutive iterations. However, due to the Transformer models tending to converge more slowly and being more

WITHHELD

susceptible to instability during the early stages of training, the use of DLB may instead lead to the problem of each iteration affecting each other and eventually being difficult to converge. To mitigate these issues, we implement a dynamic warm-up schedule for the distillation coefficient α and temperature \mathcal{T} . Specifically, the value of α is set to 0 for the first 50 training epochs, before being set to 1. In the following epoch, \mathcal{T} is linearly increased from 1 to 2, with the aim of placing more emphasis on the distribution of negative samples during later stages of training.

3. Experiments

3.1. Setup

3.1.1. Dataset

The BraTS datasets from the MICCAI Brain Tumor Segmentation Challenge contain multimodal 3D brain MRI scans annotated with ground truth segmentations of tumor regions by physicians. The datasets consist of four MRI modalities per case (T1, T1ce, T2, and FLAIR), with annotations of four tumor subregions consolidated into three sub-regions: whole tumor (WT), tumor core (TC), and enhancing tumor (ET). BraTS2018 (Menze et al., 2014a; Bakas et al., 2017c, 2018) was collected from 19 institutions and includes both low-grade and high-grade gliomas. The study focuses on high-grade gliomas, as it provides an opportunity to demonstrate the advantages of the model on small datasets. The HGG group contains 210 samples randomly split into 180 training cases and 30 testing cases. In comparison, BraTS2021 (Bakas et al., 2017c; Baid et al., 2021; Menze et al., 2014b; Bakas et al., 2017a,b) is larger and randomly split into 1200 training cases and 200 testing cases to enable a comparison of our method's performance with larger datasets.

3.1.2. Evaluation Metrics

As in previous studies (Valanarasu et al., 2021; Yan et al., 2022), we also used the Dice score and the 95% Hausdorff Distance to assess the overall accuracy of segmentation as well as the surface accuracy. The formulas of Dice and HD95 are defined in appendix section .

3.1.3. Implementation details

The images in the BraTS 2018 and 2021 datasets were re-sized to $(160 \times 224 \times 224)$ and $(128 \times 160 \times 160)$ for consistency and to fill empty slices. Data augmentation, including random flips and rotations, was applied with a 50% probability to enhance the model's fitting ability. The hyperparameters of the model were tuned via 5-fold cross-validation on the training set, and then applied to train the full training set, yielding the best model. The AdamW optimizer (Loshchilov and Hutter, 2018) with a weight decay of 10-5 was utilized with a warm-up schedule for the learning rate (Gotmare et al., 2018), where the learning rate grows linearly from 0 to 10-3 for the first 10 epochs, followed by cosine anneal (Loshchilov and Hutter, 2016) to complete the learning rate decay. In the Transformer encoder, we employed 1, 2, and 4 blocks for stage S1, S2, and S3, respectively. For the sake of fairness, in our experiments, all models were trained for an equal number of 250 epochs.

Matuian	D	WT		TC		ET		Mean	
Metrics	Params	Dice	HD95	Dice	HD95	Dice	HD95	Dice	HD95
nnUnet(3D)	$17.8 \mathrm{M}$	85.61	4.33	78.72	6.59	70.23	4.91	78.19	5.28
Trans-Unet(2D)	$96.1 \mathrm{M}$	84.42	4.91	75.12	7.18	73.97	5.07	77.84	5.72
Swin-Unet(2D)	$79.6 \mathrm{M}$	90.64	6.01	80.78	7.07	75.54	5.98	82.32	6.35
Swin-UNETR(3D)	$62.1 \mathrm{M}$	88.36	4.32	86.89	6.51	80.21	4.28	85.16	5.04
AFTer-Unet(3D)	41.5M	91.93	4.15	87.15	6.76	81.58	3.91	86.89	4.94
MAT(3D)	$\underline{11.7M}$	93.05	4.06	87.91	<u>6.09</u>	82.81	4.02	87.92	4.72

3.2. Results on BraTS2018

Table 1: Dice scores and HD95 of different methods on the BraTS2018 dataset (testing).

