# 3D Medical Multi-modal Segmentation Network Guided by Multi-source Correlation Constraint

Tongxue Zhou[*†‡], Stéphane Canu[*‡], Pierre Vera[§] and Su Ruan[†‡]

[*]INSA Rouen, LITIS - Apprentissage, Rouen 76800, France
[†]Université de Rouen Normandie, LITIS - QuantIF, Rouen 76183, France
[‡]Normandie Univ, INSA Rouen, UNIROUEN, UNIHAVRE, LITIS, France
[§] Department of Nuclear Medicine, Henri Becquerel Cancer Center, Rouen, 76038, France

*Abstract*—In the field of multimodal segmentation, the correlation between different modalities can be considered for improving the segmentation results. In this paper, we propose a multi-modality segmentation network with a correlation constraint. Our network includes N model-independent encoding paths with N image sources, a correlation constrain block, a feature fusion block, and a decoding path. The model independent encoding path can capture modality-specific features from the N modalities. Since there exists a strong correlation between different modalities, we first propose a linear correlation block to learn the correlation between modalities, then a loss function is used to guide the network to learn the correlated features based on the linear correlation block. This block forces the network to learn the latent correlated features which are more relevant for segmentation. Considering that not all the features extracted from the encoders are useful for segmentation, we propose to use dual attention based fusion block to recalibrate the features along the modality and spatial paths, which can suppress less informative features and emphasize the useful ones. The fused feature representation is finally projected by the decoder to obtain the segmentation result. Our experiment results tested on BraTS-2018 dataset for brain tumor segmentation demonstrate the effectiveness of our proposed method.

## I. INTRODUCTION

Multimodal segmentation using a single model remains challenging due to the different image characteristics of different modalities. A key challenge is to exploit the latent correlation between modalities and to fuse the complementary information to improve the segmentation performance. In this paper, we proposed a method to exploit the multi-source correlation and apply it to brain tumor segmentation task.

A brain tumor is a growth of cells in the brain that multiplies in an abnormal, uncontrollable way, which is one of the most lethal cancers in the world. Today, an estimated 700,000 people in the United States are living with a primary brain tumor, and over 87,000 more will be diagnosed in 2020[1]. Gliomas are the most common brain tumors that arise from glial cells. According to the malignant degree of gliomas [1], they can be categorized into two grades: low-grade gliomas (LGG) and high-grade gliomas (HGG), the former one tend

---

[1]NBTS: National Brain Tumor Society

to be benign, grow more slowly with lower degrees of cell infiltration and proliferation, the latter one are malignant, more aggressive and need immediate treatment, moreover, the five-year relative survival rate of gliomas is only 6.8%. Therefore, early diagnosis of brain tumors is highly desired in clinical practice for better treatment planning.

Magnetic Resonance Imaging (MRI) is commonly used in radiology to diagnose brain tumors, it is a non-invasive and good soft tissue contrast imaging modality, which provides invaluable information about shape, size, and localization of brain tumors without exposing the patient to a high ionization radiation [2]–[4]. The commonly used sequences are T1-weighted (T1), contrast-enhanced T1-weighted (T1c), T2-weighted (T2) and Fluid Attenuation Inversion Recovery (FLAIR) images. In this work, we refer to these images of different sequences as modalities. Different modalities can provide complementary information to analyze different sub-regions of gliomas. For example, T2 and FLAIR highlight the tumor with peritumoral edema, designated whole tumor. T1 and T1c highlight the tumor without peritumoral edema, designated tumor core. An enhancing region of the tumor core with hyper-intensity can also be observed in T1c, designated enhancing tumor core. Therefore applying multi-modal images can reduce the information uncertainty and improve clinical diagnosis and segmentation accuracy.

Inspired by a fact that, there is strong correlation between multi MR modalities, since the same scene (the same patient) is observed by different modalities [5]. We propose a 3D multimodal brain segmentation network guided by multi-source correlation constrain. The main contributions of our method are four folds: 1) A correlation block is introduced to discover the latent multi-source correlation between modalities, making the features more relevant for segmentation. 2) A dual attention based fusion strategy is proposed to recalibrate the feature representation along modality-wise and spatial-wise. 3) A correlation based loss function is proposed to aide the segmentation network to extract the correlated feature representation for a better segmentation. 4) The first 3D multimodal brain tumor segmentation network guided by multi-source correlation constrain is proposed.

## A. Related Work

A wide range of approaches for brain tumor segmentation, such as probability theory [5], kernel feature selection [6], belief function [7] based on [8], random forests [9], conditional random fields [10] and support vector machines [11] have been developed with success. However, brain tumor segmentation is still a challenging task due to three reasons: (1) The brain anatomy structure varies from patients to patients. (2) The variability across size, shape, and texture of gliomas. (3) The variability in intensity range and low contrast in qualitative MR imaging modalities (see Fig. 1).
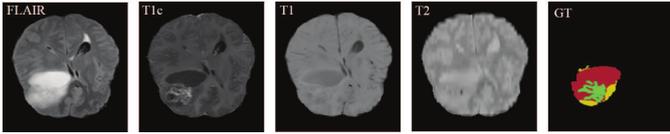


Fig. 1. Example of data from a training subject. The first four images from left to right show the MRI modalities: Fluid Attenuation Inversion Recovery (FLAIR), contrast enhanced T1-weighted (T1c), T1-weighted (T1), T2-weighted (T2) images, and the fifth image is the ground truth labels created by experts. The color is used to distinguish the different tumor regions: red: Necrotic and Non-enhancing tumor, yellow: edema, green: enhancing tumor, black: healthy tissue and background.

