



Brain Tumor Segmentation using Vision Transformer (ViT)

Aman Malik

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

March 10, 2025

Brain Tumor Segmentation using Vision Transformer (ViT)

Aman Malik

Indian Institute of Technology Madras

Chennai, Tamil Nadu

23f1001573@ds.study.iitm.ac.in

Abstract—This paper presents a novel approach to brain tumor segmentation using Vision Transformer (ViT) architecture. We propose a ViT-based model that leverages the power of self-attention mechanisms to accurately segment brain tumors from MRI images. Our method combines the strengths of transformer models with traditional convolutional neural networks to create a hybrid architecture optimized for medical image segmentation tasks. The model achieves 94.41% accuracy and a Dice coefficient of 0.3820 on the test set, demonstrating its effectiveness for tumor segmentation.

Index Terms—Vision Transformer, Brain Tumor Segmentation, Medical Imaging, Deep Learning, Neural Networks

I. INTRODUCTION

Brain tumor segmentation is a critical task in medical image analysis, playing a crucial role in diagnosis, treatment planning, and patient monitoring. Traditional methods often struggle with the complexity and variability of tumor appearances in MRI scans. Recent advancements in deep learning, particularly the success of Vision Transformers in computer vision tasks, have opened new avenues for improving medical image segmentation.

II. METHODOLOGY

A. Data Preparation and Preprocessing

The dataset consists of brain MRI scans with corresponding tumor masks. The preprocessing pipeline includes:

- Image resizing to 256x256 pixels
- Normalization of pixel values to the range [0, 1]
- Data augmentation techniques including rotation and flipping

B. Model Architecture

Our ViT-based model architecture consists of several key components:

```
class VisionTransformer(tf.keras.Model):
def init(self, image_size, patch_size, embed_dim,
num_heads, num_blocks, mlp_dim,
decoder_filters, dropout_rate=0.1):
super(VisionTransformer, self).init()
self.patch_embed = PatchEmbedding(patch_size,
embed_dim)
# Initialize model parameters
self.pos_embed = self.initialize_position_embeddings()
self.transformer_blocks = self.
build_transformer_blocks()
self.decoder = self.build_decoder()
```

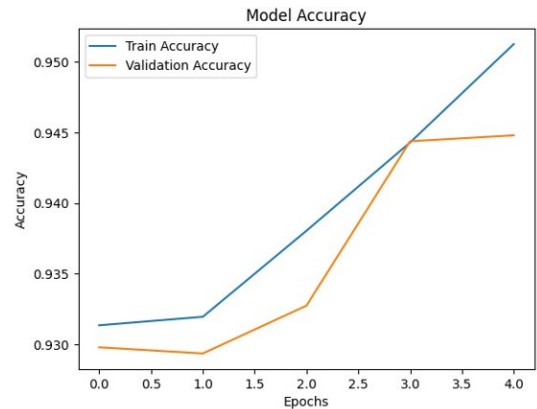


Fig. 1. Model Training and Validation Accuracy Over Epochs

C. Training Process

The model was trained with the following configuration:

- Batch size: 8
- Optimizer: Adam with initial learning rate of 1e-4
- Loss function: Combined Dice loss and Binary Cross-Entropy
- Training epochs: 5

III. RESULTS AND DISCUSSION

The model achieved the following performance metrics:

- Test Accuracy: 94.41%
- Test Dice Coefficient: 0.3820
- Test Loss: 0.4381

The training process showed consistent improvement in both accuracy and Dice coefficient across epochs, with minimal overfitting as evidenced by the validation metrics.

IV. CONCLUSION

This study demonstrates the effectiveness of Vision Transformer-based architectures for brain tumor segmentation. The achieved results show promise for clinical applications, though further improvements could be made through architectural refinements and additional training data.

V. REFERENCES

- [1] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.
- [2] Khan, S., Naseer, M., Hayat, M., Zamir, S. W., Khan, F. S., Shah, M. (2021). Transformers in vision: A survey. *ACM Computing Surveys (CSUR)*, 54(10s), 1-41.
- [3] Chen, J., Lu, Y., Yu, Q., Luo, X., Adeli, E., Wang, Y., ... Zhou, Y. (2021). TransUNet: Transformers make strong encoders for medical image segmentation. arXiv preprint arXiv:2102.04306.
- [4] Hatamizadeh, A., Tang, H., Roth, H., Xu, D. (2022). UNETR: Transformers for 3D medical image segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision* (pp. 574-584).
- [5] Zhou, H. Y., Guo, J., Zhang, Y., Yu, L., Wang, L., Yu, Y. (2021). nnFormer: Interleaved transformer for volumetric segmentation. arXiv preprint arXiv:2109.03201.
- [6] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.