RESEARCH ARTICLE

-100N. 2000-0002

OPEN ACCESS

Improving Accuracy and Efficiency of Medical Image Segmentation Using One-Point-Five U-Net Architecture with Integrated Attention and Multi-Scale Mechanisms

Muhammad Anang Fathur Rohman¹, Heri Prasetyo¹, Ery Permana Yudha¹, Chih-Hsien Hsia²

¹ Department of Informatics, Universitas Sebelas Maret, Surakarta, Indonesia

²Department of Computer Science and Information Engineering, National IIan University, Yilan, Taiwan

Corresponding author: Heri Prasetyo. (e-mail: heri.prasetyo@staff.uns.ac.id), **Author(s) Email**: Muhammad Anang Fathur Rohman (anangmuhammad245@student.uns.ac.id), Heri Prasetyo (e-mail: heri.prasetyo@staff.uns.ac.id), Ery Permana Yudha (email: erypermana@staff.uns.ac.id), Chih-Hsien Hsia (email: hsiach@niu.edu.tw)

Abstract Medical image segmentation is essential for supporting computer-aided diagnosis (CAD) systems by enabling accurate identification of anatomical and pathological structures across various imaging modalities. However, automated medical image segmentation remains challenging due to low image contrast, significant anatomical variability, and the need for computational efficiency in clinical applications. Furthermore, the scarcity of annotated medical images due to high labelling costs and the requirement of expert knowledge further complicates the development of robust segmentation models. This study aims to address these challenges by proposing One-Point-Five U-Net, a novel deep learning architecture designed to improve segmentation accuracy while maintaining computational efficiency. The main contribution of this work lies in the integration of multiple advanced mechanisms into a compact architecture: ghost modules, Multi-scale Residual Attention (MRA), Enhanced Parallel Attention (EPA) in skip connections, the Convolutional Block Attention Module (CBAM), and Multi-scale Depthwise Convolution (MSDC) in the decoder. The proposed method was trained and evaluated on four public datasets: CVC-ClinicDB, Kvasir-SEG, BUSI, and ISIC2018. One-Point-Five U-Net achieved sensitivity, specificity, accuracy, DSC, and IoU of of 94.89%, 99.63%, 99.23%, 95.41%, and 91.27% on CVC-ClinicDB; 91.11%, 98.60%, 97.33%, 90.93%, and 83.84% on Kvasir-SEG; 85.35%, 98.65%, 96.81%, 87.02%, and 78.18% on BUSI; and 87.67%, 98.11%, 93.68%, 89.27%, and 83.06% on ISIC2018. These results outperform several state-of-the-art segmentation models. In conclusion, One-Point-Five U-Net demonstrates superior segmentation accuracy with only 626,755 parameters and 28.23 GFLOPs, making it a highly efficient and effective model for clinical implementation in medical image analysis.

Keywords: Medical image segmentation; Deep Learning; One-Point-Five U-Net; Efficient Model

I. Introduction

Medical image segmentation has become a key research focus in computer-aided diagnosis (CAD) systems [1]. Its primary objective is to differentiate anatomical and pathological structures within various medical imaging modalities. Medical image segmentation is crucial for assisting clinicians in performing quantitative pathological assessments and provides a reliable foundation for clinical diagnosis [2]. However, automating the segmentation process remains challenging due to several inherent issues [3]. First, medical images typically exhibit low contrast, which blurs the boundaries between different objects [4]. Second, significant variations exist in the shape, size, and location of pathological regions across different patients [5]. Third, there is a high demand for segmentation methods that are both fast and reliable in clinical applications [6]. Fourth, obtaining a huge amount of labeled medical images is difficult due to the high cost and specialized expertise required for annotation [7]. Furthermore, variations in imaging devices and acquisition settings often result in

Manuscript received 10 May 2025; Revised 21 June 2025; Accepted 30 June 2025; Available online 5 July 2025 Digital Object Identifier (**DOI**): https://doi.org/10.35882/jeeemi.v7i3.949

Copyright © 2025 by the authors. This work is an open-access article and licensed under a Creative Commons Attribution-ShareAlike 4.0 International License (CC BY-SA 4.0).

inconsistent image quality, necessitating higher robustness and adaptability from segmentation algorithms.

With the rapid advancement of deep learning, Convolutional Neural Networks (CNNs) have become widely utilized in medical images. CNNs are known for their powerful feature representation, high degree of automation, and reliable segmentation outcomes [8]. The introduction of the Fully Convolutional Network (FCN) by [9] marked a significant milestone as it was the first end-to-end architecture designed explicitly for pixel-level prediction in segmentation applications. However, FCN still lacks precision in segmentation results due to its limited ability to capture the holistic relationship among pixels, often leading to inconsistent preservation of image structures [10]. To address this, [11] introduced a U-shaped convolutional neural network called U-Net. The encoder-decoder architecture of U-Net restores spatial information lost during downsampling, allowing it to retain fine details of small objects. U-Net performs well on small-scale datasets, making it suitable for medical image segmentation tasks. Due to these advantages, U-Net has become a benchmark in the field, inspiring numerous derivative works.

Given the high parameter count of U-Net, [12] proposed Half U-Net, an asymmetric variant of the U-Net architecture. This method integrates full-scale feature fusion from UNet3+, ghost modules, and uniform channel distribution at each level. Across three medical image datasets, Half U-Net successfully reduced the number of parameters by up to 98.6% without compromising segmentation performance, achieving results comparable to U-Net and its variants.

Numerous studies have focused on enhancing segmentation accuracy through the use of deep learning. One common approach is the incorporation of attention mechanisms. For instance, [13] proposed the Attention U-Net, which replaces hard attention with soft attention and embeds it into the skip connections and upsampling paths. This strategy enables the model to focus on relevant local features while suppressing irrelevant background noise. Similarly, [14] introduced AGU-Net, which adds attention modules at the bridge laver of the U-Net to enhance feature extraction in the segmentation process. In another study, [15] proposed the use of pixel attention and channel attention modules arranged separately. Channel attention excels at encoding global information, while pixel attention captures localized information at the pixel level. [16] further improved this approach by arranging pixel and channel attention in parallel to combine global and local information at the pixel level effectively.

