February 2022
Master's Degree Thesis

# A Study on Brain MRI Segmentation using Multi-scale Squeeze U-SegNet with Multi-global Attention

Graduate School of Chosun University

Department of Information and Communication Engineering

Chaitra Dayananda

# A Study on Brain MRI Segmentation using Multi-scale Squeeze U-SegNet with Multi-global Attention

뇌 MRI 영상에서 다중 전역 집중 기반 다중 스케일 압축 U-SegNet을 이용한 뇌영상 분할 방법에 관한 연구

February 25, 2022

Graduate School of Chosun University

Department of Information and Communication Engineering

Chaitra Dayananda

# A Study on Brain MRI Segmentation using Multi-scale Squeeze U-SegNet with Multi-global Attention

Advisor: Prof. Bumshik Lee

This thesis is submitted to Chosun University in partial fulfillment of the requirements for a Master's degree

October 2021

## Graduate School of Chosun University

Department of Information and Communication Engineering

Chaitra Dayananda

This is to certify that the master's thesis of

Chaitra Dayananda

has been approved by the examining committee for
the thesis requirement for the master's degree in
Engineering.

**Committee Chairperson:  Prof.  Jae-Young Pyun**    _____

**Committee Member:        Prof.  Jihwan Moon**       _____

**Committee Member:        Prof. Bumshik Lee**         _____

December 2021

# Graduate School of Chosun University

# Table of Contents

# List of Figures

# List of Tables

# 요약

# 뇌 MRI 영상에서 다중 전역 집중 기반 다중 스케일 압축 U-SegNet을 이용한 뇌영상 분할 방법에 관한 연구

챠이트라 다야난다
지도교수: 이범식
조선대학교 대학원
정보통신공학과

본 논문에서는 자기 공명 이미지 (MRI)에서 뇌 조직을 분할하기 위해 새로운 다중 전역 집중 기반 컨볼루션 학습을 통한 다중 스케일 압축 U-SegNet 아케텍처를 제안한다. CNN (Convolutional Neural Network) 은 의료 영상 분할에서 압도적인 성능을 보여주지만, 기존 CNN 모델에는 몇 가지 단점이 있다. 특히 인코더-디코더 기반 접근법의 사용은 유사한 low-level의 특징을 여러 번 추출하여 정보의 중복 사용을 초래한다. 또한, long-range dependency 비효율적인 모델링으로 인해, 각 의미 클래스는 정확하지 않은 특징 표현과 연관될 가능성이 높아 분할의 정확도가 낮다.

 제안된 전역 집중 모듈은 특징 추출을 조정하고 컨볼루션 신경망의 표현력을 향상시킨다. 또한 집중 기반 다중 스케일 융합 방법은 해당 전역 종속성과 로컬 특징을 통합할 수 있다. 인코더와 디코더 경로 모두에 파이어 모듈을 통합하면 모델 매개변수의 수가 크게 감소하기 때문에 계산 복잡성을 줄일 수 있다. 제안된 방법은 뇌 조직 분할을 위해 공개 데이터 세트에서 평가되었다. 실험 결과에 따르면 제안된 모델은 학습 가능한 매개 변수의 수가 크게 감소하여 뇌척수액(CSF)의 경우 94.81%, 회백질(GM)의 경우 95.54%, 백질(WM)의 경우 96.33%의 분할 정확도를 달성한다. 본 연구는 이전에 개발된 U-

SegNet 기반 분할 접근법에 비해 학습 가능한 매개 변수의 수가 4.5배 감소하면서 주사위 유사성 지수 측면에서 예측 정확도를 2.5% 향상시켜 더 나은 분할 성능을 보여준다. 이는 제안된 접근 방식이 뇌 MRI 영상의 신뢰성 있고 정확한 자동 분할을 달성할 수 있음을 보여준다.

키워드: CNN, 조직 분할, 다중 글로벌 어텐션, 뇌 MRI.

# ABSTRACT

# A Study on Brain MRI Segmentation using Multi-scale Squeeze U-SegNet with Multi-global Attention

Chaitra Dayananda
Advisor: Prof. Bumshik Lee
Department of Information and
Communication Engineering
Graduate School
Chosun University

This research work focuses on the multi-scale feature extraction with novel attention-based convolutional learning using the U-SegNet architecture to achieve segmentation of brain tissue from a magnetic resonance image (MRI). Although convolutional neural networks (CNNs) show enormous growth in medical image segmentation, there are some drawbacks with the conventional CNN models. In particular, the conventional use of encoder-decoder approaches leads to the extraction of similar low-level features multiple times, causing redundant use of information. Moreover, due to inefficient modeling of long-range dependencies, each semantic class is likely to be associated with non-accurate discriminative feature representations, resulting in low accuracy of segmentation. The proposed global attention module refines the feature extraction and improves the representational power of the convolutional neural network. Moreover, the attention-based multi-scale fusion strategy can integrate local features with their corresponding global dependencies. The

integration of fire modules in both the encoder and decoder paths can significantly reduce the computational complexity owing to fewer model parameters. The proposed method was evaluated on publicly accessible datasets for brain tissue segmentation. The experimental results show that our proposed model achieves segmentation accuracies of 94.81% for cerebrospinal fluid (CSF), 95.54% for gray matter (GM), and 96.33% for white matter (WM) with a noticeably reduced number of learnable parameters. Our study exhibits higher segmentation performance, increasing prediction accuracy by 2.5% in terms of dice similarity index while reducing the number of learnable parameters by 4.5 times in comparison to previously established U-SegNet based segmentation algorithms. This illustrates that the proposed approach can achieve reliable and precise automatic segmentation of brain MRI images.

**Keywords**: CNN, tissue segmentation, multi-global attention, brain MRI.

# 1. INTRODUCTION

The rapid progression of medical image processing techniques has benefited mankind and plays an important role in clinical diagnosis. The current advances in medical imaging help to view the human body in to diagnose and monitor medical conditions [1-2]. The imaging techniques such as ultrasound (US), magnetic resonance imaging (MRI), and X-ray imaging give image information by which the radiologist has to analyze and evaluate comprehensively in a shorter time [2]. In particular, magnetic resonance imaging (MRI) is typically favored for structural analysis as it generates images with high soft-tissue contrast and higher spatial resolution and does not entail any health hazards. Brain MRI scans are quantitatively examined to diagnose various brain disorders such as epilepsy, schizophrenia, Alzheimer's disease, and other degenerative disorders [3]. MRI is also essential to identify and localize abnormal tissues and healthy structures for diagnosis and postoperative analysis. Hence, the segmentation of these abnormal tissues from the medical images plays a vital role in the study and treatment of many diseases. Early detection of brain disorders allows the observer to follow up on the subject. Hence, the main objective is to derive better tools that help to interpret the images.

In this chapter, Section 1.1 presents a brief description of magnetic resonance imaging (MRI) in the diagnosis of various brain disorders. Section 1.2 presents the overview and motivation of the proposed work. Section 1.3 presents the research objectives and major contributions of the thesis. Section 1.4 explains the outline of the thesis.

## 1.1. A Brief Review of MRI

Magnetic resonance imaging (MRI) is a medical imaging technology that creates detailed images of the organs and tissues using a magnetic field and computer-generated radio waves [2]. Large, tube-shaped magnets are used in the majority of MRI equipment. The magnetic field momentarily realigns water molecules in the body while we lie inside an MRI machine. These aligned atoms emit tiny signals, which are used to form cross-sectional MRI pictures, similar to slices in a loaf of bread. The MRI scanner can also create three-dimensional images that may be viewed from various perspectives. MRI is a non-invasive method of examining the organs, tissues, and skeletal system by a doctor. It creates high-resolution images of the inside of the body to aid in the diagnosis of a wide range of ailments. The most common imaging test for the brain is MRI. The functional MRI of the brain (fMRI) is a unique type of MRI [3]. It generates images of blood flow to specific brain locations. It may be used to look at the structure of the brain and figure out which areas of the brain are in charge of essential functions. Damage from a head injury or illnesses like Alzheimer's disease can also be assessed using functional MRI.

## 1.2. MRI for Diagnosis of Brain Disorders

Brain magnetic resonance imaging (MRI) is a non-invasive, painless technique that gives detailed images of the brain and brain stem. The importance of MRI in brain diagnosis is well recognized, and it has been included in various new brain diagnostic criteria [4-5]. Cysts, tumors, hemorrhage, swelling, developmental and anatomical abnormalities, infections, inflammatory disorders, and blood vessel problems can all be detected using an MRI scan. High-resolution MRI may identify the existence and severity of brain atrophy, which can aid in the diagnosis of Alzheimer's

disease in vivo [6] and even show the presence of neurofibrillary tangles (NFTs), which are regarded as a hallmark pathology of the disease [7]. While MRI-measured brain atrophy is a reliable and sensitive indicator of neurodegeneration in general [8], it can also be used to identify many different types of dementia that have similar patterns of atrophy [9].

## 1.3. Overview and Motivation

The brain and nerves degenerate over time as a result of neurodegenerative illnesses. These diseases have the potential to alter personality and cause confusion. They can also damage the cells and nerves in the brain. Alzheimer's disease, for example, is a brain disorder that can develop as one becomes older [10]. They can wreak havoc on the memory system and mental processes over time. Other diseases, such as Tay-Sachs disease, are inherited and manifest themselves at a young age. The following are some other common neurodegenerative diseases:

- Alzheimer's disease
- ALS (amyotrophic lateral sclerosis)
- Parkinson's disease
- all forms of dementia

There are no treatments that can stop or reverse the disease's course, while some can temporarily alleviate symptoms. Prevention is the key to reducing the occurrence of neurological illnesses and consequently the number of global deaths in today's lifestyle. As a result, an early indicator of a higher risk of neurodegenerative disorder is a good predictor of clinical diagnosis [10]. Furthermore, neurological brain diseases are diagnosed based on the attributes such as medical test results, clinical history, and medical image acquisition. In

3

some circumstances, an expert's ability to interpret a huge number of data in a short length of time makes diagnosis more challenging [11]. The progress of medical imaging technology and computer software is being utilized to assist specialists in quickly identifying and interpreting diseases. These programs provide a diagnosis of the condition based solely on visual data [12]. Developing software for radiological image processing is, in fact, one of the most challenging tasks in the medical field. Modern diagnostic systems are created using cutting-edge computing and data processing technology since accurate illness diagnosis is dependent on both image acquisition and interpretation. Despite commercially available computer-based diagnostic systems, fully automated procedures have yet to be fully established in the literature [10-12]. As a result of this weakness, they are difficult to employ for diagnostic purposes. Brain tissue loss is the first sign of a possible onset of brain illnesses.



| Ground truth | Background | CSF | GM | WM |

Figure 1-1. The brain MRI ground truth image and its segmented structures.

MR imaging, for example, is a non-invasive and reliable approach for measuring brain tissue for the diagnosis of brain illnesses. As a result, the expert is concerned with retrieving useful information on brain tissue segmentation. This task necessitates the use of highly qualified experts. Figure 1-1 shows the brain MRI ground truth image and its segmented structures. Manual tracing of brain structures such as grey matter (GM), white matter

(WM), and cerebrospinal fluid (CSF), on the other hand, produces inconsistent results. Motivated by this limitation, we propose a CNN-based automatic segmentation of brain structures, e.g., white matter (WM), gray matter (GM), and cerebrospinal fluid (CSF) from the brain MRI, which aids in the diagnosis and effective treatment of many brain diseases.

## 1.4. Objectives

The main goal of this work is to develop an efficient CNN algorithm for the segmentation of brain tissue (GM, WM, and CSF) from brain MRI.

The contribution of this work is in the segmentation of brain tissues such as gray matter (GM), white matter (WM), and cerebrospinal fluid (CSF) from the MRI using proposed multi-scale, attention-based convolutional network.