Recently, with a strong feature learning ability, deep learning-based approaches have become more prominent for brain tumor segmentation. Cui et al. [12] proposed a cascaded deep learning convolutional neural network consisting of two sub-networks. The first network is to define the tumor region from a MRI slice and the second network is used to label the defined tumor region into multiple sub-regions. Zhao et al. [13] integrated fully convolutional neural networks (FCNNs) [14] and conditional random fields to segment brain tumor. Havaei et al. [15] implemented a two-pathway architecture that learns about the local details of the brain tumor as well as the larger context feature. Wang et al. [16] proposed to decompose the multi-class segmentation problem into a sequence of three binary segmentation problems according to the sub-region hierarchy. Kamnitsas et al. [17] proposed an efficient fully connected multi-scale CNN architecture named DeepMedic, which reassembles a high resolution and a low resolution pathway to obtain the segmentation results. Furthermore, they used a 3D fully connected conditional random field to effectively remove false positives. Kamnitsas et al. [18] introduced EMMA, an ensemble of multiple models and architectures including DeepMedic, FCNs and U-Net. Myronenko et al. [19] proposed a segmentation network for brain tumor from multimodal 3D MRIs, where variational auto-encoder branch is added into the U-net to further regularize the decoder in the presence of limited training data.

For multimodal segmentation task, exploiting the complimentary information from different modalities play an import role in the final segmentation accuracy. As presented in [20], the multi-modal segmentation network architectures can be categorized into single-encoder-based method and multi-encoder-based method. The single-encoder-based method [18],

[21] directly integrates the different multi-modality images channel-wise in the input space, while the correlations between different modalities are not well exploited. However, the multi-encoder-based method [22], allows to separately extract individual feature information by applying multiple modality-specific encoders, and to fuse them with specific fusion strategy to emphasize the useful information for the segmentation task. According to [23], multi-encoder-based method has better performance than single-encoder-based method, which can learn more complementary and cross-modal interdependent features. However, not all features extracted from the encoder are useful for segmentation. Therefore, it is necessary to find an effective way to fuse features, we focus on the extraction of the most informative features for segmentation. To this end, we propose to use the attention mechanism, which can be viewed as a tool being capable to take into account the most informative feature representation. Channel attention modules and spatial attention modules are the commonly used attention mechanisms. The former one learn a channel-wise feature representation that quantifies the relative importance of each channel's features [24]–[26]. The latter one, spatial attention modules, learn the feature representation in each position by weighted sum the features of all other positions [27]–[29]. However, the methods mentioned above evaluated the attention mechanism only on the single-modal image dataset and don't consider the fusion issue on the multi-modal medical images. In this paper, we propose to apply the attention mechanism on the multi-modality brain tumor dataset. To learn the contributions of the feature representations from different modalities, we propose a dual attention based fusion block to selectively emphasize feature representations, which consists of a modality attention module and a spatial attention module. The proposed fusion block uses the modality-specific features to derive a modality-wise and a spatial-wise weight map that quantify the relative importance of each modality's features and also of the different spatial locations in each modality. These fusion maps are then multiplied with the modality-specific feature representations to obtain a fused representation of the complementary multi-modality information. In this way, we can discover the most relevant characteristics to aide the segmentation.

For multi-modal MR brain tumor segmentation, since the four MR modalities are from the same patient, there exists a strong correlation between modalities [5]. In this paper, our goal is to exploit and utilize the correlation between modalities to improve the segmentation performance. Therefore, we first exploit the correlation between each two modalities and then utilize a loss function to guide the segmentation network to learn the correlated features to enhance the segmentation result. To the best of our knowledge, this is the first work which is capable of utilizing the latent multi-source correlation to help the segmentation.

## II. METHOD

Our network is based on our previous work [23], which used a multi-encoder based network to deal with the multi-

model fusion issue. In this paper, we aim to exploit the multi-source correlation between modalities and utilize the correlation to constrain the network to learn more effective feature so as to improve the segmentation performance. To learn complementary features and cross-modal inter-dependencies from multi-modality MRIs, we applied the multi-encoder based framework. It takes 3D MRI modality as input in each encoder. Each encoder can produce a modality-specific feature representation, at the lowest level of the network, the linear correlation block is first used to exploit the latent multi-source correlation, then a well-designed loss function is applied to guide the network to learn the effective feature information. Then, all the modality-specific feature representations are concatenated to the fusion block at each layer. With the assistance of the dual attention fusion block, the feature representations will be separated along modality-wise and space-wise, and the most informative feature is obtained as the shared latent representation, and finally it is projected by decoder to the label space to obtain the segmentation result. The pipeline of our method is described in Fig. 2.

### A. Architecture Design

It's likely to require different receptive fields when segmenting different regions in an image, a standard U-Net can't get enough semantic features due to the limited receptive field. Inspired by dilated convolution, we use residual block with dilated convolutions (rate = 2, 4) (res_dil block) on both encoder part and decoder part to obtain features at multiple scale. The encoder includes a convolutional block, a res_dil block followed by skip connection. All convolutions are $3 \times 3 \times 3$. Each decoder level begins with up-sampling layer followed by a convolution to reduce the number of features by a factor of 2. Then the upsampled features are combined with the features from the corresponding level of the encoder part using concatenation. After the concatenation, we use the res_dil block to increase the receptive field. In addition, we employ deep supervision [21] for the segmentation decoder by integrating segmentation layers from different levels to form the final network output. The proposed network architecture is described in Fig. 3.

