



# ENHANCED BRAIN TUMOR SEGMENTATION USING ADVANCED U-NET ARCHITECTURES

*M. Rahman, Dr. Umme kulsum ALo*

Dept. of Computer Science and Engineering, Islamic University, Kushtia-7003, Bangladesh

## ABSTRACT

Brain tumors are among the most aggressive types of tumors and they are very hard to accurately segment due mainly to their contours, complexity and location. This paper presents a novel approach to these problems by proposing a new U-shaped Convolutional Neural Networks (CNN) aimed at better integrating the morphology and anatomical structure of these brain tumor masses. Moreover, these target areas can include the whole tumor (WT), tumor core (TC), and tumor-enhancing region (ET). In addition, the framework contains a new feature fusion module that automatically adjusts the importance of different MR modalities for the purpose of improving lesion specific feature extraction. Further a scale aware attention embedding is incorporated into the network to obtain global feature representations of the features at different scales. Based on the use of multi-modality images, this method is able to achieve an improved segmentation performance which will greatly enhance the analysis of tumor in the brain in medical imaging.

## KeyWords

Brain tumor segmentation, Deep learning, Glioma, Magnetic resonance imaging (MRI), Medical image analysis, Multitask learning, U-Net architecture.

## 1.Introduction

Brain tumors are one of the most complex neoplasms facing mankind [1]. Their precise segmentation including various subregions is important for making diagnosis and creating treatment strategies. Magnetic resonance imaging (MRI) scans, which cannot be compared in soft tissue contrast, have become a common technique for diagnosing a brain tumor [2]. Nonetheless, conventional methods of tumor diagnosis, like interpretation of medical images by a human specialist, are very labor and time-consuming and rely heavily on the skills of surgeons or radiologists. The problem of assessing the area and shape of the tumor is easily solved with the inclusion of computer vision and deep learning, as verification times will be less, thus providing a better and quicker tool for the doctors [3].

In the recent time frame, convolutional neural networks (CNNs) have made inroads in the processing of medical images [4-5]. U-Net developed by Ronneberger et al [6]. Can be considered as the most notable one, as it utilizes an encoder-decoder architecture to compress context features and apply pixel wise classification. The growing usage of U-Net in slice brain tumor segmentation has resulted in the formulation of some advanced variants of the net [7-9]. For instance, Zhang et al [11]. presented an attention gate residual U-Net, AGResUNet, several features that attain attention gating to exclude undesired features and focus more on tumor locations. Similarly, Jiang et al. proposed a novel approach called two level cascaded U-Net allowing for progressive segmentation from coarse to fine by a multi-stage structure aimed at refining the segmentation process.

But 3D U-Net (3DUNet) works well for the volumetric MRI but it has some huge drawbacks as well. To begin with, with a small receptive field, it becomes difficult to detect or model self-similarities in medical images regardless of the distance. Secondly, the use of flat or stiff Convolutional kernels limits the scaling ability of convolutional operations [12][13]. Transformers, A type of models that initially were made to work with human languages, were shown to be much more effi-

cient because they solve every single drawback of the aforementioned networks. Doso-vitskiy et al. [15] were probably the first ones applying transformers in computer vision and proved them useful in such tasks [16-17].

Notwithstanding the benefits they harness, there are certain challenges with the transformers approach in medical imaging. MRI is multi-modal and contains information about the tumor from a variety of perspectives. Different physicians focus on different modalities during the interpretation of specific lesions which implies that these modalities are of different importance for segmentation purposes. This is a phenomenon that most deep learning approaches ignore and instead, simply stack the modalities together irrespective of their importance. Such methods overlook the diversity of the features available in multiple modalities and hence the segmentation models become less effective. Also, though convolutional layers are highly capable of capturing local information, they do not have the ability to model global information. Transformers can be used to model global features but they are quite expensive due to the complexity of their architecture [18].

In response to such issues, the study puts forth an updated model for a 3D MRI-based brain tumor segmentation. The segmentation of subregions as whole tumor (WT), tumor core (TC), and enhancing tumor (ET) is performed through a step-by-step approach within a multitask framework. Aiming to overcome this issue, the authors propose a feature fusion module that embeds and reorganizes effective integration features, including the interaction of components emphasizing dynamic learning. In addition, an integrated scale-aware attention mechanism is employed in the segmentation network to bridge various scales of representation while modeling information across the whole image. Such a combination of aspects addresses the drawbacks of the traditional CNNs and transformers making possible an accurate and rapid segmentation of gliomas.