Our main contributions are summarized as follows:

- The proposed modified U-SegNet architecture is integrated with a novel global attention module (GAM). Attention is applied at both contracting and expansive paths, creating a multi-attention network. The key element in GAM is global average pooling, which provides the global context of high-level features as assistance to low-level features to obtain class-category localization.
- The proposed multi-scale input feature fusion strategy extracts the context information of high-level features at different scales, thus incorporating neighbor-scale feature information more precisely.
- The fire modules are used to replace the convolution layer to significantly reduce the number of model parameters, which results in a reduction in the model size and computation complexity; this consequently leads to a more efficient segmentation model.

The obtained outcome is supportive and reliable towards using the deep learning method and can be intended towards medical CAD system development.

## 1.5. Thesis Layout

This thesis is composed of five consecutive chapters. Following the introduction, Chapter 2 presents the survey describing various techniques for sematic segmentation of medical and natural images. Chapter 3 explains the proposed methodology and pipeline, particularly the multi-scale, attention-based features extraction with reduced model parameters. Chapter 4 discusses the experimental results in comparison with the results of the state-of-the-art algorithms. Finally, chapter 5 concludes with a brief discussion on the findings of this work.

# 2. RELATED WORKS

Deep learning is a popular branch of artificial intelligence that involves the creation of algorithms and the automatic learning of patterns through experience. Deep learning algorithms learn to make intelligent decisions based on their capacity to detect complicated patterns, which can be used in applications such as handwriting recognition, stock market analysis, image analysis, and medical diagnosis. This thesis is mostly concerned with medical diagnosis. This is essentially a segmentation task to identify the brain tissues GM, WM, and CSF that would be used in the diagnosis of brain disease such as Alzheimer's, Parkinson's, dementia, etc., using neuroimaging data.

In this chapter, we present the basic overview of deep learning technology relevant to this thesis and suitable for image-based segmentation tasks. Section 2.1 firstly present the details of various segmentation techniques, followed by the description of the algorithm along with their performance assessment. The attention-based learning for image segmentation is presented in section 2.2. Finally, literature work related to model parameter reduction is provided in Section 2.3.

## 2.1. Medical Image Segmentation

Extensive research has been conducted on medical image segmentation in the past ref. [13–15], with CNNs growing rapidly this area, driving exceptional performances in many diverse applications. Conventional CNN architectures, including FCNN [13] or U-net [14], serve as sources of inspiration for existing medical image segmentation methods. The conventional FCNN-based classification network replaces the fully connected layers with convolutional layers to predict the output dense pixels. The input image is recovered to its

original resolution by up-sampling the predictions in a single step. In addition, skip connections [14] are used in the network using intermediate function maps to boost the prediction capabilities. On the other hand, the U-net architecture consists of encoding and decoding paths with a sequence of convolutional layers with pooling and up-sampling. The features from the encoder are concatenated with the decoder layers using skip connections. Several extended U-net and FCNN models have been developed to resolve the problems associated with pixel-wise segmentation across different applications [16–19]. In [17], a patch-wise 3D U-net was proposed for brain tissue segmentation with encoding and decoding layers with randomly sampled and overlapped 3D patches ($8 \times 24 \times 24$) used for training. Unlike the U-net, a convolution operation is introduced as a transition layer between the encoder and decoder layers to give more weight to the higher-level features learned through deeper layers in the network. Pawel *et al*. [18] proposed a 3D-CNN for brain tumor segmentation, where the model was trained on 3D random patches, and features extracted by 2D-CNNs were given as an extra input to a 3D-CNN. The combination of both 3D and 2D features captures rich feature representations from a long-range 2D context in three orthogonal directions. An ensemble of 3D U-nets designed with different hyperparameters uses non-uniformly extracted patches as inputs to obtain brain tumor segmentation [18]. Badrinarayanan *et al*. [20] introduced the SegNet model, which uses pooling indices from the encoder to the up-sampling layers. Hence, it requires very few parameters and is faster to train. Looking into the complementary strengths of SegNet and U-net models, a new hybrid model, is explored, namely U-SegNet [21]. The U-SegNet incorporates the unique architectural features from both U-net and SegNet models and uses SegNet as the base architecture with a skip connection introduced between the encoder

and decoder, providing multi-scale information for better performance. Owing to the pooling indices passed at the decoder side, the U-SegNet model has faster convergence. Recent efforts to promote the discriminative capability of feature representations include a multi-scale fusion strategy [22]. Zhou *et al*. [23] redesigned the skip connections in U-net++ [23] by enabling flexible feature fusion in decoders, thus resulting in an improvement over the restrictive skip connections in U-net [4] that require fusion of only same-scale feature maps. A small drawback in U-net++ is that the number of parameters increases owing to the employment of dense connections [24, 25]. Deep supervision is also employed to balance the decline in segmentation accuracy caused by pruning [26]. Zhao *et al*. [27] proposed a pyramid network that utilizes global learning at various scales by grouping feature maps produced by multiple dilated convolution blocks. The collection of contextual multi-scale information can also be obtained by performing pooling operations [28]. Cheng *et al*. [29] proposed a context encoder network (CE-net) that adopts a pre-trained ResNet block in the feature encoder. CE-net involves a newly proposed dense atrous convolution block, and residual multi-kernel pooling is integrated into the ResNet-modified U-net structure to capture more high-level features and preserve more spatial information. Although these approaches assist in capturing targets at different measures, contextual dependencies for all image regions are uniform and non-adaptive. Hence, these approaches neglect the contrast between local and contextual representations for different categories. Moreover, these multi-context representations were manually composed and lacked the flexibility to form multi-context representations. Further, the ignored contrast information causes long-range object associations in the entire image to be leveraged in these strategies, which is of focal interest in many medical image segmentation problems.

## 2.2. Attention Based Learning

The attention mechanisms highlight key local regions in the feature maps and discard unrelated data carried by the generated feature maps. The attention modules act as crucial components of a network that wants to gather global information. The inclusion of the attention blocks demonstrated very successful outcomes in various vision problems, such as image classification [30], image captioning [31], or image question-answering [32]. Recently, many researchers have shown interest in self-attention, as they offer a greater opportunity to model long-range dependencies while retaining computational and statistical performance [33–36]. Zhao *et al*. introduced a point-wise spatial attention network, where each position on the feature map is connected to all the other feature maps through a self-adaptively learned attention mask [37]. Dong *et al*. [38] proposed attention gates (AGs) and used them for the segmentation of the pancreas. The AGs highlight the salient features while suppressing the irrelevant region from the raw input pixel. AGs make better use of intermediate characteristics, reducing the need for cascaded models [39]. To incorporate local and global-dependent features, Wang *et* al. [40] employed a basic focus module with three convolutional layers. A similar attention method with two convolutional layers integrated with a U-net architecture was proposed in [41]. For better extraction of relevant features, focus gate modules are integrated with the skip connection in the decoding path of the U-net in [39]. Schlemper *et al*. [41] proposed attention modules where the local deep attention features are fused with the global context at multiple resolutions for prostate segmentation on ultrasound images. The multi-scale self-guided attention-based approach can integrate local features with their respective global dependencies, as well as highlight interdependent

channel maps in an adaptive manner to achieve accurate segmentation of medical images [42].

## 2.3. Parameters Regularization

Most deep learning-related studies have considered increasing network accuracy as their primary objective. However, the computational burden of a significant number of parameters and deep architecture becomes a crucial issue. Recent studies have shown that most deep neural networks are over-parametrized, resulting in redundancy in deep learning, leading to inefficient use of memory and computing resources. In these large parameter spaces, various compression techniques, such as shrinking, factorizing, or compressing pre-trained networks, are applied to minimize redundancy and obtain smaller models [43–46]. In [44], singular value decomposition (SVD) was used for a pre-trained CNN architecture to achieve lower-order parameter estimates for model compression. Network pruning methods [43, 45] have been widely studied to achieve compressed CNN models. The parameters of the pre-trained model below a certain threshold are replaced with zeros in the network pruning method to produce sparse matrices. Most of the previous works [45, 46] introduced network-pruning-based methods to decrease the network complexity and reduce the overfitting of the model. Network quantization is proposed to decrease the data dynamic range from 32 to 8 or 16 bits, which further compresses the pruned network by reducing the number of bits required to represent each weight [47]. To efficiently operate on compressed deep learning models, Son *et al*. [48] proposed an efficient inference engine (EIE), a specialized accelerator that accomplishes customized sparse matrix-vector multiplication and performs weight sharing without efficiency loss. To reduce the number of parameters and computational work

in the CNN, various methods based on factorizing the convolution kernel have been introduced [49]. The depth-wise separable convolutions used in SqueezeNet [13, 50] are a form of factorizing convolution that separates convolution across channels rather than convolution within channels. As in the MobileNetV1 architecture, the profoundly separable convolution networks realized with quantization require a special attention module [51]. The special hardware for CNNs has been considered by many methods aimed at minimizing computation time [50].

# 3. PROPOSED METHOD

It is known that the SegNet [20] and U-net [4] are the most widely used deep learning models for image segmentation [21, 25]. In SegNet, pooling indices in the up-sampling process are utilized to compensate for the missing spatial information and lead to faster convergence of the model [20]. U-net uses skip connections from the encoder to the decoder and shows a better segmentation performance [4]. However, multi-stage cascaded CNN approaches are more suitable because the performance of a single SegNet and U-net-based segmentation method is not sufficiently accurate when there are variations in the structure and intensity of the target tissues [25]. However, multiple cascaded networks can produce a significantly large number of model parameters, thus leading to the redundant use of computational resources. To overcome this problem, U-SegNet [21] uses both skip connections and pooling indices to combine feature maps from the encoder to the decoder and localize these feature up-sampling, respectively. As a result, pooling indices make the U-SegNet converge faster, and the skip connection improves the segmentation accuracy. Although U-SegNet shows better segmentation performance, the segmented outputs are still blurry, and the network is insensitive to the fine details of the image [24]. To achieve a better segmentation of brain tissues when training on a limited set, the network needs to extract more discriminative features [38]. However, U-SegNet is slightly insufficient to capture better features because numerous pooling operations in the U-SegNet model produce low-resolution feature maps. Due to the inherent complexity, a large number of layers, and the massive amounts of data required, deep learning models are very slow to train and require a lot of computational power, which makes them very time- and resource-intensive. The model which

can provide improved segmentation results while training on limited or less data is considered to be a high potential network [52]. Meanwhile, due to data scarcity, the need to develop a model which could be trained efficiently on less data is very crucial [53]. Motivated by this problem, we propose a novel global attention mechanism using a U-SegNet architecture, where the proposed architecture is designed with a multi-scale guided multi-global attention module. The multi-scale input features at each encoding layer encode both the global and local contexts. Moreover, the proposed novel global attention at the encoder and decoder can filter irrelevant information and focus on the most relevant details needed for the MRI segmentation task. Besides, the model is prone to lose local details when complete image information is employed as an input to train the network. We also propose the use of a patch-wise splitting of each input slice to resolve this problem, which is used to train the model and provide better segmentation accuracy. Finally, we adopt fire modules that comprise a squeeze layer consisting of only $1 \times 1$ convolution filters followed by an expansion layer with a combination of $1 \times 1$ and $3 \times 3$ filters, which reduces the number of learnable parameters and computational requirements, and results in a smaller efficient model.

