

# A-DenseUNet: Adaptive Densely Connected UNet for Polyp Segmentation in Colonoscopy Images with Atrous Convolution

Sirojbek Safarov

Gachon University

Taeg Keun Whangbo (✉ [taegkeunw@gmail.com](mailto:taegkeunw@gmail.com))

Gachon University

---

## Research

**Keywords:** Image segmentation, Convolutional neural networks, Colonoscopy, Polyp segmentation, Deep learning, Attention, Dilated convolution

**Posted Date:** February 5th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-158417/v1>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at Sensors on February 19th, 2021. See the published version at <https://doi.org/10.3390/s21041441>.

# A-DenseUNet: Adaptive Densely Connected UNet for Polyp Segmentation in Colonoscopy Images with Atrous Convolution

Sirojbek Safarov<sup>1</sup>, Taeg Keun Whangbo<sup>2\*</sup>

<sup>1</sup>*Department of IT Convergence Engineering, Gachon University, South Korea;*

<sup>2</sup>*Department of Computer Science, Gachon University, South Korea;*

*email: taegkeunw@gmail.com*

## Abstract

Colon carcinoma is one of the leading causes of cancer-related death in both men and women. Automatic colorectal polyp segmentation and detection in colonoscopy videos help endoscopists to identify colorectal disease more easily, making it a promising method to prevent colon cancer. In this study, we developed a fully automated pixel-wise polyp segmentation model named A-DenseUNet. The proposed architecture adapts different datasets, adjusting for the unknown depth of the network by sharing multiscale encoding information to the different levels of the decoder side. We also used multiple dilated convolutions with various atrous rates to observe a large field of view without increasing the computational cost and prevent loss of spatial information, which would cause dimensionality reduction. We utilized an attention mechanism to remove noise and inappropriate information, leading to the comprehensive re-establishment of contextual features. Our experiments demonstrated that the proposed architecture achieved significant segmentation results on public datasets. A-DenseUNet achieved a 90% Dice coefficient score on the Kvasir-SEG dataset and a 91% Dice coefficient score on the CVC-612 dataset, both of which were higher than the scores of other deep learning models such as UNet++, ResUNet, and U-Net, for segmenting polyps in colonoscopy images.

*Keywords: Image segmentation, Convolutional neural networks, Colonoscopy, Polyp segmentation, Deep learning, Attention, Dilated convolution*

## Introduction

The third most common form of cancer worldwide for both men and women is colorectal cancer, and its prevalence is increasing every year [1]. The primary cause of colorectal cancer is the growth of glandular tissue in the colonic mucosa. Precise and earlier determination of polyps from virtual colonoscopy screenings are of great significance for the avoidance and timely treatment of colon cancer [2]. However, manual detection depends on proficient endoscopists, and it takes a long time. Recent surveys have shown that more than 25% of polyps in patients undergoing colonoscopy are not detected [3]. The late diagnosis of missed polyps can lead to a low survival rate for colon

cancer patients [4]. Computer-aided detection (CAD) systems are used to detect and segment polyps from endoscopic images and video screenings, which allows endoscopists to focus their attention on the polyps displayed on the screen and act as a second viewer. This can decrease the likelihood of overlooked polyps [5].

Designing an accurate CAD system is challenging because of the high cost of labeled medical datasets for training and testing. Polyps have a wide range of colors, sizes, shapes, appearances, or combinations of these features. There are similar inter-classes and various intra-classes for four different polyp classes: adenoma, hyperplastic, serrated, and mixed. In addition, background objects are very similar; for example, the background mucosa can mix with a polyp or stool [43]. Even though these factors make the polyp segmentation task challenging, we surmise that there is still a great prospect to create such systems for medical use.

In recent years, deep-learning-based techniques have achieved significant success in the computer vision domain [44]–[46], and interest in applying deep learning to endoscopic image segmentation has grown. In particular, encoder-decoder-based methods such as U-Net [7], UNet++ [35], SegNet [47], and fully convolutional networks (FCNs) [19] have been commonly used for semantic segmentation. These networks down-sample the image several times to capture the required feature maps and up-sample once or multiple times to enable effective localization [7][19]. Furthermore, skip connection strategies have been successful in saving fine-grained information and improving the efficiency of the network, even on complicated datasets.

Recent research has shown that the attention mechanism has been commonly used to preserve the dependency of features in certain computer vision tasks such as object detection [54], image classification [52][53], and image segmentation [48]–[51]. The attention method enables the model to attend more closely to essential features without any external supervision, and it can avoid identical feature maps at various scales to lead to better feature representation. The attention mechanism improves network efficiency over traditional methods with or without multiscale features.

