# Swin UNETR++: An Efficient Framework for Brain Tumor Segmentation

**Xiwei Deng(36920231153187)[1*], Xinyi He(36920231153192)[1],**
**Yu Meng(36920231153223)[1], Xun Guan(33320231150348)[2], Mengying Zhu(36920231153271)[1]**

[1]Institute of Artificial Intelligence, Xiamen University, 361005, China
[2]School of Electronic Science and Engineering,Xiamen University, 361005, China
{xiweideng,xinyihe,yumeng,xunguan,zhumengying}@stu.xmu.edu.cn

## Abstract

Brain tumor segmentation aims to isolate and delineate glioma tissues from healthy brain tissues within magnetic resonance imaging, a critical step for effective diagnosis and therapeutic strategy formulation in neurooncology. However, MRI-based tumor segmentation is challenging due to the complex morphology and unclear boundaries of the tumors. To address such limitations, we propose Swin UNETR++ for fully automated tumor segmentation. Our approach utilizes the Swin Transformer structure and a multi-scale feature fusion strategy, enabling the network to capture more contextual information and high-level semantic details, thus improving segmentation accuracy and robustness. The proposed method outperforms the current SOTA model in the BraTS 2021 datasets, with segmentation accuracies of 92.7% for the whole tumor, 91.2% for the tumor core, and 87.6% for the enhanced tumor, which improve upon the TransBTS, by 1.6%, 1.4%, and 0.8% respectively.

## Introduction

Gliomas are the most common and deadly type of brain tumor(Pereira et al. 2016), accounting for approximately 80% malignant brain tumors(Mlynarski et al. 2019). They can be classified into low-grade gliomas (LGG) and high-grade gliomas (HGG), with HGG being more aggressive and invasive. Glioblastoma multiforme (GBM)is the most malignant form of glioma among astrocytic tumors and can be fatal if not detected early.

Over the last two decades, magnetic resonance imaging (MRI) has become a popular tool for the inspection of different brain disorders due to its ability to provide internal observations of brain tissues with high spatial and temporal resolutions. In clinical practice, radiologists rely on various MRI sequences to make a comprehensive diagnosis of gliomas. These sequences typically include T1, T2, FLAIR, and T1ce. Each sequence provides different information about the tumor, allowing for a more comprehensive analysis of different subregions within the brain tumor.

Segmentation of glioma lesions is a crucial step in computer-aided diagnosis, surgery, radiotherapy, and chemotherapy planning for gliomas. However, segmenting

---

*Corresponding author

brain tumors in high-precision magnetic resonance images is challenging due to their complex and variable shapes, unclear boundaries, and low contrast. Currently, manual segmentation by radiologists is the main approach, but it is time-consuming and lacks reproducibility. Therefore, many efforts have been made to segment semi- or fully automatic glioma to improve both efficiency and potential accuracy(Işın, Direkoğlu, and Şah 2016) Earlier efforts included intensity thresholding, edge, and region-based methods.(Hiralal and Menon 2016) The intensity thresholding method categorizes pixels based on intensity ranges(Sujji, Lakshmi, and Jiji 2013); the edge-based approach classifies pixels as edged or non-edged using filters(Soltanian-Zadeh and Windham 1997); the region-based approach groups neighboring pixels with high similarity while dividing pixels with significant dissimilarity.(Ilunga-Mbuyamba et al. 2017)

However, these methods face challenges when applied in scenarios that require fully automatic processing, as they often require setting of initial seed points, thresholds, and iteration termination conditions. On the other hand, automatic segmentation algorithms, which require no human interaction, offer high segmentation speed and reproducible results, facilitating the development of end-to-end applications for glioma. Automatic segmentation is a major research direction in glioma segmentation, with improving segmentation accuracy being a key challenge. Several algorithms based on machine learning and deep learning have shown promising results in brain tumor segmentation tasks. The latest research focuses on the development of deep learning segmentation algorithms, particularly a novel TransBTS brain tumor segmentation method based on the encoder-decoder structure. Furthermore, researchers have conducted numerous multi-parameter MRI quantifications to mitigate the uncertainties often encountered in deep learning approaches(Yang et al. 2023). These advances aim to improve the reliability and performance of brain tumor segmentation techniques in medical imaging analysis. These methods not only perform well, but also have the ability to continuously learn and update themselves in a data-driven manner, enabling fully automatic segmentation that can adapt to the complex needs of diverse clinical medical applications.