The rest of the paper focuses on the following issues: the key elements concerning the proposed network architecture including feature fusion and the attention mechanism are presented in Section 2. In Section 3, the paper provides some experimental results including description of the datasets used, descriptions of evaluation metrics used, details of the ablation studies performed, and the achieved results in comparison with other state-of-the-art methods. Finally, Section 4 concludes the paper summarizing the findings addressed in the paper as well as suggested future works.

## 2.Method

### 2.1 Proposed Network Framework

The brain tumor segmentation model which is proposed to be fully automated 3D is built on a multi-branch architecture as explained in this paper. To begin with, a novel Multimodal Reorganization Module (MRM) is employed in each branch to dynamically assign a degree of significance to each modality and compensate for their unique deficits. This allows the segmentation of different lesion regions at various stages through more effective unification of multi modal data. The restructured feature maps are fed into the interspersed U-net (specifically SC-UNet) to produce the final results. The entire layout of proposed network structure is shown in Figure 1.

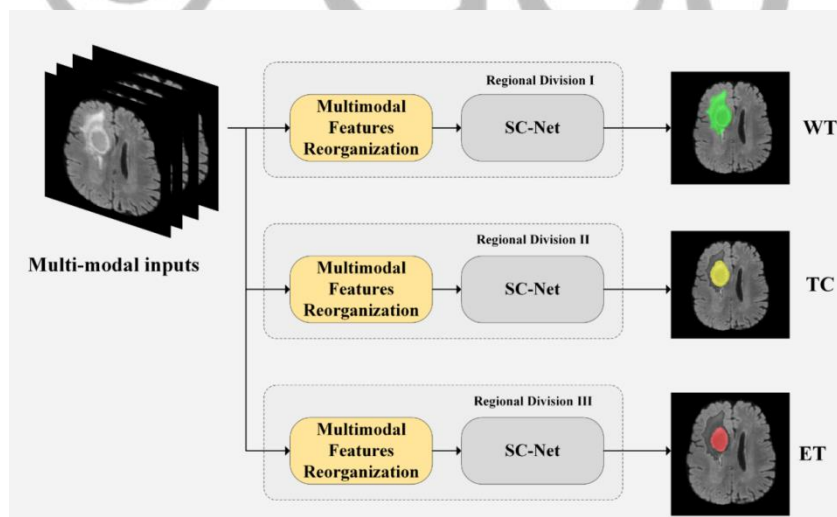


Fig. 1. The framework of overall network.

### 2.2 Multimodal Feature Reorganization Module

In order to capture more information and utilize the complementary features exhibited within multi modal MRI images, this paper proposes a multi modal feature reorganization module based on adaptive feature selection from SK-Net network concept [19]. The detailed structure of this module is illustrated in Figure 2.

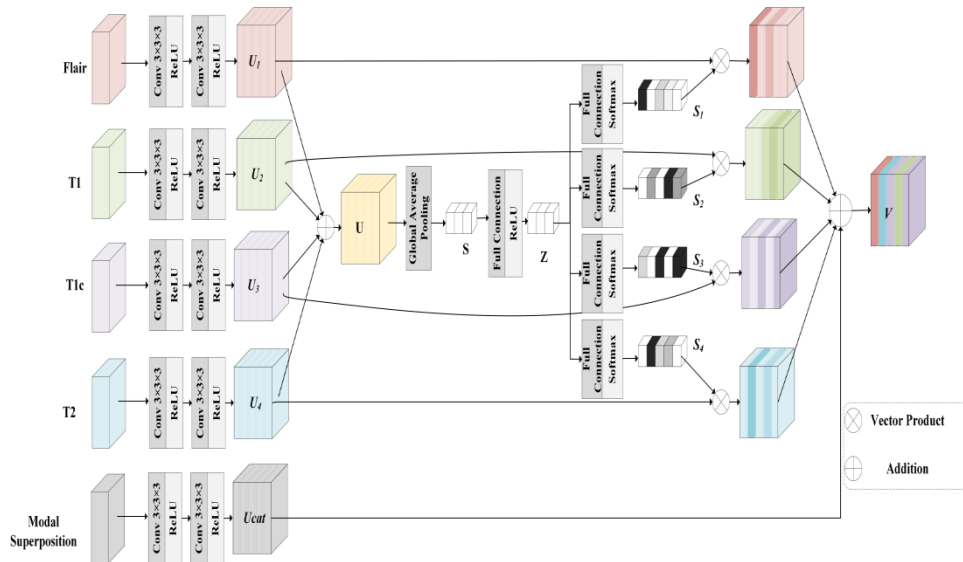


Fig.2. The structure schematic of multimodal reorganization module.