Another technique often employed to improve accuracy is multi-scale feature extraction. By

processing information at different scales, models can capture global context and fine details more effectively. [17] proposed a Multi-scale Depthwise Separable Convolution architecture that extracts multi-scale features while maintaining model efficiency through depthwise separable convolutions. Likewise, [18] introduced the Multi-scale Residual Attention (MRA) module, which combines multi-scale processing, residual connections, and attention mechanisms to capture and exploit features at different scales without suffering from degradation due to increased network depth.

Building upon the insights from the aforementioned studies, we propose a novel architecture, One-Point-Five U-Net, which aims to improve segmentation accuracy while maintaining computational efficiency. This architecture integrates the strengths of U-Net and Half U-Net, with the following modifications: (1) Replacing standard convolutions in the encoder with ghost modules as used in Half U-Net; (2) Inserting MRA modules between ghost modules in the encoder to enhance multi-scale attention-based feature extraction; (3) Incorporating Enhanced Parallel Attention (EPA) modules in the skip connections to emphasize essential features during feature fusion; (4) Adding a Convolutional Block Attention Module (CBAM) in the bridge layer to further refine salient features; and (5) Replacing standard convolutions in the second decoder with Multi-scale Depthwise Separable Convolutions (MSDC) to maintain efficiency while capturing multi-scale information. Through these innovations, One-Point-Five U-Net is expected to deliver accurate and robust medical image segmentation, providing effective support for clinical disease diagnosis and prevention.

II. Methodology

This research consists of six main stages: (1) collecting and preprocessing the dataset, (2) designing the neural network architecture, (3) performing hyperparameter tuning, (4) training the model, (5) evaluating the model, and (6) comparing the results with previous methods.

A. Dataset

This study utilizes four publicly available medical imaging datasets for segmentation tasks. The CVC-ClinicDB dataset includes 612 colonoscopy images with polyp annotations, each with a resolution of 288 × 368 pixels [19]. The Kvasir-SEG dataset contains 1,000 endoscopic images with segmentation masks and varying resolutions [20]. The Breast Ultrasonic Images (BUSI) dataset comprises 780 breast ultrasound images [21]. However, only 647 images labelled as benign or malignant were used, excluding the normal class due to the absence of regions of interest. Lastly, the ISIC2018 dataset comprises 3,694 dermoscopic

Manuscript received 10 May 2025; Revised 21 June 2025; Accepted 30 June 2025; Available online 5 July 2025 Digital Object Identifier (**DOI**): https://doi.org/10.35882/jeeemi.v7i3.949

Copyright © 2025 by the authors. This work is an open-access article and licensed under a Creative Commons Attribution-ShareAlike 4.0 International License (CC BY-SA 4.0).

skin lesion images with corresponding masks, divided into training, validation, and test sets, with a resolution of 1022 × 767 pixels [22].



Fig. 1. Image augmentation results on CVC-ClinicDB dataset (a) original, (b) vertical flip, (c) horizontal flip, (d) random rotation, (e) random brightness, and (f) grid distortion.

After collecting the datasets, all images were resized to 256 × 256 pixels to balance memory efficiency. Data augmentation was applied to the training sets of the CVC-ClinicDB, Kvasir-SEG, and BUSI datasets to increase data variability and reduce overfitting [23]. For the CVC-ClinicDB and KvasirSEG, the augmentation techniques included vertical flip, horizontal flip, random rotation, random brightness adjustment, and grid distortion. While for BUSI, the techniques included 90° rotation, 180° rotation, vertical flip, and horizontal flip. As a result of these augmentation strategies, the training sets increased to 2,940 images for CVC-ClinicDB, 4,800 images for Kvasir-SEG, and 3,108 images for BUSI. Examples of augmented images from the CVC-ClinicDB dataset are shown in Fig. 1.

B. Neural Network Architecture Design

In this study, we propose a novel deep learning architecture called One-Point-Five U-Net, a modified version of the standard U-Net framework. This architecture integrates two types of decoders. The first decoder retains the basic concept of U-Net, in which each level applies a convolution operation followed by an upsampling process [11]. In contrast, the second decoder employs the Half-UNet decoder structure, which utilizes full-scale feature fusion to merge encoder and decoder features. This full-scale fusion directly combines features from the encoder and decoder without requiring memory-intensive concatenation operations. The architecture of One-Point-Five U-Net is depicted in Fig. 2.

Inspired by the Half-UNet design, the encoder and decoder in the proposed model are designed with an equal number of filters at each level. Furthermore, standard convolutions in the encoder are replaced with Ghost Modules to reduce computational cost compared to basic convolution operations [24]. Each Ghost Module consists of a pointwise convolution followed by a depthwise convolution, referred to as a "cheap operation." The structure of the Ghost Module is shown in Fig. 3(a).



Fig. 2. Overview of the proposed Method, One-Point-Five U-Net. This architecture comprises five main components: Ghost Module, MRA, EPA, CBAM, and MSDC.

Manuscript received 10 May 2025; Revised 21 June 2025; Accepted 30 June 2025; Available online 5 July 2025 Digital Object Identifier (**DOI**): https://doi.org/10.35882/jeeemi.v7i3.949



Fig. 3. Architecture of (a) Ghost Module and (b) Multi-Scale Residual Attention (MRA).

To further enhance model performance, we introduce a Multi-scale Residual Attention (MRA) module between the two Ghost Modules. MRA combines three mechanisms: multi-scale processing, attention, and residual connections, as depicted in Fig. 3 (b). Initially, a 1×1 convolution reduces the input feature map dimension $H \times W \times C/8$ (X₁) as shown in Eq. (1). Then, three parallel convolutional paths implement the multi-scale concept. The first path (X_2) uses a standard 3×3 convolution formulated in Eq. (2). The second path (X_3) follows the InceptionV2 approach by decomposing a 5×5 convolution into two consecutive 3×3 convolutions to reduce parameter count, as described in Eq. (3). Similarly, the third path (X_4) decomposes a 7×7 convolution into three consecutive 3×3 convolutions as shown in Eq. (4). The outputs of all three paths (X_5) are then merged to integrate information across multiple receptive fields like Eq. (5).