In this work, we propose a novel deep learning method called Adaptive Densely Connected UNet with atrous convolution (A-DenseUNet) for medical image segmentation. The core assumption behind our proposed method is that the model can accumulate different levels of semantic information on a network to obtain global multiscale features. In addition, the proposed architecture uses dilated convolution to capture finer details and to eliminate scratching artifact issues. We evaluated our model using two public datasets: Kvasir-SEG [39] and CVC-612 [40]. We compared our results with those of popular deep learning models such as U-Net [7], wide U-Net [35], ResUNet [8], and UNet++ [35]. The results indicate that our method boosts performance and achieves better results than other methods.

In summary, this study makes the following contributions:

- We designed a new robust U-Net-based encoder-decoder network structure that uses dense connections as a powerful encoder model and accomplishes an adaptable image segmentation algorithm to integrate deep and superficial features, which can directly combine multiscale features to boost segmentation performance.
- We utilized an attention mechanism that fuses derived information from various modules and focuses on core information by removing noise and irrelevant regions.

- Our method uses dense blocks, residual blocks, transition blocks, and atrous convolution block capabilities, and it improves the outcome of the colorectal polyp segmentation compared to other state-of-the-art methods. Our model obtained good results with small datasets.
- We evaluated our model on the Kvasir-SEG and CVC-612 datasets, and the experimental results show that it achieved the highest intersection over union (IoU) and Dice coefficient.

The remainder of this paper is organized as follows. In Section II, we review some existing related studies. In Section III, we present our proposed densely connected deep learning architecture. We present the experimental settings and qualitative and quantitative analysis of the semantic segmentation results in Sections IV and V, respectively. In Section VI, we discuss the experimental results. Finally, we conclude this paper in Section VII.

## Related Work

Over the past two decades, the detection and classification of gastrointestinal (GI) tract diseases and the creation of effective, robust methods to automatically detect polyps in colonoscopy images and videos have been active scientific areas. The performance of machine learning-based polyp detection and segmentation software has come close to that of high-level endoscopists.

Some earlier studies used the texture and color details of polyps to create handcrafted descriptors [10][11][12][13]. For instance, Karkanis et al. [10] utilized a supplemented sliding window scheme and color wavelet texture information as descriptors to designate polyps from colonoscopy images and videos. Subsequently, researchers used spatio-temporal, edge, intensity, and shape features to detect polyps automatically. For example, Hwang et al. [14] used elliptical shape information to detect polyps automatically, whereas Wang et al. [15][16] presented edge cross-section profiles. To improve the detection performance, some methods combine two or more features [17][18]. Tajbakhsh et al. [18] integrated local intensity variation patterns and global geometric constraints to detect polyps. Although these methods achieved significant progress, they still suffer from inferior detection accuracy. The primary cause of the low accuracy level is the limited representation ability of handcrafted features to deal with both the low-level inter-class variety between hard mimics and polyps and the high-level intra-class variety of polyps.

Recent deep convolutional neural networks (CNNs) have shown noticeably better results in many biomedical image analysis domains, including object detection [26][27][28][29], classification [23][24][25], and semantic segmentation [7][30][31][32]. Some researchers have attempted to use CNNs to manage the automated polyp detection domain. For instance, Tajbakhsh et al. [33] suggested a CNN architecture for polyp detection that takes low-level handcrafted information as input and utilizes a group of CNNs to learn the shape, color, and temporal features of polyps. However, this model learned temporal and spatial information using various networks that may limit the discrimination capability. Thus, most features from colonoscopy videos have not been fully explored.

The new generation of CNNs uses transposed convolution layers to generate a probability map in image segmentation tasks. Long et al. [19] proposed the fully convolutional network (FCN)

method, which achieved state-of-the-art semantic segmentation results. FCN obtains the segmentation results without post-processing steps by using pixel-to-pixel and end-to-end training. Ronneberger et al. added modifications and extensions to the FCN to develop the U-Net [7] architecture. U-Net integrates high-resolution spatial feature maps with high-level contextual information for medical image segmentation. Inspired by these approaches, several researchers have proposed models to solve segmentation issues in a wide variety of areas [8][9][20].

The majority of studies published in the sphere of polyp segmentation achieved significant results only on special datasets, and test cases often utilized small validation and training datasets [21][22]. Furthermore, some of the scientific work focuses only on a particular type of polyps, and some of them employ non-public datasets, which makes it difficult to compare and reproduce the results. Consequently, the ML models cannot yet achieve similar or better results than endoscopists. There is an opportunity to enhance the efficiency of CAD systems, making major improvements and producing more effective and reliable architectures for polyp segmentation.