## 3.1. Proposed Pipeline

First, the MRI datasets with the corresponding ground truth are prepared. For each MRI scan with the dimension of height $\times$ width $\times$ slices (H $\times$ W $\times$ S), we pad zeros to the H $\times$ W of each slice and resize to a dimension of $256 \times 256$. Then, 48 slices are extracted starting from the 10-th slice with an interval of 3 slices [54]. Furthermore, each slice is divided into four uniform non-overlapping patches, and these patches are given as input to the proposed model for training. Figure 3-1 shows the overall framework of the proposed

14

Figure 3-1. Overall framework of the proposed algorithm.

method. The architecture of the proposed method is discussed in detail as follows: (i) encoder path, (ii) decoder path, and (iii) global attention module (GAM), (iv) fire module, (v) uniform patch-wise input, and (vi) classification layer. As shown in Fig. 3-2, the features extracted using $1 \times 1$ and $3 \times 3$ filters are fused to form a multi-scale input representation. These multi-scale data with the feature maps from a previous network layer are provided as input to the GAM at the encoder side. The GAM at the decoder can capture discriminative information and concentrate on relevant features while performing up-sampling operations. Thus, each network layer contains two separate attention modules which concentrate on extracting enhanced representations of features and generate an accurate segmentation network. In addition, the convolution blocks at the encoder and decoder layers are replaced with fire modules which reduce the number of model parameters and create a smaller network.

## 3.2. Encoder Path

Figure 3-2 shows the architecture of the proposed method with the encoder and decoder paths. The fire module replaces the convolution operation in the proposed method, significantly reducing the number of learnable parameters and computational complexity. The fire module was originally used for

SqueezeNet [13] to reduce the complexity of AlexNet [50]. In this study, we incorporate it with our proposed multi-scale U-SegNet architecture for segmentation. Let us consider $x_l$ as the input sample, where $l$ represents the index of the network layer. The convolution output for the squeeze block is computed as (1).

$$o_{squeeze(l)} = f(x_l * w_l^{1 \times 1} + b_l),  \tag{1}$$

where $o_{squeeze(l)}$ is the squeeze layer output of the fire module and $w_l^{1 \times 1}$ is the kernel weight, where the subscript $[1 \times 1]$ represents the size of the convolution kernel associated with the respective layer and $b_l$ is used as a bias term. * represents the convolution operation. The convolution output is fed to the standard Rectified Linear Unit (ReLU) activation function $f(\cdot)$.

The squeeze output is fed into the expanding module. The expanding module consists of two parallel convolutions with kernel sizes of $3 \times 3$ and $1 \times 1$. Furthermore, the output from these parallel convolutions is concatenated to form the fire module output and is expressed as (2).

$$o_{expand(l)} = \text{Concat}\left[f(o_{squeeze(l)} * w_l^{1 \times 1} + b_l), f(o_{squeeze(l)} * w_l^{3 \times 3} + b_l)\right],  \tag{2}$$

where $o_{expand(l)}$ is the fire module output of the $l$-th network layer, and Concat() is a concatenate function. As shown in Fig.3-2, the encoder path consists of a sequence of fire modules whose output is applied to the GAM as input. The GAM also receives input $ms_l$, obtained from multi-scale input feature fusion. In the multi-scale layer, the input $x_l$ is down-sampled using max-pooling with a stride of $2 \times 2$, as in (3). The max-pooled input is followed by a convolution of $1 \times 1$ and $3 \times 3$ filters separately. These convolved outputs

Figure 3-2. Overview of our multi-scale squeeze U-SegNet with multi global attention for brain MRI segmentation. Solid blue boxes represent the convolution block with dimensions as kernel width×height×filters followed by ReLU.

are concatenated to form multi-scale feature maps, as shown in (3). The multi-scale information and the fire module output and are fed as input to the GAM at each encoding layer, as in (4).

$$ms_l = Concat\ [f(m_l, w_l^{1 \times 1}), f(m_l, w_l^{3 \times 3})]\ . \tag{3}$$

$$GAM_l = GAM(ms_l, o_{expand(l)})\ . \tag{4}$$

The output from GAM is given to the max-pooling layer to reduce the dimensionality and focus on the fine details of the feature map, as expressed in (5). Pooling indices are stored at each encoder layer so that the decoder uses the information to up-sample the feature maps. The output at each encoder layer is referred to as the encoding unit (down-sampling unit) and is obtained using (5).

$$encoder(l)\ =\ Maxpool(GAM_l, 2)\ . \tag{5}$$

## 3.3. Decoder Path

Similar to the encoding path, the decoder path in the proposed method uses transposed fire modules to reduce the number of model parameters. The main component of the decoder path is the up-sampling unit. Each up-sampling unit consists of a transposed fire module. The transposed fire module consists of a 1×1 transposed convolution. The output from the 1×1 transposed convolution is fed into two 3×3, and 1×1 kernel-sized parallel transposed convolutions that are concatenated to form the output transpose fire module, as in the down-sampling unit. The decoder is integrated with attention gates, which can highlight the salient features. The feature maps extracted from the $l$-th (high-level) and $(l-1)$-th (low-level) encoding layers are used as the input signal and gating signal to the attention module, respectively. Thus, the feature map

18

obtained from encoding layers containing contextual information is computed using the GAM to eliminate unrelated responses. The GAM output is concatenated with the feature map of the corresponding up-sampling layer, as expressed in (6). Hence, the attention-based skip connections in the proposed architecture help in extracting the most informative data from the encoder, utilized by the decoder to make more accurate predictions. These skip connections use both high-and low-resolution feature information and focus on the most relevant information while performing up-pooling operations.

$$decoder\ (l)\ =\ Concat[\ GAM(x_l, x_{l-1}), Unpool\ (x_{l-1}, Pool_{idx(l-1)})]\ . \qquad (6)$$

The output from each decoder layer can be obtained using (6), where $Pool_{idx(l-1)}$ is the pooling index passed from the encoder layer to the decoder layer to recover spatial information of the feature maps while performing un-pooling operations at the decoder.

## 3.4. Global Attention Module (GAM)

A distinctive brain signal processing system for human vision is the visual attention process. This is a tool for a human to pick relevant information instantly from a large amount of information using sources of limited attention. In deep learning, the attention process is similar to that of human visual attention. Its main objective is to determine the most relevant data from a vast amount of knowledge to accomplish its goal (tissue segmentation). The attention mechanism improves network performance by suppressing function activations that are not important to the task. To do this, we propose a novel global attention module with self-attention in an efficient manner. As a guide from low-level features to assess class localization, the GAM on each encoder and decoder layer enables global context details.

Figure 3-3. Schematic of the proposed global attention module (GAM) integrated into the proposed brain segmentation architecture.

Figure 3-3 shows the proposed architecture of GAM, which is integrated with our proposed brain segmentation architecture. The $x_l$ is the output feature map from the $l$-th encoding layer (low-level features). The $x_{l+1}$ collected from a coarser scale serves as a vector of the gating signal and is applied to select the target area for each pixel. The $\alpha_l$ is the tensor coefficient that preserves activation by suppressing irrelevant feature responses associated with the target task. The operation of GAM is the element-wise addition of the feature map with the tensor coefficient from the $l$-th encoding layer, and the output of GAM is obtained using (7).

$$x_{lout} = \alpha_l + x_l .\tag{7}$$

In the case of multiple semantic groups, learning multidimensional coefficients of attention is suggested. As guidance for low-level features to incorporate local features into the global context, global average pooling provides global context information. The global information produced from the high-level feature is fed to a 1×1 convolution with the ReLU activation function. To obtain weighted low-level characteristics, it is multiplied by 1×1 convoluted low-level features. To obtain the tensor coefficient of attention, we used multiplicative attention. To extract pixel localization specific to the class

20

of the high-level feature index, the tensor coefficient is up-sampled and added with low-level features. The tensor coefficient of attention is obtained as follows:

$$\alpha_l = upsample\{(GAP[x_l] * W_x + b), (x_{l+1} * W_g + b)\},\qquad(8)$$

where $W_x$ and $W_g$ are the weight values associated with the input and gating signals, respectively, b is the bias term, and GAP is a function of the global average pooling. The feature maps obtained from the attention module $x_{lout}$, which contains contextual information, were concatenated with the feature map of the corresponding decoding layer forming skip connections. These skip connections use both high- and low-resolution features; they focus on the most relevant information while performing up-sampling operations.

## 3.5. Fire Module

The fire module was initially introduced in SqueezeNet [13] to identify CNN architectures with fewer parameters while maintaining competitive accuracy. The fire module in SqueezeNet is composed of a squeeze convolution layer with only 1×1 filters, feeding into an expand layer with both 1×1 and 3×3 convolution filters. The number of 1×1 filters in the squeeze layer is set to be less than the total number of 3×3 and 1×1 filters in the expand layer, so the squeeze layer helps in limiting the number of input channels to the 3×3 filters in the expand layers. Owing to the benefits of the fire module in reducing the learnable parameters, we used the same design of the fire module and integrated it into our proposed encoder and decoder architecture. Figure 3-4 shows a schematic representation of the fire module applied to the proposed architecture. Figure 3-4(a) and -(c) show the encoder and decoder sides of U-net [4] using a normal convolution layer, with each convolution block

21

containing $F_{in}$ filters, and takes a feature map of height×width×channels (H×W×C) as input. Figure 3-4(b) and - (d) show the architecture of the fire modules at the encoder and decoder paths in the proposed method. Likewise, the fire module of SqueezeNet, the fire module in the proposed method consists of two parts: (i) the squeeze layer and (ii) the expand layer. As shown in Fig. 3-4(b), the squeeze module consists of one convolution layer with a kernel size of 1×1 and an output channel equal to $F_{in}/4$, where $F_{in}$ is the number of convolution filters in the conventional U-net [4]. The squeeze output is fed into the expanding module. The expanding module consists of two parallel convolutions with kernel sizes of 3×3 and 1×1, each convolution with $F_{in}/2$ output channels.



Figure 3-4. Schematic of fire module. (a) & (c) show the convolution for the encoder and decoder side in U-net [4] respectively; (b) & (d) show our corresponding squeeze U-SegNet at encoder and decoder side using squeeze and expand layers to reduce the number of parameters.

Furthermore, the output from these parallel convolutions is concatenated to form a fire module output, where $F_{out} = F_{in}$. Hence, as mentioned in [13], the proposed method maintains the number of filters in the squeeze layer, which

is less than the total number of filters in the expand layers, which results in a significant reduction in the total number of network parameters.

## 3.6. Uniform Input Patch

The brain MRI scan of each subject constitutes the dimensions H×W×S. Some of the starting and ending slices of the brain MRI scan do not provide much useful information, as analyzed in the previous research [54], and the consecutive slices would share almost the same information. Hence, to exclude these non-informative slices and reduce the multiple training of consecutive slices, we selected 48 slices with a gap of 3 slices, which ensured the presence of slices with more as well as less information for model training. Each of the extracted slices was resized to 256×256. The partitioning of a slice with individual patches improves localization because the trained network can better concentrate on local details in a patch. Therefore, each slice was divided into four uniform patches using our proposed method. Therefore, the dimensions of each partitioned patch were 128×128 in the proposed method. These patches are fed into the training of the model and the predicted segmentation results are obtained for the test data.

## 3.7. Classification Layer

The final decoder layer consists of a 1×1 convolutional layer with softmax activation to predict a reconstructed segmentation map. The output contains four predicted classes: gray matter (GM), white matter (WM), cerebrospinal fluid (CSF), and background. The proposed model accepts the input image and produces the corresponding learned representation. Based on this feature representation, the input image is classified into any of the four output classes. The cross-entropy loss is used to measure the proposed model losses, as in

(10). The softmax layer learns representations from the decoder and interprets them in the output class. The probability score $y'$ is assigned to the output class. If we define the number of output classes as c, we obtain as (9) as follows:

$$y' = \frac{exp^{decoder(l)}}{\sum_{j=1}^{c} exp^{decoder(l)_j}} \; , \tag{9}$$

and the cross-entropy loss function is used to calculate the network cost function as in (11):

$$L(y, y') = \sum_{i=1}^{c} y_i \, log(y_i'), \tag{10}$$

where for each class of $i$, the ground truth and predicted distribution score are $y$ and $y'$, respectively.