The attention mechanism in MRA utilizes a Squeeze-and-Excitation (SE) block to emphasize relevant features after concatenation (X_6) , formulated in Eq. (6). A 3×3 convolution is applied to restore the feature map dimensions to match the input size (X_7) . Eq. (7). The result is then added to the original input (X) features, following the concept of residual

connections as formulated in Eq. (8). This residual addition helps mitigate gradient vanishing and information degradation caused by convolutional operations.

$$X_1 = Conv_{1x1}(X) \tag{1}$$

$$X_2 = Conv_{3x3}(X_1) \tag{2}$$

$$X_3 = \operatorname{Conv}_{3\times 3} \left(\operatorname{Conv}_{3\times 3}(X_1) \right) \tag{3}$$

$$X_4 = Conv_{3x3} \left(Conv_{3x3} \left(Conv_{3x3} (X_1) \right) \right)$$
(4)

$$X_5 = [X_2, X_3, X_4]$$
(5)

$$X_6 = SE(X_5) \tag{6}$$

$$X_7 = COnV_{3X3}(X_6)$$
 (1)

$$\hat{U}_{MRA} = Conv_{1x1}(X) + X_7 \tag{8}$$

We experiment with three types of SE blocks to determine the optimal configuration: Spatial and Channel Squeeze & Excitation (scSE), Channel Squeeze and Spatial Excitation (sSE), and Spatial Squeeze and Channel Excitation (cSE). The scSE block integrates the functions of both cSE and sSE by recalibrating feature maps across both spatial and channel dimensions simultaneously, which can be seen in Fig. 4 (a) [25]. The outputs of cSE and sSE are

- Comm



Manuscript received 10 May 2025; Revised 21 June 2025; Accepted 30 June 2025; Available online 5 July 2025 Digital Object Identifier (DOI): https://doi.org/10.35882/jeeemi.v7i3.949



Fig. 5. Architecture of (a) Enhanced Paralel Attention (EPA), (b) Simple Pixel Attention (SPA), and Pixel Attention (PA).

combined through element-wise addition. This schema is formulated by Eq. (9).

$$\widehat{U}_{scSE} = \widehat{U}_{cSE} + \widehat{U}_{sSE} \tag{9}$$

In the cSE block, spatial information is reduced using Global Average Pooling (GAP), generating a global descriptor for each channel. The excitation process then uses two fully connected layers, followed by ReLU (δ) and sigmoid activations (σ), respectively. The resulting attention weights are multiplied element-wise by the input feature map to recalibrate the importance of each channel. The structure of the cSE block is illustrated in Fig. 4 (b). cSE can be formulated in Eqs. (10), (11), and (12).

$$z_{k} = \frac{1}{H \times W} \sum_{i}^{H} \sum_{j}^{W} u_{k}(i, j),$$
(10)

$$\hat{z} = \sigma \big(W1 \cdot \delta (W2 \cdot z) \big), \tag{11}$$

$$\widehat{U}_{cSE} = \left[\sigma(\widehat{z}_1) \cdot u_1, \sigma(\widehat{z}_2) \cdot u_2, \dots, \sigma(\widehat{z}_C) \cdot u_C\right].$$
⁽¹²⁾

where *z* represents the global average pooled values and *k* denotes the channel index, so z_k represents the average value of the *k*-th channel and $H \times W$ denotes the spatial dimensions of the feature map. *W*1 and *W*2 are the weights of the fully-connected layers and *C* represents the number of channels. The variable (i, j)indicates the spatial coordinates for each pixel in the height (i) and width (j) dimensions of the feature map.

In the sSE block, channel-based features are generated by applying a 1×1 convolution at each spatial location, followed by a sigmoid activation function (σ). This operation is defined in Eqs. (13) and (14), where *W* represents the convolution weights and *U* is the input feature map. These features are then spatially scaled and recalibrated to highlight relevant regions, formulated by Eq. (15). The structure of the sSE block is illustrated in Fig. 4(c).

$$q = W_{sq} * U , \qquad (13)$$

$$\hat{q}_{ij} = \sigma(q_{ij}), \tag{14}$$

$$\widehat{U}_{SSE} = \left[\widehat{q}_{1,1} \cdot \widehat{u}_{1,1}, \dots, \widehat{q}_{1,j} \cdot \widehat{u}_{1,j}, \dots, \widehat{q}_{H,W} \cdot \widehat{u}_{H,W}, \right].$$
(15)

This architecture integrates an Enhanced Parallel Attention (EPA) module into the skip connections. Placing the EPA here allows the model to filter and selectively emphasize salient spatial features from the encoder. EPA combines multiple attention mechanisms in parallel enhance model to performance. It comprises three main components: Simple Pixel Attention (SPA), Channel Attention (CA), and Pixel Attention (PA). EPA is depicted in Fig. 5(a). The outputs of these attention modules are concatenated along the channel dimension and processed through a multi-layer perceptron (MLP). This MLP reduces the channel dimension to match that of the input. The final output is added to a shortcut identity to retain the original information. The EPA is mathematically formulated in Eq. (16).

$$\hat{U}_{EPA} = x \oplus Conv_{1x1} \left(\delta \left(Conv_{1x1} \left(Cat(\hat{U}_{SPA}, \hat{U}_{PA}, \hat{U}_{CA}) \right) \right) \right)$$
(16)

The SPA module is designed to extract locationdependent features such as texture or intensity variations in medical images. As depicted in Fig. 5 (b), it consists of two branches: PF_s , which extracts features using point-wise (*Conv*_{1x1}) and 3×3 convolution layers as shown in Eq. (17), and PA_s, which generates a pixel-level gating signal using a sigmoid function formulated in Eq. (18). The output of SPA is obtained through element-wise multiplication of these two branches. It can be calculated using Eq. (19).

$$PF_s = Conv(Conv_{1x1}(x)), \tag{17}$$

$$PA_s = \sigma(Conv_{1x1}(x)), \tag{18}$$

$$\widehat{U}_{PA} = PF_s \otimes PA_s \tag{19}$$

The PA module, illustrated in Fig. 5 (c), captures finegrained information at each pixel location in the feature maps. It employs two sequential 1×1 convolution layers, followed by ReLU (δ) and sigmoid activation functions (σ), respectively. These layers generate pixel-level attention maps, enabling the model to focus on spatially relevant features. Pixel Attention is calculated using Eq. (20).