### B. Correlation Constrain Block

Inspired by a fact that, there is strong correlation between multi MR modalities, since the same scene (the same patient) is observed by different modalities [5]. From Fig. 4 presenting joint intensities of the MR images, we can observe a strong correlation in intensity distribution between each two modalities. To this end, it's reasonable to assume that a strong correlation also exists in latent representation between modalities. Therefore, we introduce a Correlation Constrain (CC) block, which consists of a Linear Correlation (LC) block (see Fig. 5) to discover the latent correlation and a correlation loss to constrain the correlation between modalities. For simplicity, we present the CC block using two modalities. The input modality $\{X_i, ..., X_n\}$, where $n = 4$, is first input to the independent encoders (with learning parameters $\theta$) to

learn the modality-specific representation $Z_i(X_i|\theta_i)$. Then, a network with two fully connected network with LeakyReLU, maps the modality-specific representation $Z_i(X_i|\theta_i)$ to a set of independent parameters $\Gamma_i = \{\alpha_i, \beta_i\}$, $i = 1, ..., n$. Finally the linear correlation representation of $j$ modality $F_j(X_j|\theta_j)$ can be obtained via linear correlation Equation 1.

Since we have four modalities, and each two modalities have a strong linear correlation, we only need to learn three pairs of correlation expressions from each two modalities. Then, the Kullback–Leibler divergence (Equation 2) is used as the correlation loss to constrain the distributions between the estimated correlation representation and the original feature representation, which enables the segmentation network to learn the correlated feature representation to improve the segmentation performance.

$$F_j(X_j|\theta_j) = \alpha_i \odot Z_i(X_i|\theta_i) + \beta_i, (i \neq j) \qquad (1)$$

$$L_{correlation} = \sum_{x \in X} P(x) log \frac{P(x)}{Q(x)} \qquad (2)$$

where $P$ and $Q$ are probability distributions of $Z_i$ and $F_j$, respectively, which defined on the same probability space $X$.

### C. Dual Attention based fusion strategy

The purpose of fusion is to stand out the most important features from different source images to highlight regions that are greatly relevant to the target region. Since different MR modalities can identify different attributes of the target tumor to provide complementary information. In addition, from the same MR modality, we can learn different content at different locations. Inspired by the attention mechanism [27], we propose a dual attention based fusion block to enable a better integration of the complementary information between modalities, which consists of a modality attention module and a spatial attention module, the architecture is described in Fig. 6.

The individual feature representations learned by four encoders ($Z_1$, $Z_2$, $Z_3$, $Z_4$) are first concatenated to obtain the input feature representation $Z = [Z_1, Z_2, Z_3, Z_4]$, $Z_k \in R^{H \times W}$. Note that, in the lowest level of the network, there are four modality-specific feature representations ($Z_1$, $Z_2$, $Z_3$, $Z_4$), in the other levels, the upsamlping layer in the decoder path is also concatenated with the modality-specific feature representations to obtain the input feature representation $Z = [Z_1, Z_2, Z_3, Z_4, Z_5]$, $Z_k \in R^{H \times W}$, for simplicity, in the following, we describe the fusion block with the four modality-specific feature representations.

In the modality attention module, a global average pooling is first performed to produce a tensor $g \in R^{1 \times 1 \times 4}$, which represents the global spatial information of the feature representation, with its $k^{th}$ element

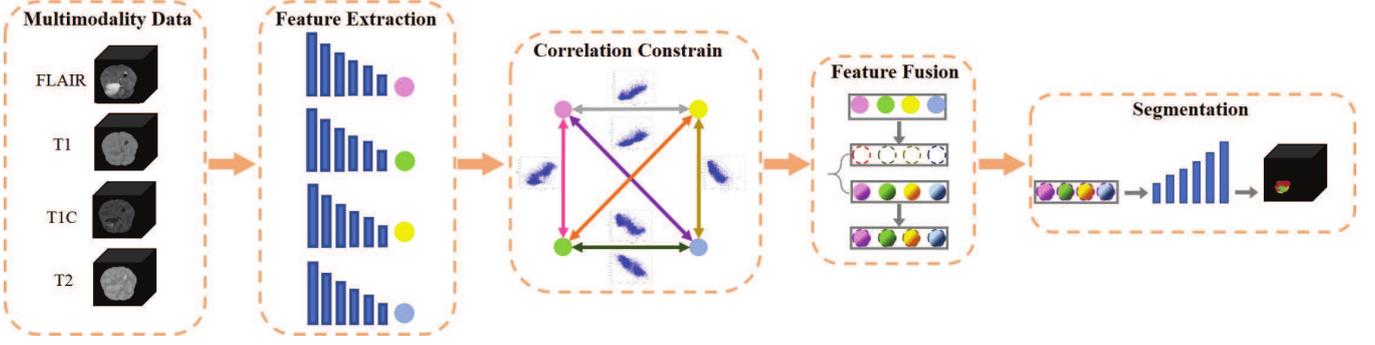$$g_k = \frac{1}{H \times W} \sum_i^H \sum_j^W Z_k(i, j) \qquad (3)$$