The goal of the research in this article is to design an automatic brain tumor segmentation system based on Convolutional Neural Networks (CNNs). This system focuses on
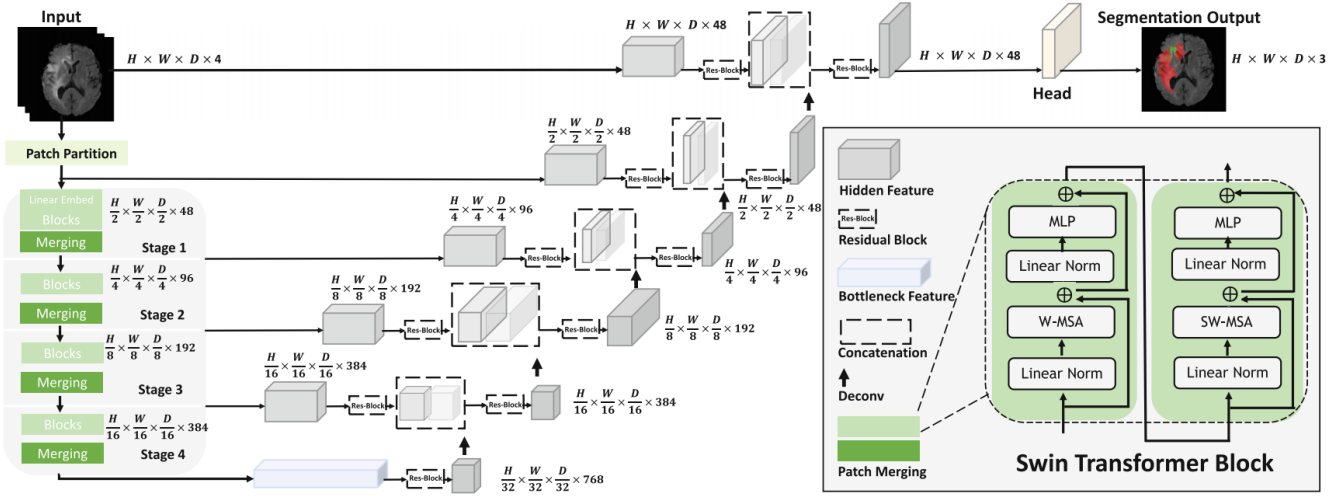
Figure 1: The overview of our proposed Swin UNETR++

multi-modal MRI image segmentation and aims to achieve comprehensive brain tumor segmentation. By providing automatic segmentation, this system can assist clinicians in making treatment decisions, better serve patients and doctors, and make significant contributions to artificial intelligence in medical and healthcare in the country.

## Related work

With the advancement of medical imaging technology, there has been a growing interest in developing algorithms for brain tumor segmentation in medical image analysis. Segmentation of brain tumors in MRI has become a major topic in the medical imaging field, aiding in diagnosis and treatment. However, low efficiency, low accuracy, and poor robustness present challenges. As a result, numerous classic medical image segmentation methods have been developed domestically and internationally to address these issues(Zhu and Shen 2019; Ramesh et al. 2021).

### Traditional segmentation methods

Ahilan A proposed a multi-threshold segmentation method based on optimization techniques to extract regions of interest and compress DICOM images using an improved lossless prediction algorithm(Ahilan et al. 2019). A S et al. proposed an improved Sobel edge detection algorithm with eight directions, which improved the edge detection performance of brain tumor MRI images(AS and Gopalan 2022). Khosravanian et al. proposed a Fuzzy Kernel Level Set (FKLS) algorithm based on fuzzy c-means, kernel mapping, and symmetry analysis for brain tumor image segmentation(Khosravanian et al. 2022).The previously mentioned methods, which are based on thresholding, region-based, graph theory, edge detection, active contours, and model-based medical image segmentation, were proposed relatively early. They have demonstrated advantages in small-scale private and older publicly available datasets but may not effectively meet the increasing demands for the segmentation of the latest clinical MRI brain tumor data.