Manuscript received 10 May 2025; Revised 21 June 2025; Accepted 30 June 2025; Available online 5 July 2025 Digital Object Identifier (**DOI**): https://doi.org/10.35882/jeeemi.v7i3.949

$$\widehat{U}_{PA} = x \otimes \sigma \left(Conv_{1x1} \left(\delta \left(Conv_{1x1}(x) \right) \right) \right)$$
(20)

The CA module extracts global contextual information and adjusts the importance of feature channels. It uses both GAP and GMP to encode distinctive patterns from the input (*x*). These pooled representations are passed through shared fully connected layers, followed by a sigmoid activation (σ) to generate channel-wise attention weights. CA is mathematically defined in Eq. (21).

$$\widehat{U}_{CA} = x \otimes \sigma \left(Conv_{1x1} \left(Act \left(Conv_{1x1} \left(GAP(x) \right) \right) \right) \right)$$
(21)

At the bridge layer of the network, a Convolutional Block Attention Module (CBAM) is added. Placing CBAM at the network's deepest point allows the model to enhance the global context before the decoder begins the spatial reconstruction process. CBAM combines channel and spatial attention mechanisms. While channel attention emphasizes the relevance of each feature channel, spatial attention focuses on highlighting significant spatial regions within the feature map [26]. The structure of CBAM is depicted in Fig. 6.



Fig. 6. Structure of the CBAM.



Fig. 7. Architecture of MSDC.

The second decoder substitutes basic convolutions with a Multi-scale Depthwise Convolution (MSDC). MSDC extracts features using three parallel depthwise convolutions with different kernel sizes, as illustrated in Fig. 7. Each convolution is followed by batch normalization and an activation function to ensure stable and nonlinear representations. The outputs of the three branches are then combined through element-wise addition to integrate multi-scale spatial information. This approach enables the model to capture diverse spatial features in the input image effectively. The final layer of the neural network applies a 1×1 convolution with a single filter and a sigmoid activation function to generate the binary prediction map. The 1×1 kernel is used to minimize the number of parameters. A single filter produces a grayscale output, and the sigmoid function maps pixel values to the [0, 1] range. Pixels above 0.5 are classified as foreground, while those below are treated as background.

C. Hyperparameter Tuning

In this study, hyperparameter tuning was conducted using the training set from the CVC-ClinicDB dataset. Hyperparameter tuning is identifying the optimal values of hyperparameters to achieve the best model performance [27]. The hyperparameters tuned in this study include the number of filters, activation functions, bias values, loss functions, and types of SE blocks.

Experiments on the number of filters were performed with three configurations: 16, 32, and 64 filters. Next, different activation functions were evaluated, including ReLU, LeakyReLU, ELU, SELU, and GELU. Subsequently, various loss functions were tested, including Binary Cross-Entropy (BCE), Dice Loss (DL), a combination of BCE and DL, Tversky Loss, and Focal Tversky Loss. Lastly, experiments were conducted on the type of Squeeze-and-Excitation (SE) block used in the architecture. The SE block variants evaluated in this study were channel SE (cSE), spatial SE (sSE), and spatial and channel SE (scSE).

D. Model Training

After completing the hyperparameter tuning process, the neural network was trained using the designated training set, where the input images served as inputs and the target images as ground truth labels. A learning rate scheduling strategy was applied, in which the learning rate was reduced by a factor of 10 if the loss value did not improve for 10 consecutive epochs. This strategy aimed to enhance model performance and prevent overfitting.

E. Model Evaluation

The trained model was subsequently evaluated using the testing set from each dataset to assess its performance in segmentation. This evaluation involved comparing the predicted image generated by the model with the ground-truth. Five standard performance metrics were employed to achieve a comprehensive assessment: sensitivity, specificity, accuracy, Dice similarity coefficient (DSC), and Intersection over Union (IoU). Sensitivity, also known as recall, measures the model's ability to identify positive pixels accurately and is calculated using Eq. (22).

Manuscript received 10 May 2025; Revised 21 June 2025; Accepted 30 June 2025; Available online 5 July 2025 Digital Object Identifier (**DOI**): https://doi.org/10.35882/jeeemi.v7i3.949

$$SE = \frac{TP}{TP + FN}$$
(22)

where TP represents true positives and FN denotes false negatives. Specificity, defined in Eq. (23), measures the model's ability to identify negative pixels correctly.

$$SP = \frac{TN}{TN + FP}$$
(23)

where TN is the count of true negatives and FP is the number of false positives.

Overall classification accuracy is computed using Eq. (24).

$$ACC = \frac{TP + TN}{TP + FP + FN + TN}$$
(24)

which measures the proportion of correctly predicted pixels across both classes. The DSC, provided in Eq. (25), is the harmonic mean of precision and recall, emphasizing the balance between FP and FN.

$$DSC = \frac{2TP}{2TP + FP + FN}$$
(25)

In segmentation tasks, the IoU, as defined in Eq. (26), quantifies the overlap between the predicted and ground-truth masks.

$$IoU = \frac{TP}{TP + FP + FN}$$
(26)

_ . .

- - -

The model was also analyzed based on the number of parameters. The total parameter count reflects the model's memory efficiency and structural complexity. A model with fewer parameters indicates a lightweight design.

III. RESULTS

A. Experimental Setup

In this stage, a series of comprehensive experiments were conducted to evaluate the performance of the proposed One-Point-Half U-Net model in medical image segmentation tasks. All experiments were implemented using the Python programming language and the TensorFlow framework, executed on the Kaggle platform with TPU v3-8 acceleration. The CVC-ClinicDB dataset was used for both hyperparameter tuning and model ablation experiments to identify the optimal configuration for model construction. Once the optimal configuration was obtained, the model was trained and evaluated using the training and testing sets from each dataset.

B. Hyperparameter Tuning

The initial set of experiments involved hyperparameter tuning on the number of filters, the type of activation function, the loss function, and the type of Squeeze-and-Excitation (SE) block utilized. Additionally, the tuning process was conducted sequentially. Table 2 compares the initial hyperparameters and optimal hyperparameters obtained as the tuning process progresses through various stages.

Table 2.	Comparison	between	initial a	and op	timal
hyperpa	arameters.				