Fig. 2. The pipeline of the proposed method, consisting of feature extraction, correlation constrain and fusion fusion block, 4 color circles represent 4 modality feature representations.
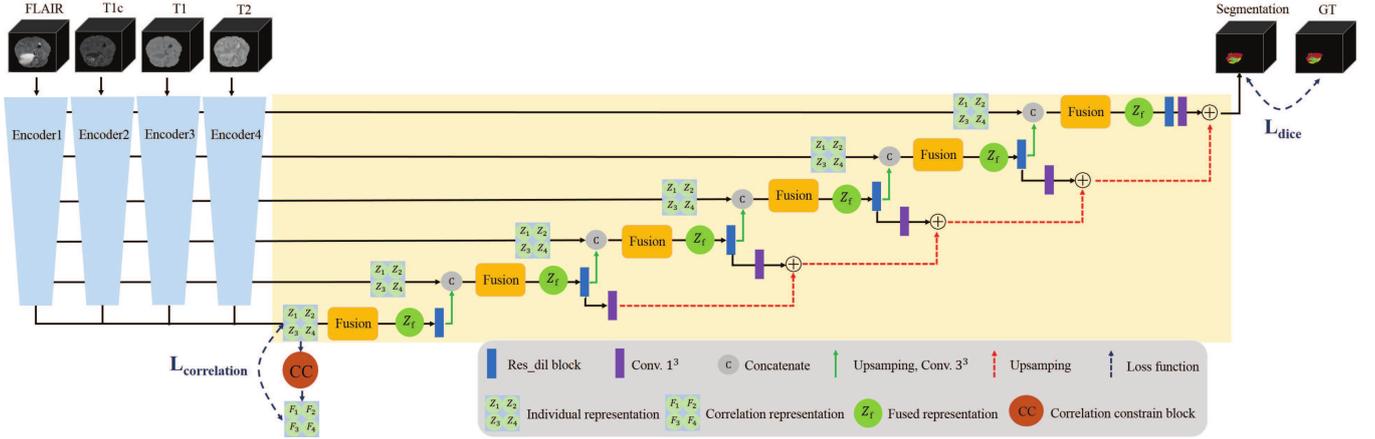


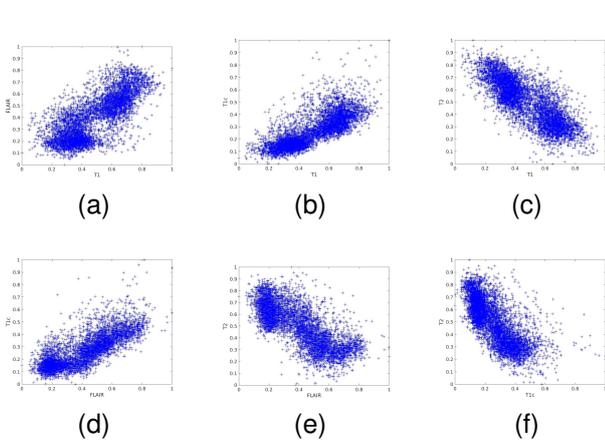Fig. 3. Overview of our proposed segmentation network framework.



Fig. 4. Joint intensity distributions of MR images: (a) T1-FLAIR, (b) T1-T1c,(c) T1-T2, (d) FLAIR-T1c, (e) FLAIR-T2, (f) T1c-T2. The intensity of the first modality is read on abscissa axis and that of the second modality on the ordinate axis.
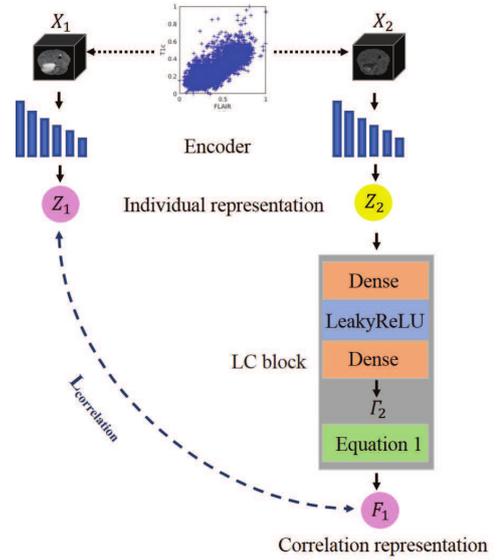


Fig. 5. Architecture of Correlation Constrain (CC) block, which consists of Linear Correlation (LC) block and a correlation constrain loss.

Then two fully-connected layers are applied to encode the modality-wise dependencies, $\hat{g} = W_1(\delta(W_2 g))$, with $W_1 \in R^{4 \times 2}$, $W_2 \in R^{2 \times 4}$, being weights of two fully-connected layers and the ReLU operator $\delta(\cdot)$, $\hat{g}$ is then passed through the sigmoid layer to obtain the modality-wise weights, which will

be applied to the input representation $Z$ through multiplication to achieve the modality-wise features $Z_m$, and the $\sigma(\hat{g}_k)$

indicates the importance of the $i$ modality of the feature representation.

$$Z_m = [\sigma(\hat{g_1})Z_1, \sigma(\hat{g_2})Z_2, \sigma(\hat{g_3})Z_3, \sigma(\hat{g_4})Z_4,] \quad (4)$$

In the spatial attention module, the feature representation can be considered as $Z = [Z^{1,1}, Z^{1,2}, ..., Z^{i,j}, ..., Z^{H,W}]$, $Z^{i,j} \in R^{1 \times 1 \times 4}$, $i \in 1, 2, ..., H$, $j \in 1, 2, ..., W$, and then a convolution operation $q = W_s \star Z$, $q \in R^{H \times W}$ with weight $W_s \in R^{1 \times 1 \times 4 \times 1}$, is used to squeeze the spatial domain, and to produce a projection tensor, which represents the linearly combined representation for all modalities for a spatial location. The tensor is finally passed through a sigmoid layer to obtain the space-wise weights, $\sigma(q_{i,j})$ indicates the importance of the spatial information $(i, j)$ of the feature representation.

$$Z_s = [\sigma(q_{1,1})Z^{1,1}, ..., \sigma(q_{i,j})Z^{i,j}, ..., \sigma(q_{H,W})Z^{H,W}] \quad (5)$$

Finally, the learned fused feature representation is obtained by adding the modality-wise feature representation and space-wise feature representation.