We compared our proposed 3D Medical Attention Transformer (MAT) model to other medical image segmentation models on the BraTS2018 dataset in Table 1. We selected several established models such as nnUnet (Isensee et al., 2021), Trans-Unet (Chen et al., 2021), Swin-Unet (Cao et al., 2021), Swin-UNETR (Hatamizadeh et al., 2022a), AFTer-Unet (Yan et al., 2022) for fair comparison. These models have shown state-of-the-art results in various medical image segmentation tasks and are widely used in the field.

Table 1 shows that all models have good segmentation performance in the Whole Tumor (WT) region, likely due to clear distinction between tumor and non-tumor regions and larger relative segmentation volume. For the Enhancing Tumor (ET) region, the most challenging to segment, 3D models incorporating the Transformer module achieved high Dice scores. The inclusion of the depth axis and long-range dependencies improves the model's understanding of the overall image, leading to better segmentation. 3D models outperformed 2D models in HD95 metrics due to problematic jagged edges when combining 2D slices into 3D images. MAT leverages 3D convolution to extract features and inputs feature maps with richer semantic information into the Axial Transformer module, achieving true global long-range dependency modeling, which resulting in superior performance. Overall, MAT achieved a mean Dice of 87.92% and a mean HD95 of 4.72 on the BraTS2018, surpassing AFTer-Unet in nearly all metrics.

In addition, the comparison of parameters demonstrates the superiority of our 3D axial attention algorithm as it reduces model complexity while preserving high segmentation performance. The MAT model also has a lower computational cost and can be trained on a single Tesla T4 (16GB) GPU, compared to other Transformer-based models.

3.3. Results on BraTS2021

Table 2 compares the performance of MAT with other medical image segmentation models on the BraTS2021 dataset. The comparison includes nnUnet (Isensee et al., 2021), Trans-Unet (Chen et al., 2021), Swin-Unet (Cao et al., 2021), Swin-UNETR (Hatamizadeh et al., 2022a), and AFTer-Unet (Yan et al., 2022). It is important to note that BraTS2021 contains more MRI data than BraTS2018, which may pose challenges for the lightweight MAT model.

Despite this, the results show that MAT's performance is comparable to other larger models. Notably, MAT performed significantly better than the other models in the TC (tumor core) region, with a Dice score improvement of 1.27% and a reduction of 0.20

WITHHELD

Motrics	W	Т	TC		ET		Mean	
Metrics	Dice	HD95	Dice	HD95	Dice	HD95	Dice	HD95
nnUnet(3D)	92.14	7.33	89.56	3.94	83.67	4.02	88.46	5.10
Trans-Unet(2D)	91.73	7.92	85.47	6.02	81.25	4.68	86.15	6.21
Swin-Unet(2D)	93.51	7.51	90.64	5.61	85.34	4.18	89.83	5.77
$\frac{\text{Swin-UNETR}(3D)}{\text{Swin-UNETR}(3D)}$	94.74	7.02	89.91	3.75	85.41	3.41	90.02	4.73
AFTer-Unet(3D)	93.47	6.95	90.48	3.76	85.31	3.29	89.75	4.67
MAT(3D)	93.21	7.13	<u>91.91</u>	3.56	85.05	3.61	<u>90.06</u>	4.77

Table 2: Dice scores and HD95 of different methods on the BraTS2021 dataset (testing).

in HD95 score. This highlights MAT's strong ability to learn with a limited number of parameters. Furthermore, on average, MAT is the best model under the Dice metric, and only slightly below the AFTer-Unet at 0.10 according to HD95.

3.4. Ablation study of warm-upped DLB

MAT with original DLB

MAT with warm-upped DLB

	BraTS	5 2018	BraTS 2021	
Mean Dice	Training	Testing	Training	Testing

 $79.22_{\pm 2.64}$

 $87.92_{\pm 2.01}$

 82.76 ± 1.49

90.02+0.98

 $82.92_{\pm 1.48}$

90.06±1.47

 77.13 ± 2.49

 87.65 ± 1.85

Table 3: Ablation study of DLB with warm-up schedule met hod on BraTS2018 and BraTS2021 Datasets through dice score.