Hyperparameter	Initial	Optimal
Batch Size	16	16
Number of filters	64	64
Number of Epochs	60	60
Optimizer	Adam (Adaptive Moment Estimation)	Adam (Adaptive Moment Estimation)
Learning Rate	0.001 (divided by 10 every 10 epochs if validation loss doesn't improve)	0.001 (divided by 10 every 10 epochs if validation loss doesn't improve)
Activation Function	ReLU	GELU
Loss Function	DL	DL + BCE
Type of SE block	cSE	cSE

The initial experiment aimed to determine the optimal number of convolutional filters that balances model performance and complexity. We evaluated three configurations with initial filter counts of 16, 32, and 64 on the CVC-ClinicDB dataset. The results, presented in Table 1, indicate that increasing the number of filters

Table 1. Hyperparameter tuning	results.
--------------------------------	----------

Hyperparameter		SE	SP	ACC	DSC	loU	Parameter
Filters	16	83,87%	97,92%	94,70%	87.89%	78.40%	68,183
	32	93,04%	99,42%	98,73%	93.40%	87.81%	195,295
	64	94,66%	99,59%	99,08%	94.78%	90.22%	626,755
	ReLU	93.71%	99.60%	99.00%	94.54%	89.79%	626,755
	Leaky ReLU	93.61%	99.49%	99.09%	94.64%	89.98%	626,755
Activation Function	ELU	94.74%	99.49%	99.07%	94.62%	89.96%	626,755
	SELU	95.33%	99.47%	99.09%	94.64%	90.03%	626,755
	GELU	94.66%	99.59%	99.08%	94.78%	90.22%	626,755
	BCE	94.33%	99.56%	99.14%	94.86%	90.31%	626,755
	DL	94.66%	99.59%	99.08%	94.78%	90.22%	626,755
Loss Function	BCE+DL	94.89%	99.63%	99.23%	95.41%	91.27%	626,755
	TL	94.52%	99.61%	99.13%	94.82%	90.30%	626,755
	FTL	94.33%	99.65%	99.26%	95.26%	91.02%	626,755
	cSE	94.89%	99.63%	99.23%	95.41%	91.27%	626,755
SEBlock	sSE	95.39%	99.40%	99.08%	94.97%	90.49%	625,780
	scSE	95.04%	99.61%	99.30%	95.36%	91.19%	626,880

Manuscript received 10 May 2025; Revised 21 June 2025; Accepted 30 June 2025; Available online 5 July 2025 Digital Object Identifier (**DOI**): https://doi.org/10.35882/jeeemi.v7i3.949



Fig. 8. Training and Validation Loss Graph on (a) CVC-ClinicDB, (b) KvasirSeg, (c) BUSI, and (d) ISIC2018.

Method	SE	SP	ACC	DSC	loU	Parameter
Ablation 1	94.33%	99.41%	98.92%	93.92%	88.61%	520,175
Ablation 2	94.00%	99.32%	98.76%	93.17%	87.50%	320,499
Ablation 3	95.05%	99.53%	99.09%	94.61%	90.15%	625,421
Ablation 4	94.12%	99.40%	98.78%	93.50%	88.35%	549,495
Ablation 5	82.63%	96.76%	94.99%	80.87%	68.24%	123,329
One-Point-Five U-Net	94.89%	99.63%	99.23%	95.41%	91.27%	626,755

 Table 3. Result of the ablation model.

improves segmentation accuracy but also increases the model's total number of trainable parameters. The configuration with 64 filters yielded the highest accuracy among the tested options. Therefore, this value was selected and fixed for all subsequent experiments to ensure a robust performance baseline.

The second experiment aimed to identify the most effective activation function for the network's non-linear transformations. Keeping the filter count at 64, we tested five different activation functions: ReLU, LeakyReLU, ELU, SELU, and GELU. As summarized in Table 1, the choice of activation function had a significant impact on the model's accuracy, while it did not alter the number of parameters. Among the tested functions, GELU demonstrated superior performance by yielding the highest segmentation accuracy. Consequently, GELU was selected as the default activation function for the remainder of the tuning process.

The third experiment was designed to determine the most effective loss function for optimizing the model during training, which is crucial for addressing the class imbalance commonly encountered in medical segmentation. We evaluated five different loss functions: Binary Cross-Entropy (BCE), Dice Loss (DL), a combined BCE + DL, Tversky Loss (TL), and Focal Tversky Loss (FTL). The results in Table 1 indicate that the choice of loss function has a significant impact on the final segmentation accuracy. The combined BCE + DL loss function achieved the highest accuracy scores, proving most effective for this task.

The final hyperparameter experiment was conducted to determine the optimal Squeeze-and-Excitation (SE) block variant for our architecture. Three different types were evaluated: scSE, sSE, and cSE. Based on the results presented in Table 1,the type of SE block influenced both the model's accuracy and the number of parameters. The cSE block variant achieved the highest segmentation accuracy. As a result, the cSE block was adopted as the definitive attention mechanism in the final proposed model.

C. Model Ablation

Ablation studies were conducted to evaluate the individual contributions of specific components within a model by systematically removing or modifying them. The proposed One-Point-Five U-Net integrates four core modules: MRA, CBAM, EPA, and MSDC. A series of ablation experiments were designed to isolate the influence of each module and explore their interactions. The conducted experiments included the following configurations: (1) Ablation 1: One-Point-Five U-Net w/o MRA, (2) Ablation 2: One-Point-Five U-Net w/o CBAM, (3) Ablation 3: One-Point-Five U-Net w/o EPA, (4) Ablation 4: One-Point-Five U-Net with MSDC replaced by a Ghost module in the decoder, (5) Ablation 5: One-Point-Five w/o MRA, EPA, CBAM, and replacing MSDC with a ghost module in the decoder.

Each ablation was evaluated to determine the performance drop or improvement compared to the whole model. The results of these ablation experiments are summarized in Table 3. These findings validate that the incorporation of MRA, CBAM, EPA, and MSDC modules significantly enhances the model's segmentation capability. Based on the experimental results, the full version of the One-Point-Five U-Net model demonstrated the best performance, achieving a

Manuscript received 10 May 2025; Revised 21 June 2025; Accepted 30 June 2025; Available online 5 July 2025 Digital Object Identifier (**DOI**): https://doi.org/10.35882/jeeemi.v7i3.949

Copyright © 2025 by the authors. This work is an open-access article and licensed under a Creative Commons Attribution-ShareAlike 4.0 International License (CC BY-SA 4.0).