The procedure starts with the feature extraction from the four MRI modalities using 3D convolutional layers with 16 kernels of size three and ReLU activation. These features are then combined using element-wise summation to form an initial set of feature maps. The next step incorporates global context information by means of global average pooling to distill spatial dimensions into statistical features. In order to provide a more interactive and intuitive way of weighting each modality, the pooled features are sent to a fully connected layer, which melds them with a ReLU activation function creating feature representation vectors that yield modality-specific attention weights. These weights are then used to modulate the contribution of each modality by means of a spatial attention mechanism, which highlights relevant features for each task.

To obtain reorganized feature maps, the original modality features are fused with the modality adjusted features through summation. This re-organization handles the problem by treating the network to exploit complementary and shared information from different modalities while focusing on the critical lesion features. Incorporating the reorganized features into the segmentation network, the inserted module improves accuracy and stability of evaluation of the lesion segmentation resulting in good outcomes in different regions around the tumor.

## 2.3 SC-UNet Split Network

### 2.3.1 Encoders

The encoded image goes through an encoding circuit, which is designed to extract features from the input while moving from a low-level positional and contouring information into a more high-level semantic information. It has five encoding blocks, of which first four have residual convolution modules and downsampling operations, while the fifth block has only a residual convolution module. The basic architectural block comprises two  $3 \times 3 \times 3$  convolution layers from which each connects to the other via a residual connection, the two layers each has a stride which is set at one. The first layer of the convolution uses 32 kernels and with each down sampling the number of kernels is doubled. IN, together with a Leaky ReLU activation function, follows each convolution operation to ensure such stability together with provisions for non-linearity. A  $3 \times 3 \times 3$  convolution operation with a stride of 2 is used to reach any desired level of downsampled images and after all five encoder layers have been the feature map size attained is  $1/32$  of the image dimension set at the input level to ensure compact images whilst preserving the information.

### 2.3.2 Decoders

To put it more simply, the decoder takes the feature maps and makes them clearer by turning the low representations into high. Every single decoding block has an up-sampling module and a convolutional structure. The convolution module is very simple, as can be seen in Figure 4(b), two  $3 \times 3 \times 3$  convolution layers with a stride of 1

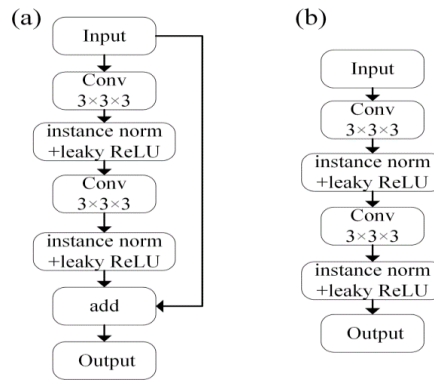


Fig.4. The structure schematics of (a) residual convolution module, and (b) convolution

So as to be consistent with the encoder, instance normalization and Leaky ReLU are performed behind every convolution layer. As for the spatial resolution, the up-sampling module has a  $2 \times 2 \times 2$  transposed convolution layer with a stride of 2, the up-sampling module enables spatial resolution to be restored. To get the final desired high-resolution outputs, the decoder goes through four up sampling stages in which the feature maps are taken and evaluated. The last convolutional with dimensions of  $1 \times 1 \times 1$  and a stride of 1 is the segmentation layer that determines the output probability map of the model. These maps are fed into an activation layer which finally produces a dilation of the end-to-end output and provides the vertices of the tumor regions.

### 2.3.3 Scaled Cross-Attention Module

In a typical setup of U-Net, skip connections carry the features extracted by the encoder and passed them directly to the decoder, reducing the amount of low-level feature information that might otherwise be lost during the information flow. However, such connections do not enhance the features interactions at varying resolutions. To resolve this limitation, a Scaled Cross-Attention (SC) module is incorporated in this work. This module is supposed to ensure a full range of multi-scale feature fusion while enjoying the global modeling ability of transformer architectures. The fusion of multi-scale features is performed using the transformer model which is known for its sequential processing of one sequence to another. This is done by turning the feature maps of several resolutions into sequences of one dimension. These sequences are joined so that the self-attention mechanism may capture the dependencies among the scales. This way, the SC module improves the communication of features on the different scales which in turn enhances the representation of semantic information in terms of the degree of the generalization attained.