## Proposed Method

The A-DenseUNet architecture is based on UNet++ [9] and densely connected convolutional networks (DenseNets) [34], utilizing the strength of U-Net [7] and DenseNet [34]. The proposed A-DenseUNet architecture takes advantage of dense blocks, atrous convolution, residual blocks, attention blocks, and restrictive skip connections.

### Overview

The proposed segmentation architecture utilizes the U-Net [7] concept which includes an encoder block on the left and a decoder block on the right. Figure 1 depicts the entire structure of the proposed method. It takes a training dataset  $X$  that consists of  $N$  sample images  $x$ :  $X = x_1, x_2, \dots, x_N$ , with corresponding  $Y = y_1, y_2, \dots, y_N$ . Then, each ground truth pixel  $i$  of any given sample  $y$  is  $y \in [0, 1]$ . We feed our network with a  $224 \times 224 \times 3$  image and obtain a  $224 \times 224 \times 1$  output segmentation mask. In the encoding path, an input image passes through a dense block that includes a combination of atrous convolutional, rectified linear unit (ReLU), and batch normalization layers. The dense block is followed by a transition block that contains a pooling layer that reduces the size of the feature map after each successive dense block. In the decoding path, transposed convolution is used to increase the feature map size back to the original size. After a very deep encoder path, there may be a loss of essential details. To handle such a problem, UNet++ [9] introduces restrictive skip connections that combine the encoding part with the output of up-sampling through channel concatenation.

We applied skip connections to unify various depths of U-Nets into one structure and used an attention mechanism to filter irrelevant information from the features. The depth of the network,  $s$ , is 5, which means we used a down-sampling approach five times and halved its feature map size each time. After five down-samples, we obtained the final  $7 \times 7$  spatial feature map. We used the attention mechanism to create a relationship between the various model information at different

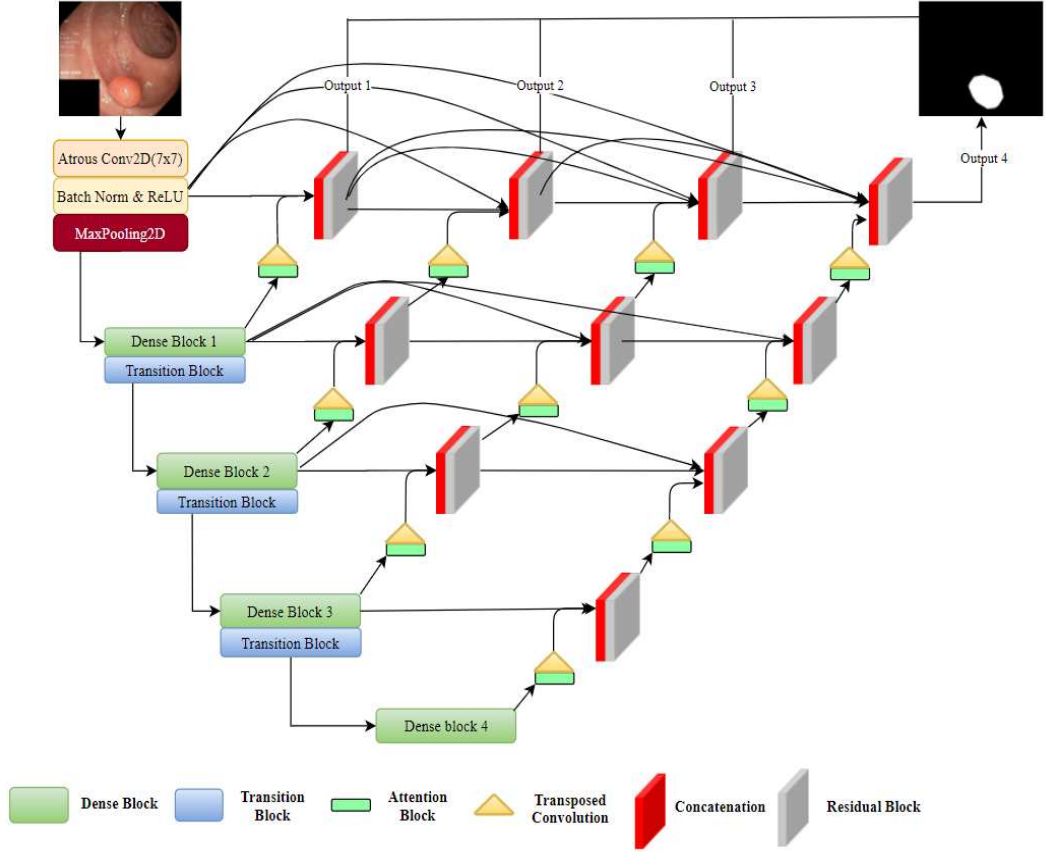


Figure 1. Block diagram of the proposed A-DenseUNet architecture: DenseNet is used as an encoder, Transposed convolution is performed for up-sampling between levels.