Table 3 presents the results of ablation experiments to evaluate the efficacy of the proposed warm-upped DLB method on the BraTS2018 and BraTS2021 datasets. The parameters \mathcal{T} and α were set to 3 and 1, respectively, as per the sets in (Shen et al., 2022). The results show improved performance of the MAT model with the warm-upped DLB method, with increases of 1.68% and 1.89% for the training and validation sets on the BraTS2018 dataset. The absence of warm-up led to decreased performance. The warm-upped DLB method was still effective on the BraTS2021 dataset, but with smaller improvement, indicating its greater efficacy for smaller datasets.

4. Conclusion

We propose 3D Medical Axial Transformer(MAT), a 3D end-to-end framework for brain tumor image segmentation that utilizes the axial Transformer and self-distillation scheme. The design of MAT enables efficient learning of semantic information while maintaining a lightweight architecture, making it suitable for clinical applications. Our experiments on brain tumor datasets demonstrate the superiority of MAT over previous related methods.

References

- Ujjwal Baid, Satyam Ghodasara, Suyash Mohan, Michel Bilello, Evan Calabrese, Errol Colak, Keyvan Farahani, Jayashree Kalpathy-Cramer, Felipe C Kitamura, Sarthak Pati, et al. The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification. arXiv preprint arXiv:2107.02314, 2021.
- S. Bakas, H. Akbari, A. Sotiras, M. Bilello, M. Rozycki, and J. Kirby. Segmentation labels and radiomic features for the pre-operative scans of the tcga-gbm collection. *The Cancer Imaging Archive*, 2017a. doi: 0.7937/K9/TCIA.2017.KLXWJJ1Q.
- S. Bakas, H. Akbari, A. Sotiras, M. Bilello, M. Rozycki, and J. Kirby. Segmentation labels and radiomic features for the pre-operative scans of the tcga-lgg collection. *The Cancer Imaging Archive*, 2017b. doi: 10.7937/K9/TCIA.2017.GJQ7R0EF.
- Spyridon Bakas, Hamed Akbari, Aristeidis Sotiras, Michel Bilello, Martin Rozycki, Justin S Kirby, John B Freymann, Keyvan Farahani, and Christos Davatzikos. Advancing the cancer genome atlas glioma mri collections with expert segmentation labels and radiomic features. *Scientific data*, 4(1):1–13, 2017c.
- Spyridon Bakas, Mauricio Reyes, Andras Jakab, Stefan Bauer, Markus Rempfler, Alessandro Crimi, Russell Takeshi Shinohara, Christoph Berger, Sung Min Ha, Martin Rozycki, et al. Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the brats challenge. *arXiv preprint arXiv:1811.02629*, 2018.
- Prashant Bhat, Elahe Arani, and Bahram Zonooz. Distill on the go: Online knowledge distillation in self-supervised learning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 2678–2687, 2021.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.
- Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning Wang. Swin-unet: Unet-like pure transformer for medical image segmentation. arXiv preprint arXiv:2105.05537, 2021.
- Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L Yuille, and Yuyin Zhou. Transunet: Transformers make strong encoders for medical image segmentation. arXiv preprint arXiv:2102.04306, 2021.
- Ozgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger. 3d u-net: learning dense volumetric segmentation from sparse annotation. In *International conference on medical image computing and computer-assisted intervention*, pages 424–432. Springer, 2016.