### 2.3.4 Feature Extraction and Transformer-Based Global Interaction

The features extracted from the encoder are represented in a tensor  $FF$  of dimension  $C \times H \times W \times D$ , where  $C$  is the number of channels, and  $H$ ,  $W$ , and  $D$  are the spatial dimensions. To adapt these features for transformer input, the multi-scale feature maps  $F_1$ ,  $F_2$  and  $F_3$  are divided into  $N$  blocks, each of size  $4 \times 4 \times 4$ . Let  $P_i$  denote the  $i$ -th block in the feature map. Each block  $P_i$  is linearly projected using a matrix  $W_L$  and flattened into a one-dimensional vector:

$$S_i = \text{Flatten}(W_L \cdot P_i), i \in \{1, 2, \dots, N\} \quad (1)$$

The patch sequences from all scales are concatenated to form the complete patch sequence  $S$ :

$$S = [S_1; S_2; \dots; S_N] \quad (2)$$

The sequence dimension  $D$  is set to 32 in this study. Since  $SS$  lacks positional information, a learnable positional encoding  $P$  is added to each element:

$$S' = S + P \quad (3)$$

where  $P \in R^{N \times D}$  is the positional encoding. This augmented sequence  $S'$  serves as the input to the transformer modules.

### 2.3.5 Transformer Module Structure

The input sequence  $S'$  is processed through three transformer modules, each consisting of multi-head self-attention (MSA) and a feedforward neural network (FFN). For the input sequence  $S'$ , queries ( $Q$ ), keys ( $K$ ), and values ( $V$ ) are computed using learnable linear projection matrices  $W_Q, W_K, W_V$  and

$$Q = W_Q \cdot S', K = W_K \cdot S', V = W_V \cdot S' \quad (4)$$

The self-attention score matrix  $SA$  is computed as:

$$SA = \text{Softmax} \left( \frac{Q \cdot K^T}{\sqrt{d_k}} \right) \quad (5)$$

where  $d_k$  is the dimension of the key vectors, the output of the multi-head attention mechanism for each head can be given by,

$$O_h = SA \cdot V \quad (6)$$

and the outputs from all heads are concatenated and linearly transformed:

$$O=W_o \cdot [O1; O2; \dots; Oh] \tag{7}$$

The attention output is added to the input sequence via a residual connection, followed by layer normalization:

$$S_{att} = \text{LayerNorm}(S'+O) \tag{8}$$

The FFN, consisting of two fully connected layers and a ReLU activation, refines the representation:

$$\text{FFN}(S_{att}) = W2 \cdot \text{ReLU}(W1 \cdot S_{att}) \tag{9}$$

The final output of the transformer module is:

$$S_{out} = \text{LayerNorm}(S_{att} + \text{FFN}(S_{att})) \tag{10}$$

### 2.3.6 Reconstruction and Integration

The output sequence  $S_{out}$  is divided back into its original scales and converted to feature maps using an inverse linear projection  $W_L^{-1}$

$$F_r = \text{Reshape}(W_L^{-1} \cdot S_{out}) \tag{11}$$

The reconstructed feature  $F_r$  maps are passed to the decoder through skip connections, ensuring effective integration of global and multi-scale features for precise segmentation.

## 3. EXPERIMENTAL RESULTS AND ANALYSIS

To verify the effectiveness of the proposed algorithm for brain tumor segmentation, a series of comparative experiments are designed in this paper.

### 3.1 EXPERIMENTAL SETUP

For the training and testing of the proposed model in this work, the BraTS 2023 [ ] public dataset was used. This dataset consists of MRI scans of 1251 patients diagnosed with gliomas and includes low-grade gliomas (LGG) and high-grade gliomas (HGG). Each case includes preoperative MR images in four different modalities: T1-weighted (T1), T1-contrast enhanced (T1ce), T2-weighted (T2) and T2 fluid-attenuated inversion recovery (FLAIR) and also includes tumor segmentation maps labeled by qualified experts. An example of the patient's image data is shown in Fig.6.