$$Z_f = Z_m + Z_s \quad (6)$$

From Fig. 6, we can observe the target tumor's characteristics in the four independent feature representations are not obvious, however, the modality attention module stands out the different attributes of the modalities to provide complementary information, for example, the FLAIR modality highlights the edema region and T1c modality stand out the tumor core region. In the spatial attention module, all the locations related to the target tumor region are highlighted, In this way, we can discover the most relevant characteristics between modalities. Furthermore, the proposed fusion block can be directly adapted to any multi modal fusion problem.
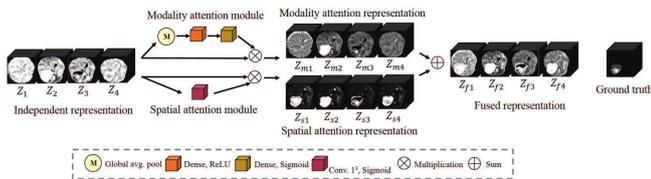


Fig. 6. Proposed dual attention fusion block. The individual feature representations ($Z_1$, $Z_2$, $Z_3$, $Z_4$) are first concatenated, then they are recalibrated along modality attention module and spatial attention module to achieve the modality attention representation $Z_m$ and spatial attention representation $Z_s$, final they are added to obtain the fused feature representation $Z_f$.

## III. DATA AND IMPLEMENTATION DETAILS

### A. Data

The datasets used in the experiments come from BraTS 2018 dataset. The training set includes 285 patients, each patient has four image modalities including T1, T1c, T2 and FLAIR. Following the challenge, four intra-tumor structures have been grouped into three mutually inclusive tumor regions: (a) whole tumor (WT) that consists of all tumor tissues, (b) tumor core (TC) that consists of the enhancing tumor, necrotic and non-enhancing tumor core, and (c) enhancing tumor (ET). The provided data have been pre-processed by organisers: co-registered to the same anatomical template, interpolated to the same resolution ($1mm^3$) and skull-stripped. The ground truth have been manually labeled by experts. We did additional pre-processing with a standard procedure. The N4ITK [30] method is used to correct the distortion of MRI data, and intensity normalization is applied to normalize each modality of each patient. To exploit the spatial contextual information of the image, we use 3D image, crop and resize it from $155 \times 240 \times 240$ to $128 \times 128 \times 128$.

### B. Implementation Details

Our network is implemented in Keras with a single Nvidia GPU Quadro P5000 (16G). The models are optimized using the Adam optimizer(initial learning rate = 5e-4) with a decreasing learning rate factor 0.5 with patience of 10 epochs, to avoid over-fitting, early stopping is used when the validation loss isn't improved for 50 epoch. We randomly split the dataset into 80% training and 20% testing.

### C. The choices of loss function

The network is trained by the overall loss function as follow:

$$L_{total} = L_{dice} + \lambda L_{correlation} \quad (7)$$

where $\lambda$ is the trade-off parameters weightig the importance of each component, which is set as 0.1 in our experiment.

For segmentation, we use dice loss to evaluate the overlap rate of prediction results and ground truth.

$$L_{dice} = 1 - 2 \frac{\sum_{i=1}^{C} \sum_{j=1}^{N} p_{ic} g_{ic} + \epsilon}{\sum_{i=1}^{C} \sum_{j=1}^{N} p_{ic} + g_{ic} + \epsilon} \quad (8)$$

where $N$ is the set of all examples, $C$ is the set of the classes, $p_{ic}$ is the probability that pixel $i$ is of the tumor class $c$ and $p_{i\bar{c}}$ is the probability that pixel $i$ is of the non-tumor class $\bar{c}$. The same is true for $g_{ic}$ and $g_{i\bar{c}}$, and $\epsilon$ is a small constant to avoid dividing by 0.

### D. Evaluation metrics

To evaluate the proposed method, two evaluation metrics: Dice Score and Hausdorff distance are used to obtain quantitative measurements of the segmentation accuracy:
1) Dice Score: It is designed to evaluate the overlap rate of prediction results and ground truth. It ranges from 0 to 1, and the better predict result will have a larger Dice value.

$$Dice = \frac{2TP}{2TP + FP + FN} \quad (9)$$

where $TP$ represents the number of true positive voxels, $FP$ represents the number of false positive voxels, and $FN$ represents the number of false negative voxels.
2) Hausdorff distance (HD): It is computed between boundaries of the prediction results and ground-truth, it is an

indicator of the largest segmentation error. The better predict result will have a smaller HD value.

$$HD = \max\{sup_{r \in \partial R} d_m(s, r), sup_{s \in \partial S} d_m(r, s)\} \quad (10)$$

where $\partial S$ and $\partial R$ are the sets of tumor border voxels for the predicted and the real annotations, and $d_m(v, v)$ is the minimum of the Euclidean distances between a voxel $v$ and voxels in a set $v$.

## IV. EXPERIMENT RESULTS

We conduct a series of comparative experiments to demonstrate the effectiveness of our proposed method and compare it to other approaches. In Section IV-A1, we first perform an ablation experiment to see the importance of our proposed components and demonstrate that adding the proposed components can enhance the segmentation performance. In Section IV-A2, we compare our method with the state-of-the-art methods. In Section IV-B, the qualitative experiment results further demonstrate that our proposed method can achieve a promising segmentation result.

### A. Quantitative Analysis

To prove the effectiveness of our network, we first did an ablation experiment to see the effectiveness of our proposed components, and then we compare our method with the state-of-the-art methods. All the results are obtained by online evaluation platform[2].