### Traditional machine learning methods

Saxena et al. introduced a brain tumor MRI image segmentation method that employs a sliding window mechanism and fuzzy c-means clustering(Saxena, Kumari, and Pattnaik 2021). Initially, the method preprocesses brain tumor MR images, completes texture feature extraction and classification, and then uses a sliding window mechanism to localize tumor regions. Finally, the fuzzy clustering algorithm of c-means is applied to remove erroneously classified windows, resulting in brain tumor segmentation. Experimental results indicate that this method exhibits comparable or superior accuracy and Dice scores in brain tumor segmentation compared to other existing methods, including deep learning-based approaches.

### Deep learning methods

In recent years, deep learning has shown remarkable potential in computer vision and has made significant advancements in brain tumor segmentation. Wang et al. proposed TransBTS, a novel brain tumor segmentation method based on an encoder-decoder structure(Wang et al. 2021). The encoder utilizes 3D CNN to extract volumetric spatial features, while the decoder employs the Transformer model with embedded features for progressive upsampling and generating tumor segmentation results. Liu et al. developed Swin Transformer, a layered Transformer based on the Shifted Window computation representation. Luu et al. introduced an extended nn-Unet-based brain tumor MRI image segmentation method(Luu and Park 2021), which incorporates group normalization and a larger Unet network, achieving first place in the BraTS 2021 Brain Tumor Segmentation Challenge Task 1: Brain Tumor Segmentation in mpMRI scans. Deep learning-based methods have consistently emerged as winners of the Brain Tumor Segmentation Challenge Task 1 at prestigious MICCAI conferences in recent years(Jiang et al. 2020; Isensee et al. 2021; Hatamizadeh et al. 2021).

## Proposed Solution

In this section, we will introduce Swin UNETR ++, as shown in Figure 1.Our model is optimized based on Swin UETR(Hatamizadeh et al. 2021). Swin UNETR++ consists of encoder, decoder and skip connections. The basic unit of Swin-Unet is the Swin Transformer block(Liu et al. 2021). The input to our model is 3D multi-modal MRI images with 4 channels. The Swin UNETR++ creates non-overlapping patches of the input data and uses a patch partition layer to create windows with a desired size for computing the self-attention. The encoded feature representations in the Swin transformer are fed to a CNN-decoder via skip connection at multiple resolutions. Final segmentation output consists of 3 output channels corresponding to ET, WT and TC subregions.

### Encoder

For the encoder, to transform the inputs into sequence embeddings, the medical images are split into nonoverlapping patches with patch size of $2 \times 2 \times 2$. By this partition approach, the feature dimension of each patch becomes to $2 \times 2 \times 2 \times 4 = 32$. Furthermore, a linear embedding layer is applied to projected feature dimension into arbitrary dimension (represented as C). The transformed patch tokens pass through several Swin Transformer blocks and patch merging layers to generate the hierarchical feature representations. Specifically, patch merging layer is responsible for downsampling and increasing dimension, and Swin Transformer block is responsible for feature representation learning.

The self-attention is computed into nonoverlapping windows that are created in the partitioning stage for efficient token interaction modeling. Figure 2 shows the shifted windowing mechanism for subsequent layers. Subsequently, in subsequent layers of l and l + 1 in the encoder, the outputs are calculated as

$$\hat{z}^l = \text{W-MSA}(\text{LN}(z^{l-1})) + z^{l-1} \tag{1}$$

$$z^l = \text{MLP}(\text{LN}(\hat{z}^l)) + \hat{z}^l \tag{2}$$

$$\hat{z}^{l+1} = \text{SW-MSA}(\text{LN}(z^l)) + z^l \tag{3}$$

$$z^{l+1} = \text{MLP}(\text{LN}(\hat{z}^{l+1})) + \hat{z}^{l+1} \tag{4}$$

Here, W-MSA and SW-MSA are regular and window partitioning multi-head self-attention modules respectively $\hat{z}^l$ and $\hat{z}^{l+1}$ denote the outputs of W-MSA and SW-MSA; MLP and LN denote layer normalization and Multi-Layer Perceptron respectively. For efficient computation of the shifted window mechanism, we leverage a 3D cyclic-shifting (Liu et al. 2021) and compute self-attention according to