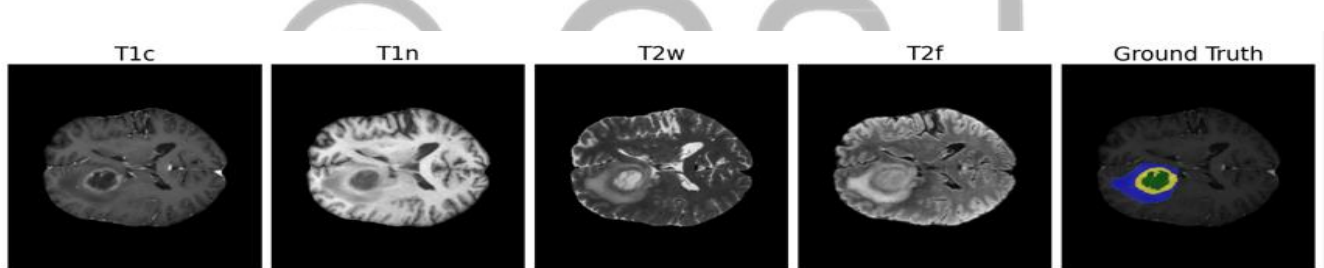


Fig.6. Four modal MRI images and physician-labeled brain tumor image (Ground Truth) of a patient.

The network model was developed using PyTorch deep learning framework with Python version 3.9 and the calculations were made on NVIDIA GeForce RTX 4090 GPU. The BraTS 2023 dataset was divided into training and testing datasets with a ratio of 85:15. The input image blocks were scaled down to the size of 128×128×128 voxels. The model was tuned in accordance with the Adam optimizer with the initial rate of 0.0001, weight decay rate of 0.00001 and 50% dropout of 0.5. The loss function was specified and encapsulated within the regularization term of 0.000001. A batch size of 2 was specified and the maximum number of training epochs was set to 300. The training process included the early stopping method for the purpose of model overfitting and performance enhancement.

### 3.2 Evaluation Metrics

The segmentation algorithm identifies and labels the WT, TC, and ET regions from patients' multimodal MRI images. These regions are hierarchically nested in a structure such that the WT includes the TC and the TC includes the ET. To measure the performance of the model's segmentation, three metrics are used: the Dice coefficient, sensitivity and 95% Hausdorff distance. Out of these, the Dice coefficient is foremost among all. The amount of overlap, as determinants of tumor segmentation adequacy, between tumors predicted and the tumors labeled on actual images is the role the Dice coefficient serves. On the other hand, the 95% Hausdorff distance analyzes the effectiveness of the segmentation more so at the boundaries of a tumor by estimating how much these boundaries predicted are in regard to the actual boundaries. According to sensitivity, which is sometimes referred to as re-call, it determines the portions of tumor regions which were correctly recognized by the model, signifying the effectiveness of the loss function used in the segmentation task

$$\text{Dice Coefficient (DC)} = \frac{2|P \cap G|}{|P| + |G|} \tag{12}$$

where  $P$  is the region of the lesion predicted by the model, and  $G$  is the ground truth lesion region. The symbol  $\cap$  denotes the logical "and" operation, indicating the overlapping region between  $P$  and  $G$ . The Dice Coefficient evaluates the degree of overlap between the predicted and actual regions, serving as the primary measure of segmentation accuracy.

$$\text{Sensitivity (Recall)} = \frac{|P \cap G|}{|G|} \tag{13}$$

where  $|G|$  represents the size of the ground truth lesion region. Sensitivity measures the proportion of the ground truth region correctly identified by the model, reflecting its ability to detect tumor regions comprehensively.

$$HD95(X, Y) = \max \{ P95(\sup_{x \in X} \inf_{y \in Y} d(x, y)), P95(\sup_{y \in Y} \inf_{x \in X} d(x, y)) \} \tag{14}$$

where  $X = \partial P$  and  $Y = \partial G$  are the boundaries of the predicted and ground truth lesion regions, respectively. The function  $d(x, y)$  computes the Euclidean distance between a point  $x \in X$  and a point  $y \in Y$ . The term P95 refers to the 95th percentile of the distance values, which reduces the influence of outliers. HD95 evaluates the precision of the model in delineating tumor boundaries, providing insight into the alignment of the predicted and actual tumor contours.

These metrics collectively assess segmentation quality from multiple perspectives: overlap (DC), detection (Sensitivity), and boundary alignment (HD95). Together, they provide a robust evaluation framework for the model's performance in segmenting brain tumor regions.

### 3.3 Ablation Experiments

To elucidate the feasibility of the design frameworks in the suggested approach, an ablation study is carried out on the 3DUNet baseline model by including the multimodal feature reorganization module and the scale cross-attention module. The detailed results are laid out in Table 1.