# 4. EXPERIMENTAL ANALYSIS

## 4.1. Datasets

The proposed method was evaluated using two sets of brain MRI images. (i) OASIS dataset, (ii) IBSR dataset. A detailed description of the dataset is provided as follows,

## 4.2. OASIS Dataset

The first sample included 416 T1 weighted brain MRI scans from the Open Access Series of Imaging Studies (OASIS) database [55], where information from both non-demented and demented subjects was obtained from Washington University. A T1-weighted (T1W) image is a common MR imaging pulse sequence that shows signal differences based on the intrinsic T1 relaxation time of distinct tissues. Clinically, T1-weighted images are superior at displaying normal anatomy and are mainly used for the anatomical details and pathological abnormalities of the intracranial lesions [56]. Of the 416 subjects in total, 150 were chosen for our experiments. The first 120 subjects were used for model training out of the selected data, and the remaining 30 subjects were used as test datasets. For our studies, the axial, sagittal, and coronal planes of the MRI slices were used for training and testing the proposed network. In the OASIS dataset, the size of each input axial scan was $208\times176\times176$ which corresponds to height×width×slices respectively and each scan consisted of 176 slices. It was observed that the distinguishable tissue regions were mostly found near the middle slices of the volume [54]. Often, the same information is exchanged for consecutive slices. Therefore, to remove these non-informative slices and decrease the repetitive training of consecutive slices, a sample of 48 slices, starting from the 10-th slice with an

interval of three slices, were selected for the evaluation procedure. By inserting 24 pixels of zeros at the top and bottom of the image and 40 pixels of zeros on the left and right of the image, the extracted slices were resized to the dimensions of 256×256×48. Similarly, the sagittal and coronal planes of MRI slices were also resized to 256×256 dimensions. Each input scan, therefore, consisted of 48 slices with dimensions of 256×256. During the training phase, slices of each MRI scan and their corresponding ground-truth segmentation maps were split into uniform patches. An input slice had dimensions of 256×256 and each slice was split into four patches. Therefore, in the proposed model, the dimensions of each partitioned patch were 128×128. These patches were given as input to the training model, and the predicted segmentation results were obtained for the test data.

## 4.3. IBSR Dataset

The second dataset contains MRI from the Internet Brain Segmentation Repository (IBSR) [57] dataset. The IBSR dataset includes 18 T1-weighted MRIs of 14 healthy men and 4 healthy women between 7 and 71 years of age. Pre-processing, such as skull stripping, normalization, and bias field correction, is performed on the MRIs in the IBSR. The training dataset included 12 subjects with manually annotated and confirmed ground truth labels for the remaining six subjects for testing the model. The original axial scans (256×128×256) were padded to the top and bottom of the image with 64 pixels of zeros to resize to a dimension of 256×256×256 to use the patches effectively in our proposed model. Similarly, the original coronal (256×256×128) and sagittal (128×256×256) scans were also resized to dimensions of 256×256×256 for the experiments. Table 1 summarizes the OASIS and IBSR datasets used in the experiments.

Table 1. Summary of OASIS and IBSR datasets used in our experiments.

| Dataset | No. of subjects | | | Experiment data | |
|---|---|---|---|---|---|
| | Male | Female | Total | Training set | Testing set |
| OASIS | 160 | 256 | 416 | 120 | 30 |
| IBSR | 14 | 4 | 18 | 12 | 6 |

## 4.4. Results and Discussions

The training and testing were performed on an NVIDIA GeForce RTX 3090 GPU to build the proposed network and use stochastic gradient descent to optimize the loss function. For training, we set the learning rate to 0.001, a high momentum rate of 0.99, and the number of epochs to 10. The Keras framework for the implementation of the proposed work was used. Figures 4-1 and 4-2 show the segmentation results for the axial, coronal, and sagittal planes of the OASIS and IBSR datasets, respectively. The figures show that the proposed approach achieves well-segmented performances for GM, WM, and CSF of the brain MRI on both datasets.

The axial plane shows the most informative details in the central slices of the MRI compared to the other planes. Thus, the segmentation results for the axial planes show the segmentation performance most effectively. In addition, the highlighted boxes in Fig. 4-1 and 4-2 show that the quality of sagittal and coronal images is highly promising without any difference in every detail. From the results of Fig. 4-1 and 4-2, it can be inferred that the proposed method can extract complicated pattern features from all three planes. We evaluated the performance of the proposed method using quantitative metrics. Table 2 presents the DSC, JI, Accuracy, Precision and Recall which are popular evaluation metrics for comparing the ground truth and segmented results.

|  (a)  |  (b)  |  (c)  |  (d)  |  (e)  |  (f)  |

Figure 4-1. Segmentation results for the axial, coronal, and sagittal planes of the brain MRI image (top to bottom) on the OASIS dataset using the proposed method. (a) Original input images, (b) ground truth segmentation map, (c) their predicted segmentation map obtained by using the proposed method, (d) predicted GM (binary map), (e) predicted CSF (binary map), (f) predicted WM (binary map).



|  (a)  |  (b)  |  (c)  |  (d)  |  (e)  |  (f)  |

Figure 4-2. Segmentation results for the axial, coronal, and sagittal planes of the brain MRI image (top to bottom) on the IBSR dataset using the proposed method. (a) Original input images, (b) ground truth segmentation map, (c) their predicted segmentation map.

The most extensively used metrics for measuring the performance of segmentation methods are the DSC [58] and JI [59] measurements. To compute the similarity between two sample sets for segmentation, all evaluation metrics were applied. These metrics determine the match between the predicted segmentation map and the corresponding ground-truth segmentation map. The evaluation metrics for brain tissue segmentation are defined as follows:

Table 2. The formulation of evaluation metrics

| | |
|---|---|
| Dice similarity coefficient(*DSC*) | 2.TP/(2.TP+FP+FN) |
| Jaccard Index (*JI*) | TP/(TP+FP+FN) |
| Accuracy | (TP+TN)/(TP+TN+FP+FN) |
| Precision | TP/(TP+FP) |
| Recall | TP/(TP+FN) |
| Hausdorff distance (*HD*) | $HD = max \times \begin{Bmatrix} \max\limits_{x \in X} \min\limits_{y \in Y} D\{X,Y\}, \\ \max\limits_{y \in Y} \min\limits_{x \in X} D\{X,Y\} \end{Bmatrix}$ |
| Mean squared error (*MSE*) | $MSE = \frac{1}{RC} \sum_{i=1}^{R} \sum_{j=1}^{C} (X-Y)^2$ |

where TP, TN, FP, and FN represent true positives, true negatives, false positives, and false negatives, respectively. We also assessed the segmentation performance in terms of the mean square error (MSE), which is the average square difference between the original X and predicted Y values. The Hausdorff distance (HD) [60] was used to determine the dissimilarity between two sets in a metric space. The two sets of small Hausdorff distances are almost identical. HD and MSE are computed as listed in Table 2, where D is the Euclidean distance between two pixels, and R and C are the image height and width, respectively. To compare the segmentation results of various network architectures, we performed an experiment on SegNet [20], U-net [4], U-SegNet [21], U-net++ [23], and CE-net [29] models under the same

29

experimental conditions. As shown in Fig. 4-3 and 4-4, the proposed method shows superior results in terms of the quality of the segmentation map compared to those of other conventional methods. Although the skip connections in the U-net improve feature representations by combining low-level and high-level information, they suffer from a large semantic gap between low- and high-resolution feature maps, resulting in high misclassification rates of brain tissues. Furthermore, for medical images with low contrast, blurred boundaries between different tissues, the segmentation accuracies of U-net and SegNet are significantly degraded. Because the network layers in U-net++ are connected through a series of nested, dense skip pathways, leading to redundant learning of features, they did not show good performance. In particular, it can be observed that there are misclassification results in the feature maps generated by SegNet, U-net, and U-net++ in the red boxes of Fig. 4-3(c) and 4-4(c). Although U-SegNet with pooling indices and skip connections yields better segmentation results, it fails to capture fine details, as shown in Fig. 4-4(c). From the highlighted red boxes in Fig. 4-4, it can be observed that U-SegNet fails to identify differences between WM and GM tissues, and most of the GM tissues are incorrectly predicted as WM. The CE-net extracts multi-scale information through a context encoder block for the segmentation of medical images. However, the context encoder block is employed only at the bottleneck layer of the model. Thus this multi-scale information could be irrelevant by the time it reaches the final decoder layer for classification. To overcome these limitations, we extract multi-scale information at each network layer followed by the GAM to enhance the segmentation performance by directing attention to related areas [61]. This improved segmentation can be observed in the results obtained using the proposed method.

| (a) | (b) | (c) | (d) | (e) | (f) |

Figure 4-3. Segmentation results for GM, CSF, and WM from brain MRI image using the existing methods and the proposed method on OASIS dataset: (a) original input image; (b) ground-truth segmentation map; (c) their segmentation results obtained SegNet, U-net, U-SegNet, U-net++,CE-net, and the proposed method (top to bottom); (d) CSF maps obtained by SegNet, U-net, U-SegNet, U-net++,CE-net, and the proposed method (top to bottom); (e) GM maps obtained by SegNet, U-net, U-SegNet, U-net++,CE-net, and the proposed method (top to bottom); (f) WM maps obtained by SegNet, U-net, U-SegNet, U-net++,CE-net, and the proposed method (top to bottom).

31

|  | (a) | (b) | (c) | (d) | (e) | (f) |

Figure 4-4. Segmentation results for GM, CSF, and WM from brain MRI image using the existing methods and the proposed method on IBSR dataset: (a) original input image; (b) ground-truth segmentation map; (c) their segmentation results obtained SegNet, U-net, U-SegNet, U-net++,CE-net, and the proposed method (top to bottom); (d) CSF maps obtained by SegNet, U-net, U-SegNet, U-net++,CE-net, and the proposed method (top to bottom); (e) GM maps obtained by SegNet, U-net, U-SegNet, U-net++,CE-net, and the proposed method (top to bottom); (f) WM maps obtained by SegNet, U-net, U-SegNet, U-net++,CE-net, and the proposed method (top to bottom).

Similar results were observed for the segmentations obtained from the IBSR images, as shown in Fig. 4-4. It can be observed that the proposed network obtains finer details than the other architectures. These results indicate that our proposed approach can strongly recover finer segmentation details while bypassing distractions between tissue boundary regions. The quantitative analysis of the proposed method is performed in comparison with SegNet [20], U-net [4], U-SegNet [21], U-net++ [23], and CE-net [29]. Table 3 and Table 4 show the comparative results in terms of the average and standard deviation of DSC, JI, and HD metrics, respectively. Table 5 and Table 6 show the comparative results in terms of the average and standard deviation of accuracy, precision, and recall metrics, respectively. As shown in Table 3 to Table 7 the proposed network achieves significant improvements of 10%, 3%, 2%, 2%, and 1% (in terms of DSC) over SegNet [20], U-net [4], U-SegNet [21], U-net++ [23], and CE-net [29], respectively, and obtained a lower MSE value of 0.003 on average. In addition, the maximum standard deviations for accuracy, precision, and recall are 0.092, 0.098, and 0.099, respectively, which are close to the mean values; This indicates that the pixel predicted values are fitted well to the ground truth values without much data variation. For each encoder map, SegNet [20] stores only the max-pooling indices, i.e., the maximum feature value positions in each pooling window are stored and used for up-sampling. This improves boundary delineation with 3.5 million parameters with approximately 1.4 hours of training time, requiring fewer resources among the existing methods in our proposed method.