$$\text{Attention}(Q, K, V) = Softmax(\frac{QK^T}{\sqrt{d}})V \tag{5}$$

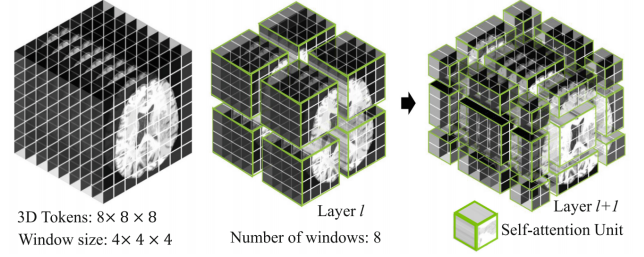In which Q, K, V denote queries, keys, and values respectively; d represents the size of the query and key.



Figure 2: Overview of the shifted windowing mechanism. Note that $8 \times 8 \times 8$ 3D tokens and $4 \times 4 \times 4$ window size are illustrated.

### Decoder

Swin UNETR++ has a U-shaped network design in which the extracted feature representations of the encoder are used in the decoder via skip connections at each resolution. At each stage i ($i \in \{0, 1, 2, 3, 4\}$) in the encoder and the bottleneck ($i = 5$), the output feature representations are reshaped into size $\frac{H}{2^i} \times \frac{W}{2^i} \times \frac{D}{2^i}$ and fed into a residual block comprising of two $3 \times 3 \times 3$ convolutional layers that are normalized by instance normalization layers. Subsequently, the resolution of the feature maps are increased by a factor of 2 using a deconvolutional layer and the outputs are concatenated with the outputs of the previous stage. The concatenated features are then fed into another residual block as previously described. The final segmentation outputs are computed by using a $1 \times 1 \times 1$ convolutional layer and a sigmoid activation function.

## Experiments

### Dataset

As with most related work, we choose the BraTS 2021 dataset which comprises four modalities of MRI images - T1, T2, T1ce and FLAIR from 2000 patients. The image size is 240×240×155 which has been resampled to isotropic 1×1×1mm resolution, stemming from multiple clinical institutions with different imaging equipment.

As shown in Figure 3, BraTS 2021 contains 4 types of segmentation labels: the blue area represents enhanced tumor (ET), the green area represents edema tissue around the tumor (ED), the red area represents necrotic tumor core (NCR), and the black area represents background. These four types of label combinations are divided into three types of subregions: the whole tumor (WT) contains ET, ED, and NCR, the tumor core (TC) contains ET and NCR, and the enhanced tumor (ET) area corresponds to the cyan, magenta, and yellow areas, respectively.

### Implementation Details

We implement our method using Pytorch 1.10, MONAI 0.9.0 and Medpy 0.4.0, which is evaluated using the 5-fold cross-validation. To save computing resources, we opt to fine-tune the swin unetr pre-trained models for 30 epochs using a single NVIDIA A100 40GB PCIe GPU with input 4-channel image size 128×128×128x4 for 472 iterations. Fur-
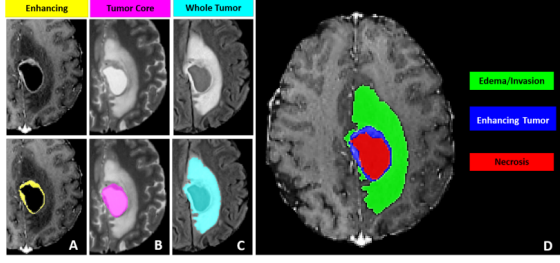
Figure 3: Glioma subregions with its corresponding labels

thermore, we utilize AdamW optimizer with learning rate of $1e^{-4}$ and weight decay of $1e^{-5}$, cosine annealing restart with warm up to adjust learning rate and sigmoid activation function.