Table 1: Ablation Experiments Results

Models	Mean Dice (%)			Mean Sensitivity (%)			Mean Hausdorff (mm)		
	WT	TC	ET	WT	TC	ET	WT	TC	ET
UNet(base)	89.79	85.13	80.90	94.08	85.46	80.77	13.61	17.47	6.45
UNet+MR	89.89	86.45	85.71	92.16	86.45	87.13	10.26	15.48	7.14
UNet+SC	90.11	87.56	86.49	93.97	87.16	88.93	12.47	11.08	5.93
<b>Proposed</b>	<b>90.95</b>	<b>87.46</b>	<b>84.98</b>	<b>95.42</b>	<b>91.82</b>	<b>87.48</b>	<b>9.43</b>	<b>10.34</b>	<b>4.53</b>

It can be observed from the table that replacing the 3D U-Net model with tumors located multi organ 3D instance segmentation model significantly improves gineering ‘methods in motivating’ segmentation performance for every elevator panel region. This enhancement is due to the reason that the multi-modality information is combined within the module through adapting modality weights according to the relevance. Particularly, there is an increase of 4.81% in the Dice score for the Enlargement averaging procedure (ET) region. This shows that the module performs well in identifying features pertaining to specific models and is quite useful for segmentation tasks which are quite difficult.

When the Scaled Cross-Attention (SC) module is integrated into the network, even better results are obtained in the validation dataset. Enhanced boundary delineation of lesions, increased sensitivity and Dice scores all result from the SC module which combined multi-scale features and enlarge the receptive field. For example, in regards to the Tumor Core (TC) and ET regions both sensitivity and Dice scores have correspondingly increased, while the Hausdorff distance has decreased showing improved boundary accuracy.

It is clear from the results that the MR and SC modules improve both weaknesses of the individual systems and therefore, the complete system performs the best. In this regard, the proposed MR and SC modules, while capable of working independently, do considerably improve specific aspects of the segmentation performance when employed in conjunction, thus boosting the overall efficacy of the system. For instance, the Dice scores for WT, TC, and ET grow to reach 90.95%, 87.46%, and 84.98%, respectively. More important is the significant drop in the Hausdorff distance for all regions with the maximum distance across all regions to the boundary of Enhancing Tumour (ET), dropping to 4.53 mm. At each motif level, detailed analyses confirm that the deployment of each module in the case of augmenting MR and SC modules into the segmentation neural network did yield improvement on models in terms of accuracy as well as robustness towards uncertainty, as a result of them being complementary.

### 3.4 Comparative Experimental Results

To validate the effectiveness of the proposed algorithm, we compare it with six state-of-the-art CNN methods: 3D U-

Net, Attention U-Net [20], U-Net++ [21], ET-Net [22], Point-UNet [23], and TransBTS [24]. All models were evaluated under the same data preprocessing conditions and hyperparameter tuning for fair comparison. The experimental results are presented in Table 2.

Table 2: Comparative Segmentation Results Across Model

Models	Mean Dice (%)			Mean Sensitivity (%)			Mean Hausdorff (mm)		
	WT	TC	ET	WT	TC	ET	WT	TC	ET
3DUNet	89.79	85.13	80.90	94.08	87.46	80.77	13.61	17.47	6.45
AttentionUNet	88.98	82.36	81.23	91.45	86.45	76.57	14.12	11.17	8.53
UNet++	90.17	84.06	79.55	94.23	85.77	80.41	10.21	19.72	7.78
ET-Net	90.08	86.39	84.73	92.54	89.46	86.15	11.57	15.49	5.19
Point-UNet	<b>89.61</b>	87.07	86.42	91.61	87.94	87.18	10.16	17.13	8.75
TransBTS	90.15	86.43	<b>86.47</b>	92.41	<b>94.92</b>	86.39	16.57	14.16	6.65
<b>Proposed</b>	<b>90.95</b>	<b>87.46</b>	84.98	<b>95.42</b>	91.87	<b>87.48</b>	<b>9.43</b>	<b>10.34</b>	<b>4.53</b>

Apart from the considered metrics, the quantitative evaluation input detailed in the previous section is re-stated in Table II as the criteria for the segmentations very similar to the Sistine Chapel Model, but the major difference is that segments are also defined for all images led to unprecedented performance in all tasks. The average Dice values, for the whole tumor, tumor core and enhancing the tumor segments come off as respectively. This indicates that the Multimodal Feature Reorganization (MR) module succeeded in enhancing the network’s capacity of assigning weights for various regions and enabled robust segmentation of WT as well as TC and ET, particularly the latter in the round. Additionally, ET inaccuracy values were 1.16%, 2.33%, and 4.08% in contraction UP models off the baseline 3D U-Net model along with WT, TC respectively.