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pretraining of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
- Kunio Doi. Computer-aided diagnosis in medical imaging: historical review, current status and future potential. *Computerized medical imaging and graphics*, 31(4-5):198–211, 2007.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
- Hanan Gani, Muzammal Naseer, and Mohammad Yaqub. How to train vision transformer on small-scale datasets? arXiv preprint arXiv:2210.07240, 2022.
- Akhilesh Gotmare, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. A closer look at deep learning heuristics: Learning rate restarts, warmup and distillation. arXiv preprint arXiv:1810.13243, 2018.
- Meng-Hao Guo, Cheng-Ze Lu, Qibin Hou, Zhengning Liu, Ming-Ming Cheng, and Shi-Min Hu. Segnext: Rethinking convolutional attention design for semantic segmentation. arXiv preprint arXiv:2209.08575, 2022.
- Ali Hatamizadeh, Vishwesh Nath, Yucheng Tang, Dong Yang, Holger R Roth, and Daguang Xu. Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images. pages 272–284, 2022a.
- Ali Hatamizadeh, Yucheng Tang, Vishwesh Nath, Dong Yang, Andriy Myronenko, Bennett Landman, Holger R Roth, and Daguang Xu. Unetr: Transformers for 3d medical image segmentation. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pages 574–584, 2022b.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016.
- Yangji He, Weihan Liang, Dongyang Zhao, Hong-Yu Zhou, Weifeng Ge, Yizhou Yu, and Wenqiang Zhang. Attribute surrogates learning and spectral tokens pooling in transformers for few-shot learning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 9119–9129, 2022.
- Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531, 2(7), 2015.
- Jonathan Ho, Nal Kalchbrenner, Dirk Weissenborn, and Tim Salimans. Axial attention in multidimensional transformers. arXiv preprint arXiv:1912.12180, 2019.
- Fabian Isensee, Paul F Jaeger, Simon AA Kohl, Jens Petersen, and Klaus H Maier-Hein. nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods*, 18(2):203–211, 2021.

- Xiaomeng Li, Hao Chen, Xiaojuan Qi, Qi Dou, Chi-Wing Fu, and Pheng-Ann Heng. Hdenseunet: hybrid densely connected unet for liver and tumor segmentation from ct volumes. *IEEE transactions on medical imaging*, 37(12):2663–2674, 2018.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 10012–10022, 2021.
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. arXiv preprint arXiv:1608.03983, 2016.
- Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. 2018.
- Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree Kalpathy-Cramer, Keyvan Farahani, Justin Kirby, Yuliya Burren, Nicole Porz, Johannes Slotboom, Roland Wiest, et al. The multimodal brain tumor image segmentation benchmark (brats). *IEEE transactions* on medical imaging, 34(10):1993–2024, 2014a.
- Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree Kalpathy-Cramer, Keyvan Farahani, Justin Kirby, Yuliya Burren, Nicole Porz, Johannes Slotboom, Roland Wiest, et al. The multimodal brain tumor image segmentation benchmark (brats). *IEEE transactions* on medical imaging, 34(10):1993–2024, 2014b.
- Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In 2016 fourth international conference on 3D vision (3DV), pages 565–571. IEEE, 2016.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing* and computer-assisted intervention, pages 234–241. Springer, 2015.
- Yiqing Shen, Liwu Xu, Yuzhe Yang, Yaqian Li, and Yandong Guo. Self-distillation from the last mini-batch for consistency regularization. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 11943–11952, 2022.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556, 2014.
- Jeya Maria Jose Valanarasu, Poojan Oza, Ilker Hacihaliloglu, and Vishal M Patel. Medical transformer: Gated axial-attention for medical image segmentation. In International Conference on Medical Image Computing and Computer-Assisted Intervention, pages 36– 46. Springer, 2021.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.

- Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Axial-deeplab: Stand-alone axial-attention for panoptic segmentation. In European Conference on Computer Vision, pages 108–126. Springer, 2020.
- Yutong Xie, Jianpeng Zhang, Chunhua Shen, and Yong Xia. Cotr: Efficiently bridging cnn and transformer for 3d medical image segmentation. In *International conference on medical image computing and computer-assisted intervention*, pages 171–180. Springer, 2021.
- Xiangyi Yan, Hao Tang, Shanlin Sun, Haoyu Ma, Deying Kong, and Xiaohui Xie. Afterunet: Axial fusion transformer unet for medical image segmentation. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pages 3971–3981, 2022.
- Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3713–3722, 2019.
- Zhengxin Zhang, Qingjie Liu, and Yunhong Wang. Road extraction by deep residual u-net. IEEE Geoscience and Remote Sensing Letters, 15(5):749–753, 2018.