*1) Effectiveness of Individual Modules:* To assess the performance of our method, and see the importance of the proposed components in our network, including dual attention fusion strategy and correlation constrain block, we did an ablation experiment, our network without dual attention fusion strategy and correlation constrain block is denoted as baseline. From Table I, we can observe the baseline method achieves Dice Score of 0.726, 0.867, 0.766 for enhancing tumor, whole tumor, tumor core, respectively. When the dual attention fusion strategy is applied to the network, we can see an increase of Dice Score and Hausdorff Distance across all tumor regions with an average improvement of 0.85% and 6.44%. respectively. The major reason is that the proposed fusion block can help to emphasize the most important representations from the different modalities across different positions in order to boost the segmentation result. In addition, another advantage of our method is using the correlation constrain block, which can constrain the encoders to discover the latent multi-source correlation representation between modalities and then guide the network to learn correlated representation to achieve a better segmentation. From the results, we can observe that with the assistance of correlation constrain block, the network can achieve the best Dice Score of 0.747, 0.886 and 0.776 and Hausdorff Distance of 7.851, 7.345 and 9.016 for enhancing tumor, whole tumor, tumor core, respectively with an average improvement of 2.21% and 9.28% relating to the baseline.

[2]https://ipp.cbica.upenn.edu/

| Methods | Dice Score | | | Hausdorff (mm) | | |
|---|---|---|---|---|---|---|
| | ET | WT | TC | ET | WT | TC |
| (1) | 0.726 | 0.867 | 0.764 | 8.743 | 8.463 | 9.482 |
| (2) | 0.733 | 0.879 | 0.765 | 8.003 | 7.813 | 9.153 |
| (3) | **0.747** | **0.886** | **0.776** | **7.851** | **7.345** | **9.016** |

The results in Table I demonstrate the effectiveness of each proposed component and our proposed network architecture can perform well on brain tumor segmentation.

*2) Comparisons with the State-of-the-art:* To demonstrate the performance of our method, we compare our proposed method with the state-of-the-art methods on Brats 2018 validation set, which contains 66 images of patients without the ground truth. Table II shows the comparison results. We have also carried out a comparison study with the state of art of methods based on U-Net.

(1) Hu et al. [31] proposed the multi-level up-sampling network (MU-Net) for automated segmentation of brain tumors, where a novel global attention (GA) module is used to combine the low level feature maps obtained by the encoder and high level feature maps obtained by the decoder.

(2) Tuan et al. [32] proposed using Bit-plane to generate a series of binary images by determining significant bits. Then, the first U-Net used the significant bits to segment the tumor boundary, and the other U-Net utilized the original images and images with least significant bits to predict the label of all pixel inside the boundary.

(3) Hu et al. [33] introduced the 3D-residual-Unet architecture. The network comprises a context aggregation pathway and a localization pathway, which encoder abstract representation of the input, and then recombines these representations with shallower features to precisely localize the interest domain via a localization path.

(4) Myronenko et al. [19] proposed a 3D MRI brain tumor segmentation using autoencoder regularization, where a variational autoencoder branch is added to reconstruct the input image itself in order to regularize the shared decoder and impose additional constraints on its layers.

The best result in BraTS 2018 Challenge is from [19], which achieves 0.814, 0.904 and 0.859 in terms of Dice Score on enhancing tumor, whole tumor and tumor core regions, respectively. However, it uses 32 initial convolution filters and a lot of memories (NVIDIA Tesla V100 32GB GPU is required) to train the model, which is computationally expensive. While our method used only 8 initial filters, and from Table II, it can be observed that our proposed method can yield a competitive results in terms of Dice Score and Hausdorff distance across all the tumor regions. We also implemented the method [19] with 8 initial filters, but the

| Methods | Dice Score | | | | Hausdorff (mm) | | | |
|---|---|---|---|---|---|---|---|---|
| | ET | WT | TC | Average | ET | WT | TC | Average |
| [31] | 0.69 | 0.88 | 0.74 | 0.77 | 6.69 | 4.76 | 10.67 | 7.373 |
| [32] | 0.682 | 0.818 | 0.699 | 0.733 | 7.016 | 9.412 | 12.462 | 9.633 |
| [33] | 0.719 | 0.856 | 0.769 | 0.781 | 5.5 | 10.843 | 9.985 | 8.776 |
| [19] | **0.814** | **0.904** | **0.859** | **0.859** | **3.804** | **4.483** | **8.2777** | **5.521** |
| Proposed | 0.705 | 0.883 | 0.783 | 0.79 | 7.27 | 5.111 | 10.047 | 7.476 |

results are not good. Compared with other methods, [33] has a better Dice Score on enhancing tumor, while our method achieves a better average Dice Score on all the tumor regions with an improvement of 1.15%, and it can also obtain an average improvement of 14.81% for Hausdorff Distance.

To visualize the effectiveness of proposed correlation constrain block, we select an example to show the feature representation of the four modalities in the last layer (before the output) of the network in Fig. 7. The first and second row show the feature representations without and with correlation constrain block, the fifth column shows the ground truth. We can observe that, the correlation constrain block can constrain the network to emphasize the interested tumor region for segmentation.
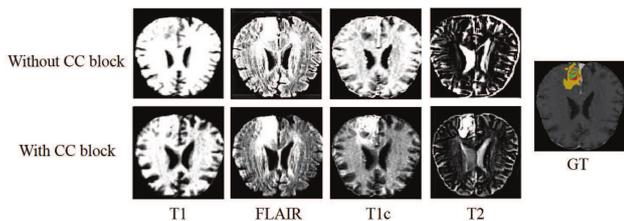
Fig. 7. Visualization of effectiveness of proposed correlation constrain block.