## Evaluations Metrics

In order to evaluate the performance of the models, we adopt Dice Similarity Coefficient(DSC),Jaccard Similarity Coefficient(JSC),Sensitivity,Specificity,and Positive Predictive Value(PPV) as follows.

DSC measures the spatial overlap between voxels of predicted segmentations and volumetric ground truths. Similar to DSC, JSC also reflects the spatial similarity of two volumes, yet is defined from a different perspective.

$$\text{DSC} = \frac{2|Y \cap \hat{Y}|}{|Y| + |\hat{Y}|} \quad (6)$$

$$\text{JSC} = \frac{|Y \cap \hat{Y}|}{|Y \cup \hat{Y}|} \quad (7)$$

where, $Y$ and $\hat{Y}$ denote the ground truths and predicted segmentations, respectively.

Sensitivity represents the proportion of predicted positive tumor labels to true tumor labels.

$$\text{Sensitivity} = \frac{TP}{P} = \frac{TP}{TP + FN} \quad (8)$$

where, TP, P, and FN represent the number of true positive examples, the number of correctly predicted positive examples, and the number of incorrectly predicted negative examples,respectively.

Specificity indicates the ratio of predicted positive background labels to true background labels.

$$\text{Specificity} = \frac{TN}{N} = \frac{TN}{TN + FP} \quad (9)$$

where, TN, N and FP stand for the number of correctly predicted negative examples, the number of true negative examples, and the number of incorrectly predicted positive examples, respectively.

PPV represents the proportion of true positive cases to all positive cases.

$$\text{PPV} = \frac{TP}{TP + FP} \quad (10)$$

where, TP and FP symbolize the number of correctly predicted positive examples and the number of incorrectly predicted positive examples,respectively.

## Loss Function

We leverage the soft Dice loss function calculated voxel-wise(Milletari, Navab, and Ahmadi 2016).

$$\text{Loss}(G, Y) = 1 - \frac{2}{J} \sum_{j=1}^{J} \frac{\sum_{i=1}^{I} G_{i,j} Y_{i,j}}{\sum_{i=1}^{I} G_{i,j}^2 + \sum_{i=1}^{I} Y_{i,j}^2} \quad (11)$$

where, G, Y, I, J, $Y_{i,j}$, and $G_{i,j}$ denote ground truth, predicted segmentation, the number of voxels, the number of classes, the output probability of class j on voxel i, and the true label in the form of one-hot encoding, respectively.

## Quantitative Evaluation

The experimental results of the 5-fold cross-validation are presented in Table 1.Figures 5 to 7 depict the trends of Dice, Sensitivity, and Specificity scores over 30 epochs, respectively.All evaluation metrics have been rounded to three decimal places.The Dice Similarity Coefficient (DSC) measures the similarity between the predicted and ground truth segmentations.The Jaccard Similarity Coefficient (JSC) quantifies the overlap between the predicted and ground truth segmentations. Sensitivity represents the ability to correctly identify positive cases, while specificity represents the ability to correctly identify negative cases. Positive Predictive Value (PPV) denotes the proportion of correctly predicted positive cases. The tumor core region is referred to as TC, the whole tumor region as WT, and the enhancing tumor region as ET.

| Regions Metrics | TC | WT | ET | avg |
|---|---|---|---|---|
| DSC | 0.912 | 0.927 | 0.876 | 0.905 |
| JSC | 0.838 | 0.864 | 0.78 | 0.827 |
| Sensitivity | 0.906 | 0.92 | 0.879 | 0.902 |
| Specificity | 1 | 0.999 | 1 | 1 |
| PPV | 0.943 | 0.944 | 0.902 | 0.93 |

Table 1: Experimental Results of 5-Fold Cross-Validation

To comprehensively evaluate the performance of the Swin UNETR++ model, we conduct a benchmark comparison with the TransBTS model, which has emerged as the winner of the BraTS Multimodal Brain Tumor Segmentation Challenge in recent years. The results of this comparison are presented in Table Table 2.The optimal results are highlighted in bold.Notably, the benchmark results demonstrate the superior performance of the proposed Swin UNETR++ model compared to the TransBTS model in terms of the three subregions of brain tumors. This improved performance can be attributed to the effective incorporation of multi-scale context information through the self-attention module embedded in the hierarchical encoder of Swin UNETR++, enabling accurate modeling of long-range dependencies.