Table 3. Comparisons between the segmentation results in terms of DSC, JI, and HD for the proposed method and conventional methods on OASIS dataset.

**OASIS**

**Axial plane**

| Method | WM | | | GM | | | CSF | | |
|---|---|---|---|---|---|---|---|---|---|
| | DSC | JI | HD | DSC | JI | HD | DSC | JI | HD |
| SegNet [20] | 0.89±0.087 | 0.80±0.096 | 4.74±0.077 | 0.86±0.069 | 0.75±0.089 | 4.69±0.053 | 0.85±0.048 | 0.74±0.068 | 4.12±0.079 |
| U-net [4] | 0.93±0.059 | 0.87±0.068 | 4.16±0.064 | 0.92±0.048 | 0.85±0.038 | 4.24±0.046 | 0.90±0.076 | 0.82±0.055 | 3.82±0.039 |
| U-SegNet [21] | 0.94±0.048 | 0.89±0.055 | 3.91±0.057 | 0.93±0.056 | 0.87±0.058 | 4.11±0.033 | 0.92±0.024 | 0.85±0.093 | 3.64±0.047 |
| U-net++[23] | 0.95±0.075 | 0.90±0.042 | 3.78±0.048 | 0.94±0.035 | 0.89±0.072 | 3.84±0.025 | 0.93±0.039 | 0.87±0.046 | 3.56±0.036 |
| CE-net[29] | 0.95±0.039 | 0.90±0.074 | 3.65±0.050 | 0.94±0.042 | 0.89±0.041 | 3.57±0.044 | 0.93±0.043 | 0.87±0.037 | 3.21±0.061 |
| Proposed | 0.97±0.025 | 0.94±0.062 | 3.13±0.037 | 0.95±0.029 | 0.90±0.033 | 3.16±0.030 | 0.94±0.022 | 0.89±0.043 | 2.44±0.038 |

**Coronal plane**

| Method | WM | | | GM | | | CSF | | |
|---|---|---|---|---|---|---|---|---|---|
| | DSC | JI | HD | DSC | JI | HD | DSC | JI | HD |
| SegNet [20] | 0.87±0.098 | 0.77±0.038 | 5.21±0.023 | 0.85±0.044 | 0.74±0.068 | 5.49±0.053 | 0.83±0.079 | 0.71±0.056 | 5.87±0.084 |
| U-net [4] | 0.94±0.065 | 0.89±0.049 | 4.88±0.042 | 0.93±0.057 | 0.87±0.056 | 4.95±0.042 | 0.92±0.063 | 0.85±0.029 | 5.34±0.073 |
| U-SegNet [21] | 0.95±0.049 | 0.90±0.029 | 4.23±0.039 | 0.94±0.081 | 0.89±0.043 | 4.48±0.088 | 0.92±0.054 | 0.85±0.047 | 4.97±0.039 |
| U-net++[23] | 0.94±0.076 | 0.89±0.073 | 4.05±0.047 | 0.93±0.042 | 0.87±0.037 | 4.29±0.044 | 0.93±0.048 | 0.87±0.036 | 4.72±0.043 |
| CE-net[29] | 0.95±0.031 | 0.90±0.026 | 3.98±0.076 | 0.94±0.038 | 0.89±0.040 | 4.17±0.071 | 0.93±0.039 | 0.87±0.029 | 4.17±0.050 |
| Proposed | 0.96±0.043 | 0.92±0.041 | 3.52±0.024 | 0.96±0.020 | 0.92±0.031 | 3.76±0.029 | 0.94±0.015 | 0.89±0.016 | 3.95±0.033 |

**Sagittal plane**

| Method | WM | | | GM | | | CSF | | |
|---|---|---|---|---|---|---|---|---|---|
| | DSC | JI | HD | DSC | JI | HD | DSC | JI | HD |
| SegNet [20] | 0.88±0.096 | 0.79±0.054 | 5.53±0.027 | 0.85±0.083 | 0.74±0.055 | 5.26±0.033 | 0.84±0.040 | 0.72±0.077 | 5.69±0.088 |
| U-net [4] | 0.94±0.068 | 0.89±0.070 | 5.11±0.030 | 0.92±0.074 | 0.85±0.030 | 5.11±0.026 | 0.93±0.058 | 0.87±0.053 | 5.21±0.079 |
| U-SegNet [21] | 0.95±0.077 | 0.90±0.049 | 4.72±0.042 | 0.93±0.066 | 0.87±0.043 | 4.67±0.042 | 0.93±0.037 | 0.87±0.060 | 4.75±0.063 |
| U-net++[23] | 0.95±0.060 | 0.90±0.036 | 4.46±0.031 | 0.94±0.038 | 0.88±0.029 | 4.32±0.019 | 0.94±0.063 | 0.89±0.041 | 4.56±0.041 |
| CE-net[29] | 0.95±0.043 | 0.90±0.064 | 4.13±0.020 | 0.94±0.025 | 0.89±0.035 | 4.25±0.034 | 0.94±0.051 | 0.89±0.062 | 4.28±0.055 |
| Proposed | 0.96±0.027 | 0.92±0.044 | 3.58±0.023 | 0.95±0.011 | 0.90±0.040 | 3.98±0.023 | 0.95±0.031 | 0.90±0.018 | 3.43±0.031 |

Table 4. Comparisons between the segmentation results in terms of DSC, JI, and HD for the proposed method and conventional methods on IBSR dataset.

**IBSR**

**Axial plane**

| Method | WM | | | GM | | | CSF | | |
|---|---|---|---|---|---|---|---|---|---|
| | DSC | JI | HD | DSC | JI | HD | DSC | JI | HD |
| SegNet [20] | 0.72±0.036 | 0.56±0.042 | 6.51±0.65 | 0.75±0.049 | 0.60±0.058 | 6.53±0.91 | 0.68±0.099 | 0.52±0.095 | 6.96±0.46 |
| U-net [4] | 0.89±0.022 | 0.80±0.034 | 5.14±0.51 | 0.91±0.017 | 0.83±0.023 | 4.87±0.51 | 0.84±0.065 | 0.72±0.079 | 5.24±0.31 |
| U-SegNet [21] | 0.90±0.043 | 0.82±0.051 | 4.76±0.39 | 0.92±0.053 | 0.85±0.028 | 4.45±0.65 | 0.84±0.029 | 0.72±0.048 | 4.84±0.18 |
| U-net++[23] | 0.88±0.085 | 0.79±0.096 | 5.37±0.36 | 0.89±0.037 | 0.80±0.049 | 5.17±0.29 | 0.83±0.058 | 0.71±0.082 | 5.34±0.64 |
| CE-net[29] | 0.89±0.055 | 0.81±0.073 | 4.98±0.84 | 0.90±0.068 | 0.82±0.083 | 4.95±0.38 | 0.82±0.037 | 0.69±0.031 | 4.74±0.93 |
| Proposed | 0.91±0.085 | 0.83±0.064 | 4.45±0.57 | 0.93±0.076 | 0.87±0.016 | 4.23±0.92 | 0.85±0.097 | 0.74±0.023 | 4.26±0.79 |

**Coronal plane**

| Method | WM | | | GM | | | CSF | | |
|---|---|---|---|---|---|---|---|---|---|
| | DSC | JI | HD | DSC | JI | HD | DSC | JI | HD |
| SegNet [20] | 0.70±0.061 | 0.54±0.051 | 6.32±0.82 | 0.73±0.037 | 0.57±0.062 | 6.21±0.84 | 0.66±0.054 | 0.49±0.086 | 6.84±0.75 |
| U-net [4] | 0.88±0.035 | 0.79±0.034 | 5.45±0.67 | 0.90±0.014 | 0.83±0.056 | 5.17±0.38 | 0.83±0.012 | 0.71±0.043 | 5.54±0.47 |
| U-SegNet [21] | 0.89±0.076 | 0.80±0.046 | 4.61±0.21 | 0.91±0.035 | 0.82±0.043 | 4.56±0.19 | 0.84±0.085 | 0.72±0.093 | 4.83±0.25 |
| U-net++[23] | 0.88±0.021 | 0.79±0.073 | 5.21±0.39 | 0.91±0.093 | 0.82±0.074 | 5.24±0.24 | 0.82±0.034 | 0.69±0.067 | 5.73±0.39 |
| CE-net[29] | 0.89±0.034 | 0.80±0.851 | 4.89±0.21 | 0.90±0.049 | 0.83±0.068 | 5.98±0.93 | 0.83±0.056 | 0.71±0.042 | 5.21±0.20 |
| Proposed | 0.90±0.039 | 0.82±0.088 | 4.24±0.43 | 0.92±0.019 | 0.86±0.035 | 4.31±0.67 | 0.84±0.078 | 0.72±0.097 | 4.55±0.12 |

**Sagittal plane**

| Method | WM | | | GM | | | CSF | | |
|---|---|---|---|---|---|---|---|---|---|
| | DSC | JI | HD | DSC | JI | HD | DSC | JI | HD |
| SegNet [20] | 0.71±0.043 | 0.55±0.039 | 6.49±0.61 | 0.74±0.073 | 0.59±0.059 | 6.36±0.76 | 0.65±0.083 | 0.48±0.092 | 6.99±0.41 |
| U-net [4] | 0.86±0.029 | 0.75±0.062 | 5.75±0.37 | 0.89±0.036 | 0.80±0.041 | 5.77±0.21 | 0.80±0.071 | 0.67±0.019 | 5.83±0.15 |
| U-SegNet [21] | 0.87±0.016 | 0.77±0.048 | 4.89±0.14 | 0.90±0.069 | 0.82±0.046 | 5.42±0.06 | 0.81±0.096 | 0.68±0.073 | 4.98±0.09 |
| U-net++[23] | 0.85±0.083 | 0.74±0.039 | 4.57±0.54 | 0.88±0.077 | 0.79±0.081 | 4.96±0.22 | 0.79±0.049 | 0.65±0.069 | 5.60±0.44 |
| CE-net[29] | 0.86±0.054 | 0.75±0.025 | 5.34±0.66 | 0.89±0.051 | 0.80±0.037 | 5.86±0.55 | 0.79±0.033 | 0.6±0.022 | 5.25±0.37 |
| Proposed | 0.88±0.035 | 0.79±0.073 | 4.63±0.36 | 0.91±0.028 | 0.83±0.083 | 5.30±0.18 | 0.82±0.053 | 0.75±0.011 | 4.12±0.66 |

Table 5. Comparisons between the segmentation results in terms of accuracy, precision and recall for the proposed method and conventional methods on OASIS dataset.