### B. Qualitative Analysis

In order to evaluate the robustness of our model, we randomly select several examples on BraTS 2018 dataset and visualize the segmentation results in Fig. 8. From Fig. 8, we can observe that the segmentation results are gradually improved when the proposed strategies are integrated, these comparisons indicate that the effectiveness of the proposed strategies. In addition, with all the proposed strategies, our proposed method can achieve the best results.

### V. DISCUSSION AND CONCLUSION

In this paper, we proposed a 3D multimodal brain tumor segmentation network guided by a multi-source correlation constrain, where the architecture demonstrated their segmentation performances in multi-modal MR images of glioma patients.

To take advantage of the complimentary information from different modalities, the multi-encoder based network is used
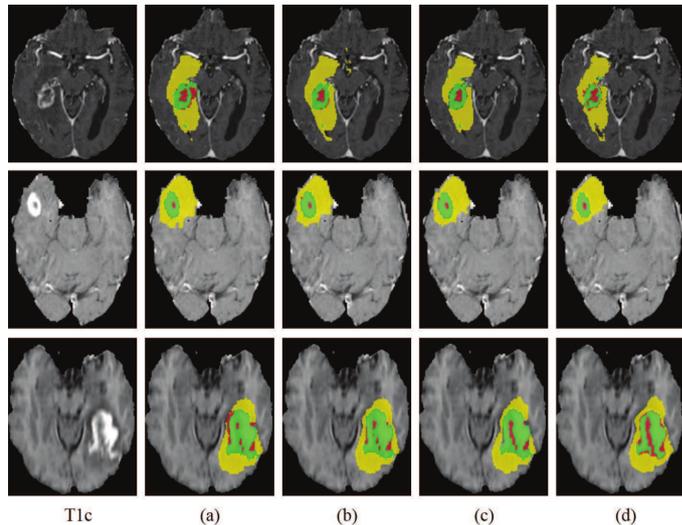
Fig. 8. Visualization of several segmentation results. (a) Baseline (b) Baseline with fusion block (c) Proposed method with fusion block and correlation constrain (d) Ground truth. Red: necrotic and non-enhancing tumor core; Yellow: edema; Green: enhancing tumor.

to learn modality-specific feature representation. Considering the correlation between MR modalities can help the segmentation, a linear correlation block is used to describe the latent multi-source correlation. Since an effective feature learning can contribute to a better segmentation result, a loss function is to guide the network to learn the most correlated feature representation to improve the segmentation. Furthermore, different MR modalities can identify different attributes of the target tumor, and each MR modality image can present different contents at different locations. To this end, inspired by an attention mechanism, a dual-attention fusion strategy is integrated to our network. The modality attention module is used to distinguish the contribution of each modality, and the spatial attention module is used to extract more useful spatial information to boost the segmentation result. The proposed fusion strategy encourages the network to learn more useful feature representation to boost the segmentation result, which is better then the simple max or mean fusion method.

The advantages of our proposed network architecture (i) The segmentation results evaluated on the two metrics (Dice Score and Hausdorff Distance) are similar to real annotation provided by the radiologist experts. (ii) The architecture are an end-to-end Deep Leaning approach and fully automatic without any user interventions. (iii) The experiment results demonstrate that our proposed method gives a very accurate result for the segmentation of brain tumors and its sub-regions even small regions, and it also achieves very competitive results with less computational complexity. In addition, our method can be generalized to other kinds of correlation (e.g. nonlinear) and applied to other kinds of multi-source images if some correlation exists between them.

As a perspective of this research, we will valid our method in different clinical scenarios. In addition, we intend to study a

more efficient correlation representation approach to describe the correlation between modalities, and apply it to synthesize additional images to cope with the limited medical image dataset.

## REFERENCES

[1] B. H. Menze, A. Jakab, S. Bauer, J. Kalpathy-Cramer, K. Farahani, J. Kirby, Y. Burren, N. Porz, J. Slotboom, R. Wiest *et al.*, "The multimodal brain tumor image segmentation benchmark (brats)," *IEEE transactions on medical imaging*, vol. 34, no. 10, pp. 1993–2024, 2014.

[2] Z.-P. Liang and P. C. Lauterbur, *Principles of magnetic resonance imaging: a signal processing perspective*. SPIE Optical Engineering Press, 2000.

[3] S. Bauer, R. Wiest, L.-P. Nolte, and M. Reyes, "A survey of mri-based medical image analysis for brain tumor studies," *Physics in Medicine & Biology*, vol. 58, no. 13, p. R97, 2013.

[4] A. Drevelegas, *Imaging of brain tumors with histological correlations*. Springer Science & Business Media, 2010.

[5] J. Lapuyade-Lahorgue, J.-H. Xue, and S. Ruan, "Segmenting multi-source images using hidden markov fields with copula-based multivariate statistical distributions," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3187–3195, 2017.

[6] N. Zhang, S. Ruan, S. Lebonvallet, Q. Liao, and Y. Zhu, "Kernel feature selection to fuse multi-spectral mri images for brain tumor segmentation," *Computer Vision and Image Understanding*, vol. 115, no. 2, pp. 256–269, 2011.

[7] C. Lian, S. Ruan, T. Denœux, H. Li, and P. Vera, "Joint tumor segmentation in pet-ct images using co-clustering and fusion based on belief functions," *IEEE Transactions on Image Processing*, vol. 28, no. 2, pp. 755–766, 2018.

[8] C. Lian, S. Ruan, and T. Denoeux, "Dissimilarity metric learning in the belief function framework," *IEEE Transactions on Fuzzy Systems*, vol. 24, no. 6, pp. 1555–1564, 2016.