Dataset	SE	SP	ACC	DSC	loU	Param	FLOPs(G)
CVC-ClinicDB	94.89%	99.63%	99.23%	95.41%	91.27%	626,755	28.23
KvasirSeg	91.11%	98.60%	97.33%	90.93%	83.84%		
BUSI	85.35%	98.65%	96.81%	87.02%	78.18%		
ISIC2018	87.67%	98.11%	93.68%	89.27%	83.06%		

Table 4. Evaluation results on all four datasets.



DSC of 95.41%, an IoU of 91.27%, and a total parameter count of 626.755. When the MRA module was removed. the model's performance declined significantly, with the DSC dropping to 93.92% and the IoU to 88.61%. This finding underscores the crucial role of the MRA module in enhancing segmentation accuracy. Furthermore, the simultaneous removal of multiple modules, specifically MRA, EPA, CBAM, and MSDC, led to a drastic performance decline, yielding only a DSC of 80.87% and an IoU of 68.24%, despite reducing the number of parameters to just 123,329. These results confirm that integrating attention mechanisms and multi-scale convolutional modules is crucial to the model's success. The absence of one or more of these modules consistently results in reduced accuracy, even reducing the model's complexity.

D. Experiment Result

After performing hyperparameter tuning and model ablation, the optimal hyperparameter configuration and the most effective model architecture were identified. The most optimal hyperparameters, as shown in Table 1, were then applied during the training phase on the CVC-ClinicDB, KvasirSeg, BUSI, and ISIC2018 datasets.. Throughout training, both training and validation losses were continuously monitored to assess the convergence behavior of the model and to detect signs of overfitting or underfitting. The loss curves depicting this process are presented in Fig. 8. In the subsequent phase, the model was evaluated using the corresponding testing sets of all four datasets to assess the generalization ability of the trained model. The evaluation metrics obtained from the testing process are summarized in Table 4. Additionally, a qualitative comparison between the predicted segmentation masks and the ground truth annotations is provided in Table 5.

E. Comparison with Previous Methods

Following the series of evaluations, the proposed neural network was compared with several existing segmentation methods, including U-Net [11], U-Net++ [28], ResU-Net [29], Half-UNet [12], Attention U-Net [13], and CMAUNeXT [30]. The comparison was conducted based on four key aspects: DSC, IoU, number of parameters, and floating-point operations (FLOPs).

The results of comparing these metrics are presented in Table 6. The results show that the proposed method, One-Point-Five U-Net, outperforms existing models in terms of DSC and IoU, indicating higher accuracy in generating segmentation predictions. Additionally, although this model is still less efficient than Half-UNet in terms of the number of parameters and FLOPs, the values obtained are still relatively low compared to other methods. This signifies that the proposed model strikes a superior balance between high accuracy and computational efficiency.

Manuscript received 10 May 2025; Revised 21 June 2025; Accepted 30 June 2025; Available online 5 July 2025 Digital Object Identifier (**DOI**): https://doi.org/10.35882/jeeemi.v7i3.949

Mothod	CVC-C	linicDB	Kvas	irSeg	BL	ISI ISIC20		2018	Doromotor	
Method	DSC	loU	DSC	loU	DSC	loU	DSC	loU	U	
U-Net [11]	89.95%	83.12%	80.58%	71.28%	77.19%	67.97%	80.41%	70.69%	31.04 M	96.48
U-Net++ [28]	85.76%	77.92%	77.08%	67.24%	72.28%	62.68%	73.32%	62.65%	19.92 M	203.21
ResU-Net [29]	87.50%	81.25%	83.61%	75.30%	78.27%	69.25%	82.53%	72.43%	32.45 M	101.94
Half-UNet [12]	89.50%	83.41%	82.40%	73.94%	79.11%	69.89%	85.64%	77.58%	221,729	9.86
Attention U-Net [13]	88.74%	81.34%	80.65%	70.96%	71.11%	60.83%	78.40%	68.20%	8.13 M	91.30
CMAUNeXT [30]	93.05%	87.44%	81.24%	71.25%	85.59%	76.67%	85.63%	76.86%	2.89 M	13.17
Proposed Method	95.41%	91.27%	90.93%	83.84%	87.02%	78.18%	90.03%	82.40%	626,755	28.23

 Table 6. Comparison with previous methods.

IV. CONCLUSION

This study introduces a novel neural network architecture for medical image segmentation, demonstrating superior performance compared to existing state-of-the-art methods. The proposed One-Point-Five architecture, named U-Net, is constructed by combining the foundational structures of U-Net and Half-UNet while integrating several advanced modules, including the MRA, EPA, CBAM, and MSDC. The integration of attention mechanisms and multi-scale feature extraction enables the One-Point-Five U-Net to enhances segmentation performance while maintaining high computational efficiency. The One-Point-Five U-Net achieves consistently high results across four benchmark datasets. Specifically, it obtains sensitivity. specificity, accuracy, DSC, and IoU of 94.89%, 99.63%, 99.23%, 95.41%, and 91.27% on the CVC-ClinicDB dataset; 91.11%, 98.60%, 97.33%, 90.93%, 83.84% on the Kvasir-Seg dataset; 85.35%, 98.65%, 96.81%, 87.02%. 78.18% on the BUSI dataset: and 87.67%. 98.11%, 93.68%, 89.27%, 83.06% on the ISIC2018 dataset. Furthermore, the One-Point-Five U-Net outperforms several baseline architectures, including U-Net, U-Net++, ResU-Net, Half-UNet, Attention U-Net, and CMAUNeXT based on DSC and Intersection over Union (IoU) metrics. In terms of model complexity, the One-Point-Five U-Net contains 626,755 parameters and requires 28.23 GFLOPs, which is relatively lightweight compared to conventional models. This computational efficiency indicates the model's practical applicability in clinical scenarios with limited hardware resources, without compromising accuracy in computer-aided diagnosis (CAD) systems.

Acknowledgment

This work was fully supported and funded by the RKAT Universitas Sebelas Maret (UNS) of the year 2025 under the research grant HIBAH KOLABORASI MITRASMART with the contract 372/UN27.22/PT.01.03/2025.

References

[1] J. Zhang et al., "Advances in attention mechanisms for medical image segmentation,"

Comput. Sci. Rev., vol. 56, p. 100721, May 2025, doi: 10.1016/j.cosrev.2024.100721.