depths. The attention blocks reduce the noise and unnecessary features, and only important information can pass to the next layer. The output of the attention block is up-sampled by transposed convolution and concatenated with the same depth output as the encoder part. After concatenation, the feature map passes the residual block (dilated convolution followed by batch normalization and ReLU activation), which allows features to converge more quickly. Other decoder blocks at levels  $s=2$  to  $s=5$  use such blocks. Finally, feature maps from all U-Net depths are agglomerated and then averaged, after which layers for  $1 \times 1$  convolution and sigmoid activation, as shown in Equation (1), are used to obtain the final segmentation map. We trained our network with the binary cross-entropy loss function based on the ground truth for the training images. Equation (2) shows the formula of the loss function, where  $y$  is the label and  $p(y)$  is the predicted probability for  $N$  pixels.

$$y = \frac{1}{1 + e^{-x}} \quad (1)$$

$$L = - \sum_{n=1}^N y_i \log(p(y_i)) + (1 - y_i) \log(1 - p(y_i)) \quad (2)$$

## Dense Units

Training deeper neural networks can increase a model's accuracy, but it can also cause degradation problems and interrupt the training process [8][34]. To solve this type of problem, Huang et

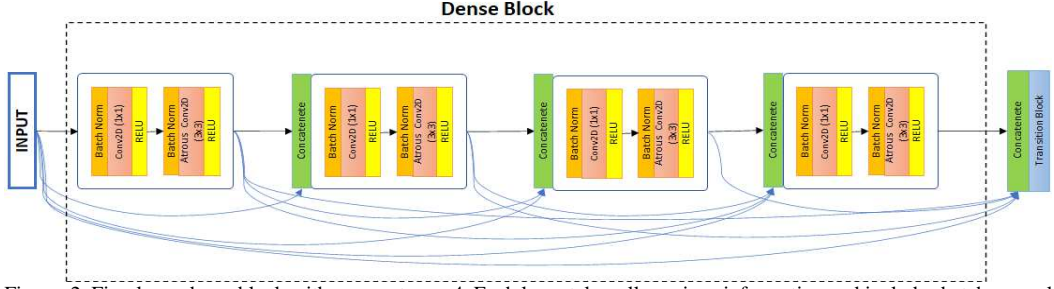


Figure. 2. Five-layer dense block with grow rate  $n = 4$ . Each layer takes all previous information and includes batch normalization, atrous convolution, and ReLU activation.

al. [34] proposed densely connected convolutional networks (DenseNets), which allow all subsequent layers to connect directly, as shown in Figure 2.

Accordingly, the  $l^{th}$  layer takes the feature maps of all previous layers,  $x_0, \dots, x_{l-1}$ , as input:

$$x_l = H_l([x_0, x_1, \dots, x_{l-1}]) \quad (3)$$

where  $[x_0, x_1, \dots, x_{l-1}]$  refers to the concatenation of the feature maps produced in layers  $0, \dots, l-1$ . Dense encoder blocks have several advantages; for example, densely connected layers have fewer output dimensions than other networks, which can help to avoid learning excessive features and reduces the time required.

Furthermore, the densely connected layers provide maximum gradient flow, and very deep neural networks alleviate the vanishing gradient problem. Based on these advantages, we used DenseNets [34] as the encoder part of our proposed method. Table 1 presents the encoder layers of the proposed architecture. The input layer takes  $224 \times 224 \times 3$ -sized images; thus, all training data were resized to fit the given size.

In the first layer of the network, using  $7 \times 7$  dilated convolution with 96 filters and a stride of two, we obtained an output feature map of  $112 \times 112 \times 96$  after the first convolution layer. Then, we used four dense blocks and three transition blocks to create the remaining encoder layers,  $s=2$  to  $s=5$ . Each dense block includes batch normalization, dilated convolution, and ReLU non-linearity and is repeated several times, as shown in Table 1, to create a deeper encoder path and obtain more robust feature maps. After each dense block (except dense block 4) is a transition block, which consists of  $1 \times 1$  convolution and  $2 \times 2$  average pooling with a stride of two. A  $1 \times 1$  convolution is used Table 1. Densely connected encoder block of the proposed A-DenseUNet architecture. Note that “ $1 \times 1, 192$  conv” corresponds to  $1 \times 1$  kernel size convolution with 192 features and a sequence of BN-Conv-ReLU layers. “[ $\ ] \times n$ ” indicates  $n$  iterations of the dense block.