Appendix A. Qualitative results on BraTS2018

Figure 2: Qualitative results of different methods on the BraTS2018 dataset. The first row is sagittal, the second row is axial, and the third row is the coronal view. The red bounding boxes serve to annotate regions of challenging segmentation, facilitating qualitative comparison of model results.

Figure 2 presents qualitative results of different models on the BraTS2018 dataset, compared withnnUnet (Isensee et al., 2021), Trans-Unet (Chen et al., 2021), Swin-Unet (Cao et al., 2021), Swin-UNETR (Hatamizadeh et al., 2022a), and AFTer-Unet (Yan et al., 2022).

We presented slices of images from three views (sagittal, axial, and coronal) to facilitate the comparison of 3D segmentation results. The most challenging regions to segment in brain tumor images are TC (yellow part + red part) and ET (red part) due to their small distinction from other tumors and small volume. The design of MAT's 3D Axial Transformer module with DLB enables it to learn global information, resulting in excellent segmentation ability for ET as seen in sagittal and coronal views. Additionally, as shown in the axial view (second row) and sagittal view (first row), MAT's predictions of gaps in tumors are more accurate due to its pixel-level long-distance dependence in three dimensions, in contrast to other models. Overall, MAT demonstrates a superior ability for multi-class tumor segmentation compared to other models, for both overall tumor region and local gaps.

Appendix B. Additional Results on BraTS

Table 4: Mean and std of the Dice scores for various methods on the BraTS2021(testing).

Dice	WT	TC	\mathbf{ET}	Mean
nnUnet(3D)	$85.61_{\pm 0.71}$	$78.72_{\pm 1.98}$	$70.23_{\pm 3.37}$	$78.19_{\pm 1.97}$
Trans-Unet(2D)	$84.42_{\pm 1.32}$	$75.12_{\pm 2.86}$	$73.97_{\pm 2.99}$	$77.84_{\pm 2.23}$
Swin-Unet(2D)	$90.64_{\pm 1.07}$	$80.78_{\pm 3.12}$	$75.54_{\pm 3.05}$	$82.32_{\pm 2.28}$
Swin-UNETR(3D)	88.36 ± 0.89	86.89 ± 2.75	80.21 ± 2.84	85.16 ± 1.92
AFTer-Unet(3D)	$91.93_{\pm 1.23}$	87.15 ± 3.09	81.58 ± 2.10	86.89 ± 2.28
MAT(3D)	$93.05_{\pm 1.49}$	$87.91_{\pm 2.47}$	$\underline{\textbf{82.81}_{\pm 2.41}}$	$87.92_{\pm 2.01}$

Table 5: Mean and std of the Dice scores for various methods on the BraTS2021(testing).

HD95	WT	TC	\mathbf{ET}	Mean
nnUnet(3D)	$4.33_{\pm 0.85}$	$6.59_{\pm 1.73}$	$4.91_{\pm 1.05}$	$5.28_{\pm 1.07}$
Trans-Unet(2D)	$4.91_{\pm 1.08}$	$7.18_{\pm 2.47}$	$5.07_{\pm 1.98}$	$5.72_{\pm 1.65}$
Swin-Unet(2D)	$6.01_{\pm 0.68}$	$7.07_{\pm 2.56}$	5.98 ± 2.48	$6.35_{\pm 1.63}$
Swin-UNETR(3D)	$4.32_{\pm 0.71}$	$6.51_{\pm 1.91}$	$4.28_{\pm 1.87}$	$5.04_{\pm 1.40}$
AFTer-Unet(3D)	$4.15_{\pm 0.88}$	$6.76_{\pm 2.39}$	$3.91_{\pm 2.04}$	$4.94_{\pm 1.61}$
<u>MAT(3D)</u>	$\underline{\textbf{4.06}_{\pm 1.01}}$	$\underline{\textbf{6.09}_{\pm 2.16}}$	$\underline{4.02_{\pm 2.10}}$	$\underline{4.72_{\pm 1.64}}$

Table 6: Mean and std of the Dice scores for various methods on the BraTS2021(testing).