- [2] X. Shu, J. Wang, A. Zhang, J. Shi, and X.-J. Wu, "CSCA U-Net: A channel and space compound attention CNN for medical image segmentation," Artif. Intell. Med., vol. 150, p. 102800, Apr. 2024, doi: 10.1016/j.artmed.2024.102800.
- [3] M. E. Rayed, S. M. S. Islam, S. I. Niha, J. R. Jim, M. M. Kabir, and M. F. Mridha, "Deep learning for medical image segmentation: State-of-the-art advancements and challenges," Informatics Med. Unlocked, vol. 47, p. 101504, 2024, doi: 10.1016/j.imu.2024.101504.
- [4] R. Wang, T. Lei, R. Cui, B. Zhang, H. Meng, and A. K. Nandi, "Medical image segmentation using deep learning: A survey," IET Image Process., vol. 16, no. 5, pp. 1243–1267, Apr. 2022, doi: 10.1049/ipr2.12419.
- [5] X. Ding, K. Qian, Q. Zhang, X. Jiang, and L. Dong, "Dual-channel compression mapping network with fused attention mechanism for medical image segmentation," Sci. Rep., vol. 15, no. 1, p. 8906, Mar. 2025, doi: 10.1038/s41598-025-93494-4.
- [6] H. Prasetyo, M. A. F. Rohman, A. W. H. Prayuda, and J.-M. Guo, "Enhancing Polyp Segmentation Efficiency Using Pixel Channel Attention HalfU-Net," in 2024 IEEE 10th Information Technology International Seminar (ITIS), IEEE, Nov. 2024, pp. 381–386. doi: 10.1109/ITIS64716.2024.10845590.
- Y. Zhang, Q. Liao, L. Ding, and J. Zhang, "Bridging 2D and 3D segmentation networks for computation-efficient volumetric medical image segmentation: An empirical study of 2.5D solutions," Comput. Med. Imaging Graph., vol. 99, p. 102088, Jul. 2022, doi: 10.1016/j.compmedimag.2022.102088.
- [8] X. Wu, S. Huang, X. Shu, C. Hu, and X.-J. Wu, "MPFC-Net: A multi-perspective feature compensation network for medical image segmentation," Expert Syst. Appl., vol. 248, p. 123430, Aug. 2024, doi: 10.1016/j.eswa.2024.123430.

Manuscript received 10 May 2025; Revised 21 June 2025; Accepted 30 June 2025; Available online 5 July 2025 Digital Object Identifier (**DOI**): https://doi.org/10.35882/jeeemi.v7i3.949

Journal of Electronics, Electromedical Engineering, and Medical Informatics Homepage: jeeemi.org; Vol. 7, No. 3, July 2025, pp: 869-880 e-ISSN: 2656-8632

- [9] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, Jun. 2015, pp. 3431–3440. doi: 10.1109/CVPR.2015.7298965.
- [10] T. Zhang, Y. Liu, Y. Fan, and M. Lu, "Improvement of park drivable area segmentation method based on STDCSeg network," Discov. Appl. Sci., vol. 7, no. 4, p. 297, Apr. 2025, doi: 10.1007/s42452-025-06767-y.
- [11] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," 2015, pp. 234–241. doi: 10.1007/978-3-319-24574-4_28.
- [12] H. Lu, Y. She, J. Tie, and S. Xu, "Half-UNet: A Simplified U-Net Architecture for Medical Image Segmentation," Front. Neuroinform., vol. 16, Jun. 2022, doi: 10.3389/fninf.2022.911679.
- [13] O. Oktay et al., "Attention U-Net: Learning Where to Look for the Pancreas," Apr. 2018, doi: https://doi.org/10.48550/arXiv.1804.03999.
- [14] M. A. F. Rohman, H. M. Akbar, A. D. A. Firdaus, and H. Prasetyo, "AGU-NET:Attention Ghost U-NetUntuk Segmentasi Penyakit Polip Berbasis Citra Biomedis," Bul. PAGELARAN Mhs. Nas. Bid. Teknol. Inf. DAN Komun., vol. 1, no. 1, pp. 44–49, 2023.
- [15] X. Qin, Z. Wang, Y. Bai, X. Xie, and H. Jia, "FFA-Net: Feature Fusion Attention Network for Single Image Dehazing," Proc. AAAI Conf. Artif. Intell., vol. 34, no. 07, pp. 11908–11915, Apr. 2020, doi: 10.1609/aaai.v34i07.6865.
- [16] L. Lu, Q. Xiong, B. Xu, and D. Chu, "MixDehazeNet: Mix Structure Block For Image Dehazing Network," in 2024 International Joint Conference on Neural Networks (IJCNN), IEEE, Jun. 2024, pp. 1–10. doi: 10.1109/IJCNN60899.2024.10651326.
- [17] Y. Dai, C. Li, X. Su, H. Liu, and J. Li, "Multi-Scale Depthwise Separable Convolution for Semantic Segmentation in Street–Road Scenes," Remote Sens., vol. 15, no. 10, p. 2649, May 2023, doi: 10.3390/rs15102649.
- [18] X. Shu, X. Li, X. Zhang, C. Shao, X. Yan, and S. Huang, "MRAU-net: Multi-scale residual attention U-shaped network for medical image segmentation," Comput. Electr. Eng., vol. 118, p. 109479, Sep. 2024, doi: 10.1016/j.compeleceng.2024.109479.
- [19] J. Bernal, F. J. Sánchez, G. Fernández-Esparrach, D. Gil, C. Rodríguez, and F. Vilariño, "WM-DOVA maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from

physicians," Comput. Med. Imaging Graph., vol. 43, pp. 99–111, Jul. 2015, doi: 10.1016/j.compmedimag.2015.02.007.

- [20] D. Jha et al., "Kvasir-SEG: A Segmented Polyp Dataset," 2020, pp. 451–462. doi: 10.1007/978-3-030-37734-2_37.
- [21] W. Al-Dhabyani, M. Gomaa, H. Khaled, and A. Fahmy, "Dataset of breast ultrasound images," Data Br., vol. 28, p. 104863, Feb. 2020, doi: 10.1016/j.dib.2019.104863.
- [22] N. Codella et al., "Skin Lesion Analysis Toward Melanoma Detection 2018: A Challenge Hosted by the International Skin Imaging Collaboration (ISIC)," Feb. 2019, [Online]. Available: http://arxiv.org/abs/1902.03368
- [23] M. A. Fathur Rohman, H. Prasetyo, H. M. Akbar, and A. D. Afan Firdaus, "ACMU-Net: An Efficient Architecture Based on ConvMixer and Attention Mechanism for Colorectal Polyp Segmentation," in 2024 IEEE International Conference on Smart Mechatronics (ICSMech), IEEE, Nov. 2024, pp. 279–284. doi: 40.4400/JCCMachCOD20.0004.40040000

10.1109/ICSMech62936.2024.10812309.