The Scale Cross Attention (SC) module improves the model in terms of feature fusion across scales thus boosting the quality of the features represented in it multiscale feature networks. The developed architecture achieved the highest structural similarity scores and the lowest Hausdorff distances (9.43 mm for WT, 10.34 mm for TC, and 4.53 mm for ET) as well as very high sensitivity scores (95.42% for WT, 91.87% for TC, and 87.48% for ET). This indicates higher accuracy of segmentation and precision of the boundaries. The results clearly demonstrate that the proposed algorithm consistently ranks better than the rest algorithms across the complexity spectrum especially in segmentation tasks that are quite complicated.

The results of Figures 7 depict the segmentation results of tumor located in the brain using 3D U-Net, TransBTS and the proposed approach. It could be observed from the figures that the contours produced by the proposed model are more refined and precise, making good use of the multimodal MRI and correctly identifying the tumor tissues. The network eliminates false-positive areas by fusing multi-scale features and extending receptive areas, resulting in an output that closely matches the expected outcomes.

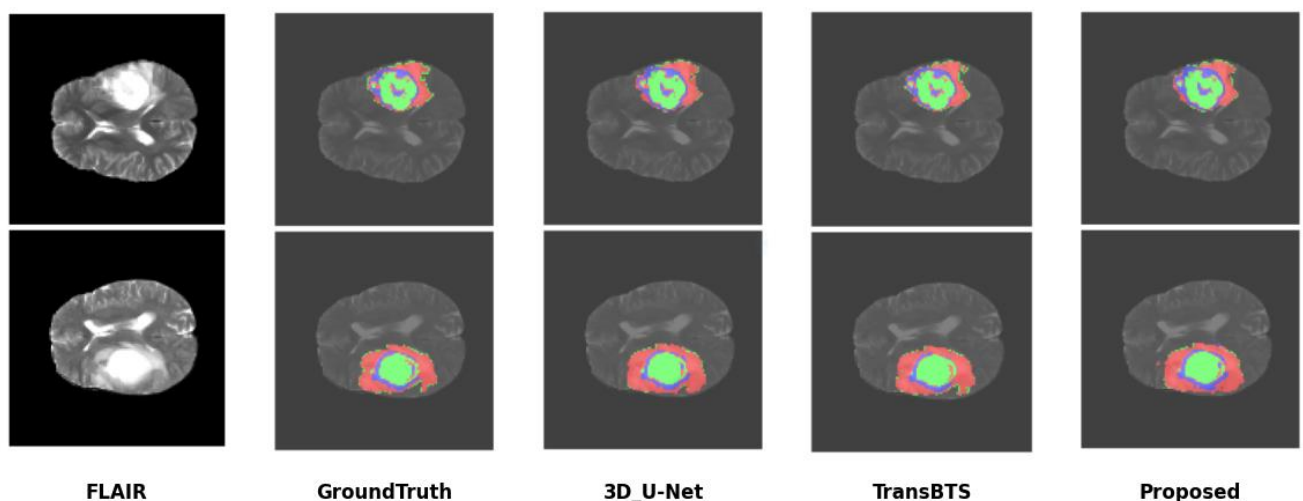


Fig.7. 2D segmentation results of different models.

## 4. Conclusion

The paper develops an enhanced model for segmenting brain tumors by embedding the Multimodal Feature Reorganization (MR) module and the Scale Cross-Attention (SC) module into the 3D U-Net architecture. The MR module allows for adaptive weighting of different MRI modalities making the extraction of features related to tumor subregions more effective. The SC module improves model's ability to incorporate multi-scale information and increases the receptive field thereby enhancing the accuracy of the segmentation and the boundaries' definition. The Dice scores, sensitivity and Hausdorff distances obtained in segmenting the Whole Tumor (WT), Tumor Core (TC) and Enhancing Tumor (ET) regions reveal the superiority of the proposed approach over state-of-the-art methods particularly when subjected to extensive experiments on the BraTS 2023 dataset. The results further indicate the robustness of the model with respect to complex, hierarchical segmentation tasks and the occurrence of false positive results is significantly reduced.

## Acknowledgment

The authors thank the laboratory staff for their generous support and assistance in providing access to lab facilities beyond regular hours.