| Method | WM | | | GM | | | CSF | | |
|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | Precision | Recall | Accuracy | Precision | Recall | Accuracy | Precision | Recall |
| **OASIS** | | | | | | | | | |
| **Axial plane** | | | | | | | | | |
| **SegNet [20]** | 0.96±0.087 | 0.97±0.043 | 0.84±0.077 | 0.95±0.069 | 0.90±0.075 | 0.95±0.005 | 0.95±0.048 | 0.90±0.068 | 0.85±0.069 |
| **U-net [4]** | 0.97±0.058 | 0.98±0.037 | 0.89±0.097 | 0.97±0.081 | 0.92±0.034 | 0.93±0.080 | 0.97±0.044 | 0.88±0.077 | 0.97±0.005 |
| **U-SegNet [21]** | 0.98±0.059 | 0.94±0.020 | 0.97±0.072 | 0.98±0.048 | 0.97±0.041 | 0.90±0.010 | 0.97±0.076 | 0.87±0.059 | 0.97±0.033 |
| **U-net++[23]** | 0.99±0.051 | 0.95±0.001 | 0.96±0.038 | 0.99±0.007 | 0.94±0.098 | 0.94±0.003 | 0.98±0.055 | 0.93±0.082 | 0.93±0.077 |
| **CE-net[29]** | 0.99±0.053 | 0.95±0.010 | 0.95±0.088 | 0.97±0.011 | 0.95±0.010 | 0.94±0.040 | 0.98±0.055 | 0.93±0.042 | 0.94±0.041 |
| **Proposed** | 0.99±0.038 | 0.96±0.078 | 0.95±0.086 | 0.99±0.018 | 0.95±0.033 | 0.95±0.000 | 0.99±0.058 | 0.93±0.037 | 0.95 ±0.019 |
| **Coronal plane** | | | | | | | | | |
| **SegNet [20]** | 0.95±0.092 | 0.96±0.098 | 0.85±0.032 | 0.96±0.056 | 0.83±0.053 | 0.92±0.045 | 0.95±0.049 | 0.88±0.066 | 0.85±0.085 |
| **U-net [4]** | 0.96±0.065 | 0.97±0.088 | 0.87±0.099 | 0.96±0.045 | 0.91±0.054 | 0.93±0.067 | 0.96±0.033 | 0.86±0.089 | 0.96±0.053 |
| **U-SegNet [21]** | 0.96±0.054 | 0.93±0.072 | 0.96±0.068 | 0.97±0.053 | 0.95±0.023 | 0.91±0.052 | 0.98±0.024 | 0.85±0.017 | 0.97±0.038 |
| **U-net++[23]** | 0.98±0.067 | 0.94±0.023 | 0.96±0.045 | 0.98±0.012 | 0.92±0.059 | 0.94±0.042 | 0.97±0.018 | 0.90±0.056 | 0.94±0.063 |
| **CE-net[29]** | 0.99±0.050 | 0.95±0.069 | 0.94±0.089 | 0.96±0.002 | 0.94±0.044 | 0.93±0.051 | 0.98±0.029 | 0.90±0.039 | 0.94±0.051 |
| **Proposed** | 0.99±0.032 | 0.96±0.041 | 0.94±0.0.66 | 0.99±0.018 | 0.95±0.031 | 0.95±0.019 | 0.99±0.014 | 0.93±0.026 | 0.95±0.024 |
| **Sagittal plane** | | | | | | | | | |
| **SegNet [20]** | 0.96±0.091 | 0.96±0.064 | 0.86±0.047 | 0.96±0.080 | 0.82±0.015 | 0.93±0.023 | 0.95±0.020 | 0.87±0.022 | 0.85±0.008 |
| **U-net [4]** | 0.97±0.062 | 0.97±0.025 | 0.88±0.010 | 0.97±0.054 | 0.92±0.033 | 0.94±0.006 | 0.96±0.054 | 0.87±0.020 | 0.96±0.079 |
| **U-SegNet [21]** | 0.97±0.047 | 0.92±0.044 | 0.94±0.032 | 0.97±0.026 | 0.94±0.053 | 0.92±0.047 | 0.98±0.067 | 0.86±0.081 | 0.95±0.083 |
| **U-net++[23]** | 0.98±0.050 | 0.93±0.057 | 0.95±0.051 | 0.96±0.048 | 0.93±0.027 | 0.94±0.059 | 0.95±0.068 | 0.92±0.061 | 0.93±0.047 |
| **CE-net[29]** | 0.98±0.044 | 0.94±0.054 | 0.94±0.070 | 0.94±0.026 | 0.92±0.005 | 0.95±0.033 | 0.97±0.010 | 0.90±0.014 | 0.94±0.065 |
| **Proposed** | 0.99±0.021 | 0.96±0.032 | 0.95±0.051 | 0.99±0.018 | 0.95±0.042 | 0.95±0.013 | 0.99±0.012 | 0.93±0.012 | 0.95±0.051 |

Table 6. Comparisons between the segmentation results in terms of accuracy, precision and recall for the proposed method and conventional methods on IBSR dataset.

**IBSR**

**Axial plane**

| Method | WM | | | GM | | | CSF | | |
|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | Precision | Recall | Accuracy | Precision | Recall | Accuracy | Precision | Recall |
| SegNet [20] | 0.96±0.053 | 0.84±0.066 | 0.62±0.039 | 0.96±0.010 | 0.52±0.086 | 0.83±0.087 | 0.96±0.060 | 0.53±0.036 | 0.55±0.036 |
| U-net [4] | 0.97±0.007 | 0.74±0.013 | 0.89±0.073 | 0.98±0.086 | 0.91±0.056 | 0.80±0.089 | 0.99±0.068 | 0.70±0.074 | 0.72±0.021 |
| U-SegNet [21] | 0.97±0.017 | 0.86±0.046 | 0.74±0.048 | 0.98±0.077 | 0.84±0.033 | 0.87±0.062 | 0.99±0.057 | 0.61±0.022 | 0.62±0.009 |
| U-net++[23] | 0.98±0.042 | 0.83±0.038 | 0.92±0.035 | 0.99±0.010 | 0.90±0.054 | 0.88±0.022 | 0.99±0.068 | 0.76±0.088 | 0.61±0.008 |
| CE-net[29] | 0.99±0.018 | 0.76±0.024 | 0.92±0.019 | 0.98±0.094 | 0.91±0.099 | 0.82±0.057 | 0.99±0.067 | 0.73±0.042 | 0.64±0.038 |
| Proposed | 0.99±0.042 | 0.85±0.036 | 0.93±0.039 | 0.99±0.013 | 0.91±0.086 | 0.87±0.039 | 0.99±0.064 | 0.85±0.065 | 0.74 ±0.061 |

**Coronal plane**

| Method | WM | | | GM | | | CSF | | |
|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | Precision | Recall | Accuracy | Precision | Recall | Accuracy | Precision | Recall |
| SegNet [20] | 0.95±0.092 | 0.85±0.098 | 0.65±0.032 | 0.94±0.056 | 0.55±0.013 | 0.85±0.015 | 0.95±0.013 | 0.60±0.076 | 0.60±0.084 |
| U-net [4] | 0.96±0.065 | 0.74±0.088 | 0.87±0.099 | 0.95±0.015 | 0.90±0.034 | 0.82±0.062 | 0.94±0.093 | 0.80±0.085 | 0.75±0.059 |
| U-SegNet [21] | 0.97±0.054 | 0.85±0.072 | 0.75±0.068 | 0.95±0.023 | 0.85±0.053 | 0.88±0.032 | 0.97±0.087 | 0.65±0.047 | 0.65±0.088 |
| U-net++[23] | 0.98±0.067 | 0.82±0.023 | 0.90±0.045 | 0.96±0.013 | 0.91±0.079 | 0.90±0.044 | 0.97±0.015 | 0.84±0.053 | 0.70±0.064 |
| CE-net[29] | 0.99±0.050 | 0.76±0.069 | 0.89±0.089 | 0.97±0.042 | 0.92±0.094 | 0.85±0.055 | 0.98±0.022 | 0.75±0.032 | 0.72±0.061 |
| Proposed | 0.99±0.032 | 0.85±0.041 | 0.92±0.0.66 | 0.99±0.015 | 0.92±0.001 | 0.90±0.069 | 0.99±0.034 | 0.85±0.016 | 0.78±0.023 |

**Sagittal plane**

| Method | WM | | | GM | | | CSF | | |
|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | Precision | Recall | Accuracy | Precision | Recall | Accuracy | Precision | Recall |
| SegNet [20] | 0.96±0.001 | 0.84±0.011 | 0.68±0.087 | 0.95±0.085 | 0.60±0.014 | 0.88±0.075 | 0.94±0.074 | 0.62±0.027 | 0.62±0.019 |
| U-net [4] | 0.96±0.042 | 0.76±0.009 | 0.88±0.017 | 0.95±0.036 | 0.91±0.052 | 0.80±0.007 | 0.95±0.045 | 0.84±0.063 | 0.77±0.065 |
| U-SegNet [21] | 0.95±0.044 | 0.84±0.078 | 0.77±0.036 | 0.96±0.052 | 0.88±0.086 | 0.85±0.075 | 0.92±0.086 | 0.70±0.085 | 0.66±0.026 |
| U-net++[23] | 0.97±0.082 | 0.80±0.056 | 0.91±0.045 | 0.97±0.074 | 0.90±0.075 | 0.91±0.045 | 0.96±0.024 | 0.88±0.074 | 0.72±0.075 |
| CE-net[29] | 0.98±0.068 | 0.79±0.043 | 0.85±0.032 | 0.94±0.052 | 0.91±0.015 | 0.88±0.085 | 0.97±0.011 | 0.79±0.095 | 0.75±0.085 |
| Proposed | 0.99±0.041 | 0.86±0.012 | 0.93±0.013 | 0.99±0.015 | 0.92±0.042 | 0.92±0.056 | 0.99±0.045 | 0.86±0.068 | 0.80±0.045 |

Table 7. Mean square error (MSE) comparison between the segmentation accuracy for the proposed method and conventional methods on OASIS and IBSR datasets.

| | Mean Square Error (MSE) | | | | | |
|---|---|---|---|---|---|---|
| | SegNet [20] | U-net [4] | U-SegNet [21] | U-net++ [23] | CE-net [29] | Proposed method |
| OASIS | 0.021 | 0.006 | 0.005 | 0.004 | 0.004 | 0.003 |
| IBSR | 0.013 | 0.008 | 0.007 | 0.005 | 0.005 | 0.004 |

SegNet tends to miss several fine details because when performing up-sampling from low-resolution feature maps, it loses adjacent information. On the other hand, U-net uses skip connections as the core of the architecture, which blends deep, coarse information with shallow, fine semantic information. A drawback of U-net is its significant memory requirement because lower-level features in the up-sampling process must be stored for further concatenation. Because U-net uses low-level feature maps for up-sampling, translation invariance is often compromised. Moreover, U-SegNet [21] tends to be insensitive to fine details, and it is evident from the difficulty in identifying boundaries between adjacent tissues, such as WM and GM. The design of atrous convolution followed by multi-kernel max-pooling in the CE-net helps to capture multi-scale information and avoids the acquisition of redundant information. However, the multi-scale feature extraction capability of the CE-net is limited to the bottleneck layer, leading to poor feature presentation at the final decoder layer. The segmentation maps generated by these existing models have a relatively low resolution because of the pooling layers in the encoder stage. Hence, to preserve the high spatial resolution, the pooling layers must be removed. However, since convolution is a rather local operation, SegNet, U-net, U-SegNet, U-net++, and CE-net models would not learn holistic features in the images without pooling layers. Our proposed

method presents a multi-scale feature fusion scheme combined with GAM as a potential solution to the problems discussed above and produces improved segmentation accuracy. The max-pooled output was filtered with 1×1, 3×3 kernels. Then, they are concatenated together and can extract the global context without losing the resolution of the segmentation map.
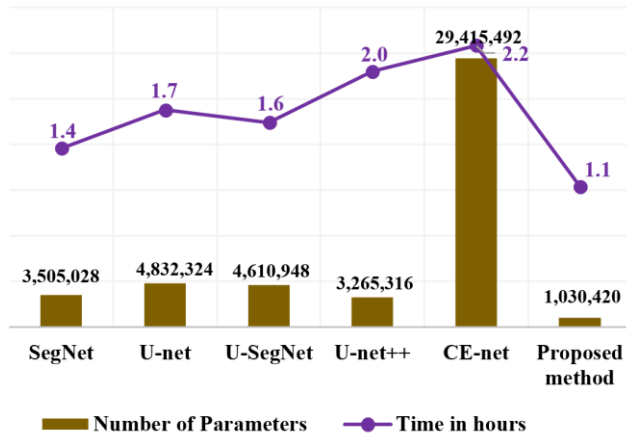


Figure 4-5. Detail data on the number of learnable parameters and computation time for the proposed and conventional methods for the OASIS dataset.