## Qualitative Evaluation

As manifested in Figure 4, we employs the Swin UNETR++ with fold = 1 to perform inference on the BraTS 2021 training set. As an illustrative example, we select the 73rd slice of
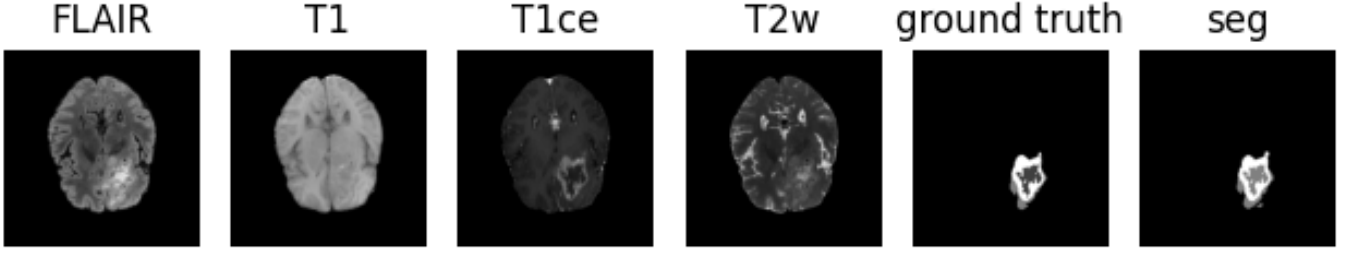
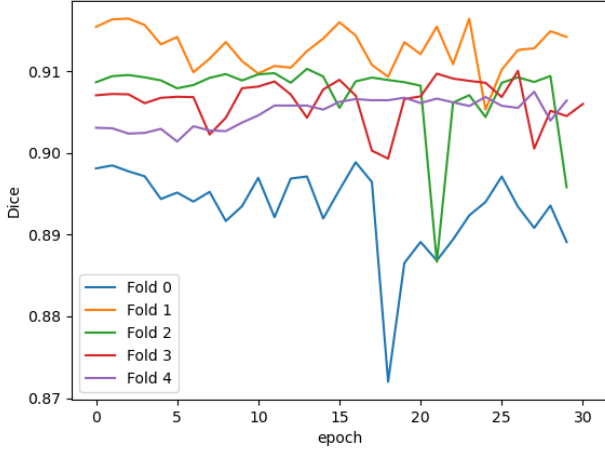Figure 4: The segmentation results for brain tumors



Figure 5: The Dice results for the 5-fold cross-validation evaluation



Figure 6: The Sensitivity results for the 5-fold cross-validation evaluation

| DSC | Our model | | | | TranBTS | | | |
|---|---|---|---|---|---|---|---|---|
| | TC | WT | ET | avg | TC | WT | ET | avg |
| Fold 0 | **0.902** | **0.92** | **0.857** | **0.893** | 0.897 | 0.910 | 0.856 | 0.883 |
| Fold 1 | **0.917** | **0.931** | **0.89** | **0.913** | 0.903 | 0.919 | 0.885 | 0.902 |
| Fold 2 | **0.912** | **0.929** | **0.881** | **0.907** | 0.898 | 0.903 | 0.866 | 0.889 |
| Fold 3 | **0.913** | **0.925** | **0.877** | **0.905** | 0.893 | 0.915 | 0.867 | 0.892 |
| Fold 4 | **0.912** | **0.927** | **0.876** | **0.905** | 0.893 | 0.915 | 0.867 | 0.892 |
| avg | **0.912** | **0.927** | **0.876** | **0.905** | 0.898 | 0.911 | 0.868 | 0.891 |

Table 2: Comparison of 5-Fold Cross-Validation Dice Similarity Coefficients

the cross-sectional section to display the input images, real labels, and output segmentation results of the four modalities. It is prone to observe that compared to the ground truth, our approach demonstrates a better overall segmentation effect on brain tumors, with only a few minor imperfections.