## References

- [1] Z. Liu, L. Tong, L. Chen, Z. Jiang, F. Zhou, Q. Zhang, X. Zhang, Y. Jin, H.J.C. Zhou, Deep learning-based brain tumor segmentation: a survey, *Complex intelligent systems*, 9 (2023), pp.1001-1026.
- [2] H. Luo, D. Zhou, Y. Cheng, S. Wang, MPEDA-Net: A lightweight brain tumor segmentation network using multi-perspective extraction and dense attention, *Biomedical Signal Processing and Control*, 91 (2024), Article 106054.
- [3] Z. Liu, J. Wei, R. Li, J. Zhou, Learning multi-modal brain tumor segmentation from privileged semi-paired MRI images with curriculum disentanglement learning, *Computers in Biology and Medicine*, 159 (2023), Article 106927.
- [4] Y. Zhang, R. Xie, I. Beheshti, X. Liu, G. Zheng, Y. Wang, Z. Zhang, W. Zheng, Z. Yao, B. Hu, Improving brain age prediction with anatomical feature attention-enhanced 3D-CNN, *Computers in Biology and Medicine*, 169 (2024), Article 107873.
- [5] Y. Peng, X. Hu, X. Hao, P. Liu, Y. Deng, Z. Li, Spider-Net: High-resolution multi-scale attention network with full-attention decoder for tumor segmentation in kidney, liver and pancreas, *Biomedical Signal Processing and Control*, 93 (2024), Article 106163.
- [6] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III* 18, Springer, (2015), pp. 234-241.
- [7] L. Fang, X. Wang, Multi-input Unet model based on the integrated block and the aggregation connection for MRI brain tumor segmentation, *Biomedical Signal Processing and Control*, 79 (2023), Article 104027.
- [8] R. Raza, U.I. Bajwa, Y. Mehmood, M.W. Anwar, M.H. Jamal, dResU-Net: 3D deep residual U-Net based brain tumor segmentation from multimodal MRI, *Biomedical Signal Processing Control*, 79 (2023), Article 103861.
- [9] M.U. Rehman, J. Ryu, I.F. Nizami, K.T.J.C.i.B. Chong, *Medicine*, RAAGR2-Net: A brain tumor segmentation network using parallel processing of multiple spatial frames, 152 (2023), Article 106426.
- [10] J. Zhang, Z. Jiang, J. Dong, Y. Hou, B.J.I.A. Liu, Attention gate resU-Net for automatic MRI brain tumor segmentation, 8 (2020), pp.58533-58545.
- [11] Z. Jiang, C. Ding, M. Liu, D. Tao, Two-Stage Cascaded U-Net: 1st Place Solution to BraTS Challenge 2019 Segmentation Task, Springer International Publishing, Cham, (2020), pp.231-241.
- [12] Z. Zhu, X. He, G. Qi, Y. Li, B. Cong, Y. Liu, Brain tumor segmentation based on the fusion of deep semantics and edge information in multimodal MRI, *Information Fusion*, 91 (2023), pp.376-387.
- [13] X. Liu, C. Yao, H. Chen, R. Xiang, H. Wu, P. Du, Z. Yu, W. Liu, J. Liu, D. Geng, BTSC-TNAS: A neural architecture search-based transformer for brain tumor segmentation and classification, *Computerized Medical Imaging*, 110 (2023), Article 102307.
- [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A.N. Gomez, Ł. Kaiser, I.s. Polosukhin, Attention is all you need, *Advances in neural information processing* 30 (2017).
- [15] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, an image is worth 16x16 words: Transformers for image recognition at scale, *arXiv preprint arXiv:11929*, (2020).
- [16] K. He, C. Gan, Z. Li, I. Rekik, Z. Yin, W. Ji, Y. Gao, Q. Wang, J. Zhang, D. Shen, Transformers in medical image analysis, *Intelligent Medicine*, 3 (2023), pp.59-78.
- [17] X. Liu, Y. Ding, Y. Zhang, J. Tang, Multi-scale local-global transformer with contrastive learning for biomarkers segmentation in retinal OCT images, *Biocybernetics and Biomedical Engineering*, 44 (2024), pp.231-246.



- [18] D. Wang, Q. Zhang, Y. Xu, J. Zhang, B. Du, D. Tao, L. Zhang, Advancing plain vision transformer toward remote sensing foundation model, *IEEE Transactions on Geoscience Remote Sensing*, 61 (2022), pp.1-15.
- [19] de Verdier, M.C., Saluja, R., Gagnon, L., LaBella, D., Baid, U., Tahon, N.H., Foltyn-Dumitru, M., Zhang, J., Alafif, M., Baig, S. and Chang, K., 2024. The 2024 Brain Tumor Segmentation (BraTS) Challenge: Glioma Segmentation on Post-treatment MRI. *arXiv preprint arXiv:2405.18368*.

© GSJ