In this way, global information can be exchanged between layers without reducing the resolution, leading to lowered blurring in the segmentation maps. In addition, the GAM at the encoder enables the presentation of global context information as a guide for low-level features to extract the original resolutions for segmentation. The GAM at the decoder shows that the combination of global features and local features is essential to discriminate brain tissues and is consistent with the results from previous studies. Furthermore, uniform input patches allow the network to concentrate better on local information. As a result of the selective integration of spatial information through uniform patches, feature maps followed by multi-scale guided multiple GAMs help in

capturing context information and can efficiently encode complementary information to segment the brain MRI accurately.

As mentioned above, we propose the use of a fire module for fewer learnable parameters while maintaining equivalent accuracy. Figure 4-5 shows the details of the number of learnable parameters and computation time consumed by the proposed method in comparison with conventional methods. Smaller models can be built by arranging a sequence of fire modules that consist of a squeeze layer that has only $1 \times 1$ convolution filters. This serves as an expansion layer that has a combination of $1 \times 1$ and $3 \times 3$ filters. The number of filters in the squeeze layer was defined to be less than the number of $1 \times 1$ and $3 \times 3$ filters in the expand layer. The $1 \times 1$ filters in the squeeze layer down-sample the input channels and decrease the parameters before they are given as an input to the expand layer. The expansion layer consists of both $3 \times 3$ and $1 \times 1$ filters. The $1 \times 1$ filters in the expand layer combine channels and perform cross-channel pooling but cannot recognize spatial structures. The $3 \times 3$ convolution filter in the expand layer identifies the spatial representation. The model becomes more descriptive by integrating these two distinct size filters while running on lower parameters. Hence, fire modules reduce the computational load by reducing the parameter maps and building a smaller CNN network that can preserve a higher degree of accuracy. The total parameters in our proposed method are one million parameters, which are 3, 5, 4.5, 3, and 28 times smaller than SegNet, U-net, U-SegNet, U-net++, and CE-net networks, respectively. The training time for the proposed method for the OASIS dataset was 50% of that of the U-net++ and CE-net methods. Compared to traditional approaches, a reduction in memory requirements would result in a substantial decrease in energy and processing requirements.

## 4.5. Ablation Study

We conducted an ablation study on the three simplified versions of the proposed modules to investigate the influence of each selection on the segmentation performance as follows: (i) Squeeze U-SegNet, (ii) Squeeze U-SegNet with multi-scale input, (iii) Squeeze U-SegNet with multi-global attention and (iv) multi-scale Squeeze U-SegNet with multi-global attention (proposed method). The Squeeze U-SegNet was obtained by replacing each convolution block with a fire module in the conventional U-SegNet. The second network proposes that the encoder of the Squeeze U-SegNet includes a multi-scale input layer. This is achieved by max-pooling the input and performing parallel convolution with 1×1, 3×3 kernels, and concatenating these multi-scale features. These fused multi-scale features are concatenated with the corresponding fire module output and fed as input for further max-pooling operations. This process is repeated for all encoding layers. The multi-scale feature module extracts neighbor scale information of global features more precisely while filtering out irrelevant information. The impact of the attention mechanism is explored in the third network, where GAMs are integrated at both the encoder and decoder, forming a multi-attention network. Finally, the multi-scale squeeze U-SegNet with multi-global attention referred to as the proposed method incorporates semantic guidance by combining all the proposed modules. Table 8 lists the results of the individual contributions of different components to segmentation performance. The fire module-based model significantly decreases the requirement of learnable parameters with reduced computation time for model training while maintaining network accuracy. We observe that, compared to the baseline squeeze U-SegNet, the performance of the models integrated with the multi-scale feature fusion input

scheme and with the multi-attention modules is improved by 1.5% and 2%, respectively. Although the multi-scale feature fusion shows a slight increase in the DSC, its contribution combined with GAM provides more network efficiency. Furthermore, the combination of both multi-scale and multiple global attention strategies boosts the performance and yields the best values in the three metrics: 96% (DSC), 91% (JI), 3.1(HD), and with the lowest MSE of 0.003. These results represent an improvement of 2% in DSC compared to the baseline U-SegNet [20], showing the efficiency of the proposed multi-scale guided multi-GAM compared to individual components.

Table 8. Detail data on the number of learnable parameters and computation time for the proposed method and its three simplified versions.

| | GM | | | WM | | | CSF | | | | Computation time(10 epochs) | #Learnable parameters |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DSC | JI | HD | DSC | JI | HD | DSC | JI | HD | MSE | | |
| 1 | 92.05 | 88.06 | 3.52 | 93.37 | 90.42 | 2.8 | 91.65 | 88.06 | 2.0 | 0.006 | 1 h | 768,788 |
| 2 | 93.44 | 89.47 | 4.81 | 94.78 | 91.90 | 4.1 | 93.32 | 90.25 | 3.0 | 0.005 | 1.04 h | 860180 |
| 3 | 94.32 | 89.25 | 5.10 | 95.12 | 91.40 | 4.2 | 94.24 | 89.22 | 3.3 | 0.004 | 1.15 h | 942164 |
| 4 | 95.54 | 91.09 | 3.21 | 96.56 | 92.05 | 3.1 | 94.86 | 90.29 | 3.0 | 0.003 | 1.12 h | 1030420 |

**1: Squeeze U-SegNet, 2: Squeeze U-SegNet with multi-scale input, 3: Squeeze U-SegNet with multi global attention, 4: Proposed**

We also investigated the effects of patch size in terms of training time and segmentation performance. The experiments were performed on the OASIS dataset for three distinct patch sizes (128×128, 64×64, and 32×32).
Table 9 lists the output of the segmentation in terms of the DSC with respect to different patch sizes. It can be observed that smaller patches result in better performance.

42

Table 9. Segmentation accuracy (%) and training time (hours) for the
proposed method with different sizes of input patches.

| Patch size | DSC | | | JI | | | Training time (hours) |
|---|---|---|---|---|---|---|---|
| | WM | GM | CSF | WM | GM | CSF | |
| 128×128 | 96.56 | 95.54 | 95.73 | 92.05 | 91.09 | 90.23 | 1.1 |
| 64×64 | 96.74 | 96.13 | 95.49 | 92.49 | 91.58 | 90.54 | 6.7 |
| 32×32 | 96.91 | 96.85 | 95.73 | 91.88 | 91.76 | 90.71 | 12.4 |

This is because smaller patches create more training data for the network to
train. Moreover, local regions can be restored more accurately. Furthermore,
when the patch size is 128×128, it takes 1.1 h to train the model, whereas the
training time doubles for 32×32 patches with almost identical accuracy. We,
therefore, concluded that a patch size of 128×128 provides a fair tradeoff
between the DSC score and the computational time needed to train the model,
based on the results in Table 9.

# 5. CONCLUSION

This paper proposes multi-scale feature extraction with novel global attention-based learning based on the U-SegNet architecture for brain MRI segmentation. The multi-scale data provide rich spatial information and improve the robustness of feature extraction. The global attention module provides the global context as guidance for low-level features to select category localization details. The squeeze and expand layers lead to the generation of one million parameters, which are 3, 5, 4.5, 3, and 28 times smaller than SegNet, U-net, U-SegNet, U-net++, and CE-net networks, respectively. Our proposed network obtains the best DSC value of 96%. The training time for the proposed method for the OASIS dataset is 50% of that of the U-net++ and CE-net methods. Our validation proves that the network operating on patch-wise input, integrated with multi-scale global attention and fire modules, will yield an efficient brain MRI segmentation model. The proposed model can be easily extended to complex network architectures owing to flexibility and adaptability with faster computation. Hence, a three-dimensional (3D) segmentation model can be devised using the extended model of the proposed architecture as future works.

# ACKNOWLEDGEMENT

I would like to express my deepest gratitude to my advisor, Prof. Bumshik Lee, for his support, patience, and encouragement throughout my graduate study. His technical and editorial advice and words of encouragement and guideless were essential to complete this research successfully.

I am expressing my deep gratitude to my family for their love, understanding, supports, and encouragement during the period of this research.

Finally, my deepest thanks to my fellow lab mates for their support throughout the course of my research.

# REFERENCES

[1]     M. Thompson and L. G. Apostolova, "Computational anatomical methods as applied to aging and dementia," *Br J Radiol*, vol. 80 Spec No 2, pp. S78-91, Dec. 2007.

[2]     Whitwell JL, Przybelski SA, Weigand SD, Knopman DS, Boeve BF, Petersen RC, and Jack CR "3D maps from multiple MRI illustrate changing atrophy patterns as subjects progress from mild cognitive impairment to Alzheimer's disease," *Brain*, vol. 130, no. Pt 7, pp. 1777–1786, Jul. 2007.

[3]     R. C. Petersen, G. E. Smith, S. C. Waring, R. J. Ivnik, E. G. Tangalos, and E. Kokmen, "Mild Cognitive Impairment: Clinical Characterization and Outcome," *Arch Neurol*, vol. 56, no. 3, pp. 303–308, Mar. 1999.

[4]     Sperling RA, Aisen PS, Beckett LA, Bennett DA, Craft S, Fagan AM, Iwatsubo T, Siemers E, Stern Y, Yaffe K, Carrillo MC, Thies B, Morrison-Bogorad M, Wagster MV, and Phelps CH, "Toward defining the preclinical stages of Alzheimer's disease: recommendations from the National Institute on Aging-Alzheimer's Association workgroups on diagnostic guidelines for Alzheimer's disease," *Alzheimers Dement*, vol. 7, no. 3, pp. 280–292, May 2011.

[5]     Tábuas-Pereira, M., Baldeiras, I., Duro, D., Santiago, B., Ribeiro, M. H., Leitão, M. J., Oliveira, C., and Santana, I "Prognosis of Early-Onset vs. Late-Onset Mild Cognitive Impairment: Comparison of Conversion Rates and Its Predictors," *Geriatrics (Basel)*, vol. 1, no. 2, Apr. 2016.

[6]     Mosconi L, Mistur R, Switalski R, Tsui WH, Glodzik L, Li Y, Pirraglia E, De Santi S, Reisberg B, Wisniewski T, and de Leon MJ., "FDG-PET changes in brain glucose metabolism from normal cognition to pathologically verified Alzheimer's disease," *Eur. J. Nucl. Med. Mol. Imaging*, vol. 36, no. 5, pp. 811–822, May 2009.

[7]     N. Srivastava and R. Salakhutdinov, "Multimodal Learning with Deep Boltzmann Machines," *J. Mach. Learn. Res.*, vol. 15, no. 1, pp. 2949–2980, Jan. 2014.

[8]     J. B. Colby, J. D. Rudie, J. A. Brown, P. K. Douglas, M. S. Cohen, and Z. Shehzad, "Insights into multimodal imaging classification of ADHD," *Front Syst Neurosci*, vol. 6, p. 59, 2012.

[9]     Charles Decarli and Steven T. DeKosky and Mony J. de Leon and Norman L. Foster and Nick and Fox and Richard Frank and Richard S. and Thies and W. Michael and Weiner and Zaven S. Khachaturian "The Use of MRI and PET for Clinical Diagnosis of Dementia and Investigation of Cognitive Impairment : *A Consensus Report*," 2004.

[10]    Liu, C.-Y. Wee, H. Chen, and D. Shen, "Inter-modality relationship constrained multi-modality multi-task feature selection for Alzheimer's Disease and mild cognitive impairment identification," *Neuroimage*, vol. 84, pp. 466–475, Jan. 2014.

[11]    P. Vemuri and C. R. Jack, "Role of structural MRI in Alzheimer's disease," *Alzheimers Res Ther*, vol. 2, no. 4, p. 23, Aug. 2010.