Dice	WT	TC	ET	Mean
nnUnet(3D)	$92.14_{\pm 0.69}$	$89.56_{\pm 1.83}$	$83.67_{\pm 2.39}$	$88.46_{\pm 1.41}$
Trans-Unet(2D)	91.73 ± 1.01	85.47 ± 1.37	81.25 ± 2.19	86.15 ± 1.40
Swin-Unet(2D)	$93.51_{\pm 0.87}$	$90.64_{\pm 1.97}$	$85.34_{\pm 2.12}$	$89.83_{\pm 1.53}$
Swin-UNETR(3D)	$94.74_{\pm 0.78}$	$89.91_{\pm 1.81}$	$85.41_{\pm 1.97}$	$90.02_{\pm 1.39}$
$\operatorname{AFTer-Unet}(3D)$	$93.47_{\pm 0.93}$	$90.48_{\pm 2.01}$	$85.31_{\pm 2.14}$	$89.75_{\pm 1.58}$
<u>MAT(3D)</u>	$\underline{93.21_{\pm 0.93}}$	$\underline{91.91_{\pm 1.92}}$	$\underline{85.05_{\pm 1.98}}$	$\underline{\textbf{90.06}_{\pm 1.47}}$

Table 7: Mean and std of the Dice scores for various methods on the BraTS2021(testing)

HD95	WT	TC	\mathbf{ET}	Mean
nnUnet(3D)	$7.33_{\pm 1.37}$	$3.94_{\pm 1.08}$	$4.02_{\pm 1.32}$	$5.10_{\pm 1.11}$
Trans-Unet(2D)	$7.92_{\pm 1.22}$	$6.02_{\pm 1.30}$	$4.68_{\pm 1.47}$	$6.21_{\pm 1.20}$
Swin-Unet(2D)	$7.51_{\pm 1.02}$	$5.61_{\pm 1.19}$	$4.18_{\pm 1.28}$	$5.77_{\pm 1.10}$
Swin-UNETR(3D)	$7.02_{\pm 0.96}$	$3.75_{\pm 0.99}$	$3.41_{\pm 1.13}$	$4.73_{\pm 0.92}$
AFTer-Unet(3D)	$6.95_{\pm 1.00}$	3.76 ± 1.34	$3.29_{\pm 0.98}$	$4.67_{\pm 1.02}$
MAT(3D)	$7.13_{\pm 1.10}$	$3.56_{\pm 1.38}$	$3.61_{\pm 1.29}$	$4.77_{\pm 1.08}$

Appendix C. Formulas of evaluation metrics

$$Dice(T, P) = \frac{2\sum_{i=1}^{I} T_i P_i}{\sum_{i=1}^{I} T_i + \sum_{i=1}^{I} P_i}$$
(5)

$$HD(T',P') = \max\left\{\max_{t'\in T'}\min_{p'\in P'} \|t'-p'\|, \max_{p'\in P'}\min_{t'\in T'} \|p'-t'\|\right\}$$
(6)

Where T_i and P_i denote the ground truth and predicted values of voxel *i*, while T' and P' denote the set of surface points of the ground truth and predicted values, respectively. HD95 is based on the calculation of the 95th percentile of the distances between boundary points in T' and P', in order to eliminate the effect of a very small subset of the outliers.

Appendix D. The output of axial attention for depth axis and width axis

$$y_{\text{depth}_{ijk}} = \sum_{d=1}^{D} \text{Softmax} \left(q_{ijk}^{T} k_{djk} + q_{(i,j,k)}^{T} r_{djk}^{q} + k_{djk}^{T} r_{djk}^{k} \right) \left(b_{djk} + r_{djk}^{v} \right)$$
(7)

$$y_{\text{width }_{ijk}} = \sum_{w=1}^{W} \text{Softmax} \left(q_{ijk}^T k_{ijw} + q_{(i,j,k)}^T r_{ijw}^q + k_{ijw}^T r_{ijw}^k \right) \left(v_{ijw} + r_{ijw}^v \right)$$
(8)