## Conclusion

In this paper, we present the Swin UNETR++ model based on multi-modal magnetic resonance images, which improves the effectiveness of brain tumor segmentation. A body of research trials were conducted to evaluate the performance of the proposed method, achieving state-of-the-art performance on BraTS 2021.

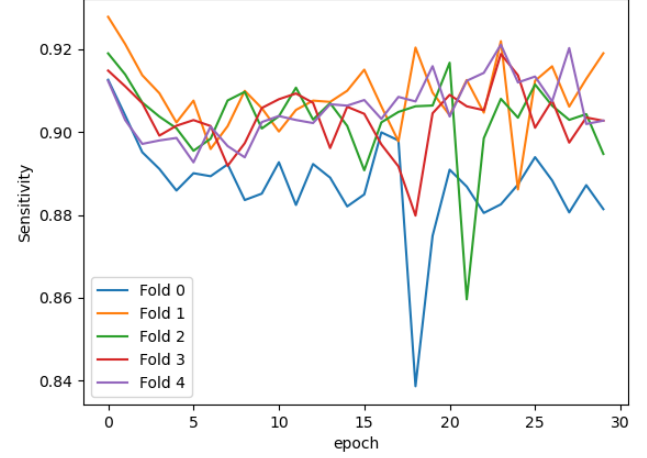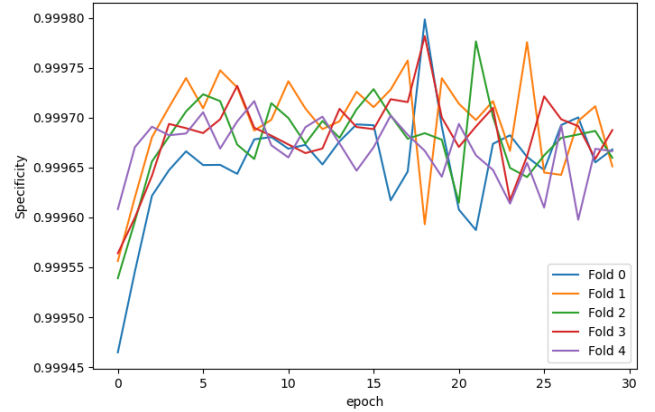In the future, the diagnosis of tumor lesions will be per-



Figure 7: The Specificity results for the 5-fold cross-validation evaluation

formed with segmented images for computer-aided diagnosis. Future research will apply the proposed method to multi-house datasets, including real-time data, to evaluate the generalization.

# References

Ahilan, A.; Manogaran, G.; Raja, C.; Kadry, S.; Kumar, S. N.; Kumar, C. A.; Jarin, T.; Krishnamoorthy, S.; Kumar, P. M.; Babu, G. C.; et al. 2019. Segmentation by fractional order darwinian particle swarm optimization based multi-level thresholding and improved lossless prediction based compression algorithm for medical images. *Ieee Access*, 7: 89570–89580.

AS, R. A.; and Gopalan, S. 2022. Comparative Analysis of Eight Direction Sobel Edge Detection Algorithm for Brain Tumor MRI Images. *Procedia Computer Science*, 201: 487–494.

Hatamizadeh, A.; Nath, V.; Tang, Y.; Yang, D.; Roth, H. R.; and Xu, D. 2021. Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images. In *International MICCAI Brainlesion Workshop*, 272–284. Springer.

Hiralal, R.; and Menon, H. P. 2016. A Survey of Brain MRI Image Segmentation Methods and the Issues Involved. 245–259. Cham: Springer International Publishing.

Ilunga-Mbuyamba, E.; Avina–Cervantes, J. G.; Garcia–Perez, A.; de Jesus Romero–Troncoso, R.; Aguirre–Ramos, H.; Cruz–Aceves, I.; and Chalopin, C. 2017. Localized active contour model with background intensity compensation applied on automatic MR brain tumor segmentation. *Neurocomputing*, 220: 84–97. Recent Research in Medical Technology Based on Multimedia and Pattern Recognition.