[12]    Bernard, Olivier and Lalande, Alain and Zotti, Clement and Cervenansky, Frederick and Yang, Xin and Heng, Pheng-Ann and Cetin, Yoonmi and Patravali, Jay and Jain, Shubham and Humbert, Olivier and Jodoin, Pierre-Marc, "Deep learning techniques for automatic MRI cardiac multi-structures segmentation and diagnosis: Is the problem solved?" *IEEE transactions on medical imaging*, vol. 37, no. 11, pp. 2514–2525, 2018.

[13]    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*, pp. 3431–3440, 2015.

[14]    O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in International Conference on Medical image computing and computer-assisted intervention. Springer, pp. 234–241, 2015.

[15]    J. Dolz, L. Massoptier, and M. Vermandel, "Segmentation algorithms of subcortical brain structures on MRI for radiotherapy and radiosurgery: a survey," *IRBM*, vol. 36, no. 4, pp. 200–212, 2015.

[16]   Fechter T, Adebahr S, Baltas D, Ben Ayed I, Desrosiers C, Dolz J. Esophagus, "Esophagus segmentation in CT via 3D fully convolutional neural network and random walk," *Medical physics*, vol. 44, no. 12, pp. 6341–6352, 2017.

[17]   Luna M., Park S.H,  "3D Patchwise U-net with Transition Layers for MR Brain Segmentation", In: Crimi A., Bakas S., Kuijf H., Keyvan F., Reyes M., van Walsum T. (eds) Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries. BrainLes 2018. Lecture Notes in Computer Science, vol 11383. Springer, Cham, 2019.

[18]   M. Pawel, D.Hervé, C. Antonio, A. Nicholas, "3D Convolutional Neural Networks for Tumor Segmentation using Long-range 2D Context", *Computerized Medical Imaging and Graphics*,vol.73, 2019.

[19]   F, Xue, T. Nicholas, P. Sohil, M. Craig, "Brain Tumor Segmentation Using an Ensemble of 3D U-nets and Overall Survival Prediction Using Radiomic Features", *Frontiers in Computational Neuroscience*, vol.14, pp. 25, 2020.

[20]   V. Badrinarayanan, A. Kendall, and R. Cipolla, "SegNet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell*., vol. 39, no. 12, pp. 2481–2495, Dec. 2017.

[21]   Kumar, Pulkit & Nagar, Pravin & Arora, Chetan & Gupta, Anubha. (2018). "U-SegNet: Fully Convolutional Neural Network based Automated Brain tissue segmentation Tool ",  *ISBI*, 2017.

[22]   Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, Alan L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2018.

[23]   Zhou, Z., Rahman Siddiquee, M. M., Tajbakhsh, N., & Liang, J, "UNet++: A Nested U-net Architecture for Medical Image Segmentation", *arXiv*, pp.3–11, 2018.

[24] Tao Lei, Risheng Wang, Yong Wan, Bingtao Zhang, Hongying Meng, Asoke K. Nandi, "Medical Image Segmentation Using Deep Learning: A Survey", *arXiv*, pp.3–11,2020.

[25] Yamanakkanavar, N.; Choi, J.Y.; Lee, B. "MRI Segmentation and Classification of Human Brain Using Deep Learning for Diagnosis of Alzheimer's disease: A Survey". *Sensors*, 20, 3243, 2020.

[26] C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu, "Deeplysupervised nets," in *Artif. Intell. Statistics*, 2015, pp. 562–570, 2015.

[27] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in Proceedings of *the IEEE conference on computer vision and pattern recognition*, pp. 2881–2890, 2017.

[28] W. Liu, A. Rabinovich, and A. C. Berg, "Parsenet: Looking wider to see better," *arXiv*: 506.04579, 2015.

[29] Gu, Z., Cheng, J., Fu, H., Zhou, K., Hao, H., Zhao, Y., Liu, J. "CE-net: Context Encoder Network for 2D Medical Image Segmentation", *IEEE Transactions on Medical Imaging*, pp.1-1.2019.

[30] Fei Wang, Mengqing Jiang, Chen Qian, Shuo Yang, Cheng Li, Honggang Zhang, Xiaogang Wang, Xiaoou Tang "Residual attention network for image classification," in Proceedings of the *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3156–3164, 2017.

[31] M. Pedersoli, T. Lucas, C. Schmid, and J. Verbeek, "Areas of attention for image captioning," in Proceedings of the *IEEE International Conference on Computer Vision*, pp. 1242–1250, 2017.

[32] Z. Xiaodong He, Jianfeng Gao, Li Deng, Alex Smola "Stacked attention networks for image question answering," in Proceedings of the *IEEE conference on computer vision and pattern recognition*, pp. 21–29, 2016.

[33] A. P. Parikh, O. Tackstr om, D. Das, and J. Uszkoreit, "A decomposable attention model for natural language inference," in In *EMNLP*, 2016.

[34] A.Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin "Attention is all you need," in *Advances in neural information processing systems*, pp. 5998–6008, 2017.

[35] Y. Nagaraj, Pardhu Madipalli, Jeny Rajan, P. Krishna Kumar, A.V. Narasimhadhan, "Segmentation of intima media complex from carotid ultrasound images using wind driven optimization technique", *Biomedical Signal Processing and Control*, vol.40, pp.462-472, 2018.

[36] Zhouhan Lin, MinweiFeng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and YoshuaBengio. "A structured self-attentive sentence embedding". *arXiv* preprint arXiv:1703.03130, 2017.

[37] H. Zhao, Y. Zhang, S. Liu, J. Shi, C. Change Loy, D. Lin, and J. Jia, "PSANet: Point-wise spatial attention network for scene parsing," in Proceedings of the *European Conference on Computer Vision (ECCV)*, pp. 267–283, 2018.

[38] S. Li, M. Dong, G. Du and X. Mu, "Attention Dense-U-net for Automatic Breast Mass Segmentation in Digital Mammogram," in *IEEE Access*, vol. 7, pp. 59037-59047, 2019.

[39] Jo Schlemper, OzanOktay , MichielSchaap , Mattias Heinrich , Bernhard Kainz , Ben Glocker , Daniel Rueckert, "Attention gated networks: Learning to leverage salient regions in medical images", *Medical Image Analysis*, vol. 53, pp. 197-207, 2019.

[40] Y. Wang, Z. Deng, X. Hu, L. Zhu, X. Yang, X. Xu, P.-A. Heng, and D. Ni, "Deep attentional features for prostate segmentation in ultrasound," in *MICCAI*, 2018.

[41] J. Schlemper, Ozan Oktay, Michiel Schaap, Mattias, Heinrich, Bernhard Kainz, Ben Glocker, Daniel Rueckert "Attention gated networks: Learning to leverage salient regions in medical images," *Medical image analysis*, vol. 53, pp. 197–207, 2019.

[42]    A. Sinha, J. Dolz, "Multi-Scale Self-Guided Attention for Medical Image Segmentation," *IEEE Journal of Biomedical and Health Informatics*, vol. 25, no. 1, pp. 121-130, Jan. 2021.

[43]    S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding," *arXiv* Prepr arXiv1510.00149, 2015.

[44]    LeCun, Y. "Generalization and network design strategies." (1989).

[45]    Nagaraj Y., Asha C.S., Hema Sai Teja A., A.V. Narasimhadhan,"Carotid wall segmentation in longitudinal ultrasound images using structured random forest", *Computers & Electrical Engineering*, vol.69, pp. 753-767, 2018.

[46]    Y. Cheng, D. Wang, P. Zhou, and T. Zhang, "A survey of model compression and acceleration for deep neural networks," *arXiv* Prepr. arXiv1710.09282, 2017.

[47]    L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in Proceedings of the *European conference on computer vision (ECCV)*, pp. 801–818, 2018.

[48]    Horowitz, and William J. Dally, "EIE: efficient inference engine on compressed deep neural network", *SIGARCH Comput. Archit*. Vol 243–254, June 2016.

[49]    F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in Proceedings of the *IEEE conference on computer vision and pattern recognition*, pp. 1251–1258, 2017.

[50]    Krizhevsky, A., Sutskever, I., & Hinton, G. E. "ImageNet classification with deep convolutional neural networks". *Communications of the ACM*, 60(6), 84–90, 2017.

[51]    T. Sheng, C. Feng, S. Zhuo, X. Zhang, L. Shen, and M. Aleksic, "A quantization friendly separable convolution for mobilenets," in 2018 1st Workshop on *Energy Efficient Machine Learning and Cognitive*

*Computing for Embedded Applications* (EMC2), 2018, pp. 14–18, 2018.

[52]    C.Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, Andrew Rabinovich "Going deeper with convolutions," *CVPR*, pp. 1-9, 2015.

[53]    Zisserman, Andrew, "Very Deep Convolutional Networks for Large-Scale Image Recognition", *arXiv* 1409.1556, 2014.

[54]    N. Yamanakkanavar and B. Lee, "Using a Patch-Wise M-Net Convolutional Neural Network for Tissue Segmentation in Brain MRI Images," in *IEEE Access*, vol. 8, pp. 120946-120958, 2020.

[55]    H. Tracy, G. John, C. Csernansky, M. Marcus, S. Daniel, and L. Randy, "Open Access Series of Imaging Studies (OASIS): Cross-Sectional MRI Data in Young, Middle Aged, Nondemented, and Demented Older Adults". *Journal of Cognitive Neuroscience*, vol.19, pp.1498-1507, 2007.

[56]    Chen Y., Almarzouqi S.J., Morgan M.L., Lee A.G. "T1-Weighted Image." In: Schmidt-Erfurth U., Kohnen T. (eds) *Encyclopedia of Ophthalmology*. Springer, Berlin, Heidelberg, 2018.

[57]    IBSR dataset. Available: https://www.nitrc.org/projects/ibsr, 2012.

[58]    L. R. Dice, "Measures of the Amount of Ecologic Association between Species," *Ecology*, 1945.

[59]    P. Jaccard, "The distribution of the flora in the alpine zone", *New Phytol.*, 1912.

[60]    R. Gunter, "Computing the Minimum Hausdorff Distance between Two Point Sets on a Line under Translation" Inf. Process. Lett, vol. 38, pp. 123-127, 1991.

[61]    C. Dayananda, JY. Choi, and B. Lee, "Multi-Scale Squeeze U-SegNet with Multi Global Attention for Brain MRI Segmentation", *Sensors*, 10, 3363, 2021.

# PUBLICATIONS

- **International Journal Papers**

  1. Chaitra Dayananda, Jae-Young Choi, and Bumshik Lee, "Multi-Scale Squeeze U-SegNet with Multi Global Attention for Brain MRI Segmentation", *Sensors*, 10, 3363, 2021. **(Published)**
  2. Chaitra Dayananda, Jae-Young Choi, and Bumshik Lee, "A squeeze U-SegNet architecture based on residual convolution for brain MRI segmentation", *IEEE Transaction on Image processing*. **(Under review)**

- **International Conference Papers**

  1. Chaitra Dayananda and Bumshik Lee, "Squeeze U-SegNet: A CNN architecture for Automatic segmentation of brain MRI", The 9th International Conference on Smart Media and Applications (*SMA 2020*), Jeju, South Korea. **(Presented)**
  2. Chaitra Dayananda, You Wonsang, Jae-Young Choi and Bumshik Lee, "Skin lesion segmentation in dermoscopic images using CNN architecture", The 12th International Conference On ICT Convergence (*ICTC 2021*), Jeju, South Korea. **(Presented)**
  3. Chaitra Dayananda, and Bumshik Lee, "Using asymmetric multi-cross convolutions for skin lesion segmentation on dermoscopic images", *CVPR, 2022*. **(Submitted)**