[9] D. Zikic, B. Glocker, E. Konukoglu, A. Criminisi, C. Demiralp, J. Shotton, O. M. Thomas, T. Das, R. Jena, and S. J. Price, "Decision forests for tissue-specific segmentation of high-grade gliomas in multi-channel mr," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2012, pp. 369–376.

[10] Y. Yu, P. Decazes, J. Lapuyade-Lahorgue, I. Gardin, P. Vera, and S. Ruan, "Semi-automatic lymphoma detection and segmentation using fully conditional random fields," *Computerized Medical Imaging and Graphics*, vol. 70, pp. 1–7, 2018.

[11] S. Bauer, L.-P. Nolte, and M. Reyes, "Fully automatic segmentation of brain tumor images using support vector machine classification in combination with hierarchical conditional random field regularization," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2011, pp. 354–361.

[12] S. Cui, L. Mao, J. Jiang, C. Liu, and S. Xiong, "Automatic semantic segmentation of brain gliomas from mri images using a deep cascaded neural network," *Journal of healthcare engineering*, vol. 2018, 2018.

[13] X. Zhao, Y. Wu, G. Song, Z. Li, Y. Zhang, and Y. Fan, "A deep learning model integrating fcnns and crfs for brain tumor segmentation," *Medical image analysis*, vol. 43, pp. 98–111, 2018.

[14] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.

[15] M. Havaei, A. Davy, D. Warde-Farley, A. Biard, A. Courville, Y. Bengio, C. Pal, P.-M. Jodoin, and H. Larochelle, "Brain tumor segmentation with deep neural networks," *Medical image analysis*, vol. 35, pp. 18–31, 2017.

[16] G. Wang, W. Li, S. Ourselin, and T. Vercauteren, "Automatic brain tumor segmentation using cascaded anisotropic convolutional neural networks," in *International MICCAI Brainlesion Workshop*. Springer, 2017, pp. 178–190.

[17] K. Kamnitsas, C. Ledig, V. F. Newcombe, J. P. Simpson, A. D. Kane, D. K. Menon, D. Rueckert, and B. Glocker, "Efficient multi-scale 3d cnn with fully connected crf for accurate brain lesion segmentation," *Medical image analysis*, vol. 36, pp. 61–78, 2017.

[18] K. Kamnitsas, W. Bai, E. Ferrante, S. McDonagh, M. Sinclair, N. Pawlowski, M. Rajchl, M. Lee, B. Kainz, D. Rueckert *et al.*, "Ensembles of multiple models and architectures for robust brain tumour segmentation," in *International MICCAI Brainlesion Workshop*. Springer, 2017, pp. 450–462.

[19] A. Myronenko, "3d mri brain tumor segmentation using autoencoder regularization," in *International MICCAI Brainlesion Workshop*. Springer, 2018, pp. 311–320.

[20] T. Zhou, S. Ruan, and S. Canu, "A review: Deep learning for medical image segmentation using multi-modality fusion," *Array*, vol. 3, p. 100004, 2019.

[21] F. Isensee, P. Kickingereder, W. Wick, M. Bendszus, and K. H. Maier-Hein, "Brain tumor segmentation and radiomics survival prediction: Contribution to the brats 2017 challenge," in *International MICCAI Brainlesion Workshop*. Springer, 2017, pp. 287–297.

[22] K.-L. Tseng, Y.-L. Lin, W. Hsu, and C.-Y. Huang, "Joint sequence learning and cross-modality convolution for 3d biomedical segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6393–6400.

[23] T. Zhou, S. Ruan, Y. Guo, and S. Canu, "A multi-modality fusion network based on attention mechanism for brain tumor segmentation," in *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2020, pp. 377–380.

[24] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.

[25] H. Li, P. Xiong, J. An, and L. Wang, "Pyramid attention network for semantic segmentation," *arXiv preprint arXiv:1805.10180*, 2018.

[26] O. Oktay, J. Schlemper, L. L. Folgoc, M. Lee, M. Heinrich, K. Misawa, K. Mori, S. McDonagh, N. Y. Hammerla, B. Kainz *et al.*, "Attention u-net: Learning where to look for the pancreas," *arXiv preprint arXiv:1804.03999*, 2018.

[27] A. G. Roy, N. Navab, and C. Wachinger, "Concurrent spatial and channel 'squeeze & excitation' in fully convolutional networks," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2018, pp. 421–429.

[28] A. G. Roy, S. Siddiqui, S. Pölsterl, N. Navab, and C. Wachinger, "'squeeze & excite' guided few-shot segmentation of volumetric images," *Medical image analysis*, vol. 59, p. 101587, 2020.

[29] J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, and H. Lu, "Dual attention network for scene segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3146–3154.

[30] B. B. Avants, N. Tustison, and G. Song, "Advanced normalization tools (ants)," *Insight j*, vol. 2, pp. 1–35, 2009.

[31] Y. Hu, X. Liu, X. Wen, C. Niu, and Y. Xia, "Brain tumor segmentation on multimodal mr imaging using multi-level upsampling in decoder," in *International MICCAI Brainlesion Workshop*. Springer, 2018, pp. 168–177.

[32] T. A. Tuan *et al.*, "Brain tumor segmentation using bit-plane and unet," in *International MICCAI Brainlesion Workshop*. Springer, 2018, pp. 466–475.

[33] X. Hu, H. Li, Y. Zhao, C. Dong, B. H. Menze, and M. Piraud, "Hierarchical multi-class segmentation of glioma images using networks with multi-level activation function," in *International MICCAI Brainlesion Workshop*. Springer, 2018, pp. 116–127.