Isensee, F.; Jäger, P. F.; Full, P. M.; Vollmuth, P.; and Maier-Hein, K. H. 2021. nnU-Net for brain tumor segmentation. In *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: 6th International Workshop, BrainLes 2020, Held in Conjunction with MICCAI 2020, Lima, Peru, October 4, 2020, Revised Selected Papers, Part II 6*, 118–132. Springer.

Işın, A.; Direkoğlu, C.; and Şah, M. 2016. Review of MRI-based Brain Tumor Image Segmentation Using Deep Learning Methods. *Procedia Computer Science*, 102: 317–324.

Jiang, Z.; Ding, C.; Liu, M.; and Tao, D. 2020. Two-stage cascaded u-net: 1st place solution to brats challenge 2019 segmentation task. In *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: 5th International Workshop, BrainLes 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October 17, 2019, Revised Selected Papers, Part I 5*, 231–241. Springer.

Khosravanian, A.; Rahmanimanesh, M.; Keshavarzi, P.; Mozaffari, S.; and Kazemi, K. 2022. Level set method for automated 3D brain tumor segmentation using symmetry analysis and kernel induced fuzzy clustering. *Multimedia tools and applications*, 81(15): 21719–21740.

Liu, Z.; Lin, Y.; Cao, Y.; Hu, H.; Wei, Y.; Zhang, Z.; Lin, S.; and Guo, B. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, 10012–10022.

Luu, H. M.; and Park, S.-H. 2021. Extending nn-UNet for brain tumor segmentation. In *International MICCAI Brainlesion Workshop*, 173–186. Springer.

Milletari, F.; Navab, N.; and Ahmadi, S.-A. 2016. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, 565–571. Ieee.

Mlynarski, P.; Delingette, H.; Criminisi, A.; and Ayache, N. 2019. 3D convolutional neural networks for tumor segmentation using long-range 2D context. *Computerized Medical Imaging and Graphics*, 73: 60–72.

Pereira, S.; Pinto, A.; Alves, V.; and Silva, C. A. 2016. Brain Tumor Segmentation Using Convolutional Neural Networks in MRI Images. *IEEE Transactions on Medical Imaging*, 35(5): 1240–1251.

Ramesh, K.; Kumar, G. K.; Swapna, K.; Datta, D.; and Rajest, S. S. 2021. A review of medical image segmentation algorithms. *EAI Endorsed Transactions on Pervasive Health and Technology*, 7(27): e6–e6.

Saxena, S.; Kumari, N.; and Pattnaik, S. 2021. Brain tumour segmentation in FLAIR MRI using sliding window texture feature extraction followed by fuzzy C-means clustering. *International Journal of Healthcare Information Systems and Informatics (IJHISI)*, 16(3): 1–20.

Soltanian-Zadeh, H.; and Windham, J. P. 1997. A multiresolution approach for contour extraction from brain images. *Medical Physics*, 24(12): 1844–1853.

Sujji, G.; Lakshmi, Y.; and Jiji, W. 2013. MRI Brain Image Segmentation based on Thresholding. *International Journal of Advanced Computer Research*, 3: 97–101.

Wang, W.; Chen, C.; Ding, M.; Yu, H.; Zha, S.; and Li, J. 2021. Transbts: Multimodal brain tumor segmentation using transformer. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part I 24*, 109–119. Springer.

Yang, Z.; Hu, Z.; Ji, H.; Lafata, K.; Vaios, E.; Floyd, S.; Yin, F.-F.; and Wang, C. 2023. A neural ordinary differential equation model for visualizing deep neural network behaviors in multi-parametric MRI-based glioma segmentation. *Medical physics*, 50(8): 4825—4838.

Zhu, W.; and Shen, Y. 2019. A region growing segmentation approach for MRI brain image processing. In *2019 IEEE 13th International Conference on Anti-counterfeiting, Security, and Identification (ASID)*, 188–191. IEEE.