- [24] K. Han, Y. Wang, Q. Tian, J. Guo, C. Xu, and C. Xu, "GhostNet: More Features From Cheap Operations," in 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, Jun. 2020, pp. 1577–1586. doi: 10.1109/CVPR42600.2020.00165.
- [25] A. G. Roy, N. Navab, and C. Wachinger, "Concurrent Spatial and Channel 'Squeeze & amp; Excitation' in Fully Convolutional Networks," 2018, pp. 421–429. doi: 10.1007/978-3-030-00928-1_48.
- [26] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "CBAM: Convolutional Block Attention Module," 2018, pp. 3–19. doi: 10.1007/978-3-030-01234-2_1.
- [27] R. Andonie, "Hyperparameter optimization in learning systems," J. Membr. Comput., vol. 1, no. 4, pp. 279–291, Dec. 2019, doi: 10.1007/s41965-019-00023-0.
- [28] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang, "UNet++: A Nested U-Net Architecture for Medical Image Segmentation," Jul. 2018, [Online]. Available: http://arxiv.org/abs/1807.10165
- [29] X. Xiao, S. Lian, Z. Luo, and S. Li, "Weighted Res-UNet for High-Quality Retina Vessel Segmentation," in 2018 9th International Conference on Information Technology in Medicine and Education (ITME), IEEE, Oct. 2018, pp. 327–331. doi: 10.1109/ITME.2018.00080.
- [30] H. Prasetyo, R. B. Ashidiqy, and U. Salamah, "CMAUNeXt: An Efficient Neural Network Based

Manuscript received 10 May 2025; Revised 21 June 2025; Accepted 30 June 2025; Available online 5 July 2025 Digital Object Identifier (**DOI**): https://doi.org/10.35882/jeeemi.v7i3.949

Copyright © 2025 by the authors. This work is an open-access article and licensed under a Creative Commons Attribution-ShareAlike 4.0 International License (CC BY-SA 4.0).

on Large Kernel and Multi-Dimensional Attention Module for Breast Tumor Segmentation," in 2024 IEEE International Conference on Smart Mechatronics (ICSMech), IEEE, Nov. 2024, pp. 89–94. doi:

10.1109/ICSMech62936.2024.10812276.

Author Biography



Muhammad Anang Fathur Rohman is a Computer Science student at Universitas Sebelas Maret (UNS) with a strong passion for Machine Learning and Deep Learning. He is the second outstanding student of Universitas Sebelas Maret (UNS) in 2025 and the

Most Outstanding Student of the Faculty of Technology Information and Data Science, UNS in 2023 and 2024. He received the Best Paper Award at the International Conference on Electronics Representation and Algorithm (ICERA 2025). Actively engaged in academic he has authored several publications research, focusing on cutting-edge deep learning techniques. With a keen interest in advancing technology through research, he continues to explore innovative solutions in artificial intelligence and its real-world applications. contacted He can be at email: anangmuhammad245@student.uns.ac.id.



Heri Prasetyo received the doctoral degree from the Department of Electrical Engineering, National Taiwan University of Science and Technology (NTUST), Taiwan, in 2015. He was awarded the Best Dissertation Award from the Taiwan Association for

Consumer Electronics (TACE) in 2015, and has received multiple Best Paper Awards including from the International Symposium on Electronics and Smart Devices 2017 (ISESD 2017), ISESD 2019, the International Conference on Science in Information Technology (ICSITech 2019), the International Conference on Smart Technology, Applied Informatics, and Engineering (APICS 2022), the International Conference on Informatics and Computing (ICIC 2023), International Conference on Computer, Control, Informatics and its Applications (IC3INA 2024), International Conference Electronics on Representation and Algorithm (ICERA 2025), and the Outstanding Faculty Award from Universitas Sebelas Maret (UNS) in 2019 and 2023. His research interests include multimedia signal processing, computational intelligence, pattern recognition, and machine learning. contacted He can be at email: heri.prasetyo@staff.uns.ac.id.



Ery Permana Yudha is a Lecturer in the Department of Informatics, Universitas Sebelas Maret (UNS), Indonesia. His research interests include Image Processing, Machine Learning, and Information Retrieval. He completed both his Bachelor's and Master's degrees in

Computer Science at Institut Teknologi Sepuluh Nopember (ITS), Surabaya, with a master's thesis on hypergraph-based feature matching. He has been active in research since 2021, including contributing to FKIP UNS's Wahana Bikons project in educational technology. He also mentors undergraduate students and frequently speaks at academic events. He has published more than 10 papers in national and international journals and conferences. He can be contacted at email: erypermana@staff.uns.ac.id.



Chih-Hsien Hsia is a Distinguished Professor at National Ilan University (NIU), where he also serves as the CEO of the AI Promotion Office and Director of the AloT Research Center. He holds two Ph.D. degrees, from Tamkang University (2010) and National Cheng Kung University

(2023), and has a distinguished academic career that includes roles as a Department Chair at NIU and a Visiting Scholar at Iowa State University. An active member of IEEE and IET, his work has been recognized with numerous accolades, including being named a Fellow of the International Association of Advanced Materials (IAAM) in 2025 and receiving multiple Outstanding Young Scholar awards. Dr. Hsia is deeply involved in the academic community, holding key leadership positions in professional societies, serving as General and Technical Program Chair for numerous international conferences, and acting as an Associate Editor for several scientific journals. His research interests include Computer Vision, Embedded Systems, and Information Technology in Education. He can be contacted at email: hsiach@niu.edu.tw.

Manuscript received 10 May 2025; Revised 21 June 2025; Accepted 30 June 2025; Available online 5 July 2025 Digital Object Identifier (**DOI**): https://doi.org/10.35882/jeeemi.v7i3.949