	Feature size	Encoder DenseNet-164 (k=48)
input	$224 \times 224 \times 3$	-
convolution 1	$112 \times 112$	$7 \times 7, 96, \text{ stride } 2$
pooling	$56 \times 56$	$3 \times 3$ max pool, stride 2
dense block 1	$56 \times 56$	$\begin{bmatrix} 1 \times 1, 192 & \text{conv} \\ 3 \times 3, 48 & \text{conv} \end{bmatrix} \times 6$
transition layer 1	$56 \times 56$	$1 \times 1$ conv
	$28 \times 28$	$2 \times 2$ average pool, stride 2
dense block 2	$28 \times 28$	$\begin{bmatrix} 1 \times 1, 192 & \text{conv} \\ 3 \times 3, 48 & \text{conv} \end{bmatrix} \times 12$
transition layer 2	$28 \times 28$	$1 \times 1$ conv
	$14 \times 14$	$2 \times 2$ average pool, stride 2
dense block 3	$14 \times 14$	$\begin{bmatrix} 1 \times 1, 192 & \text{conv} \\ 3 \times 3, 48 & \text{conv} \end{bmatrix} \times 36$
transition layer 3	$14 \times 14$	$1 \times 1$ conv
	$7 \times 7$	$2 \times 2$ average pool, stride 2
dense block 4	$7 \times 7$	$\begin{bmatrix} 1 \times 1, 192 & \text{conv} \\ 3 \times 3, 48 & \text{conv} \end{bmatrix} \times 24$

before the pooling layer to reduce the channels of the feature map. After dense block 4, a robust  $7 \times 7$  feature map is obtained, which is decoded to produce the final output.

## Adaptive Network Structure

The original U-Net obtains the final segmentation result from a fixed number of down-samples and a corresponding number of up-samples. In practice, some datasets contain images of various sizes, and there is a significant difference between the amount of information contained in various-sized images. One down-sampling and one up-sampling might be sufficient to obtain satisfactory segmentation results from small, simple data. Large, complicated datasets require multiple up- and down-samplings to obtain semantic feature maps of various regions because it is difficult to obtain global information from a small-scale network. Figure 3 represents the U-Net network with depths

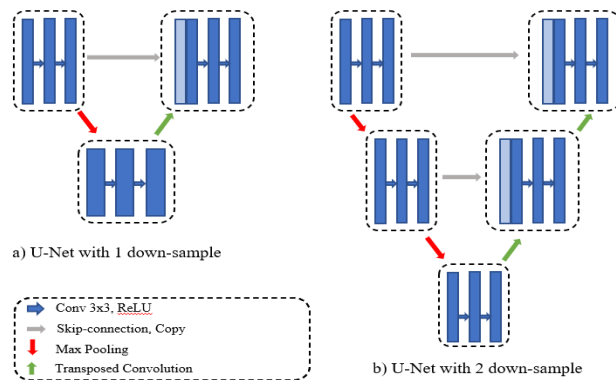


Figure. 3. Multi-depth U-Net models.

of one and two. To overcome different depth problems, Zhou et al. [35] proposed redesigning the skip connections to integrate the advantages of different depths of U-Net into one architecture.

We redesigned the skip connections in our proposed method to connect different depths in the U-Net structure. In addition, we added an efficient feature map transition and aggregated different layer characteristics. As shown in Fig. 1, horizontal dense connections and connections between each depth are added. Horizontal densely connected layers are equipped with dilated convolution, batch normalization, and residual blocks. Dense connections pass feature maps efficiently at various depths. Even though the use of various depth decoder architectures and densely connected structures increases the network size, it enhances the efficiency of the method. The final segmentation mask is achieved by averaging the output of each U-Net depth and employing a  $1 \times 1$  convolution and sigmoid classifier.

## Attention Units

Over the last few years, the attention mechanism has become very popular in various deep learning research areas, starting with natural language processing (NLP) [36]. Recently, it has been applied to computer vision tasks. The attention model has been utilized as a pixel-wise prediction model in the semantic segmentation domain [37]. It identifies the sections of the network that need more attention. Continuous use of an attention mechanism at each level allows long-range spatial dependency of feature maps. The attention block also decreases each image's computing cost to a fixed dimensional vector. Thus, the fundamental value of an attention unit is that it is straightforward



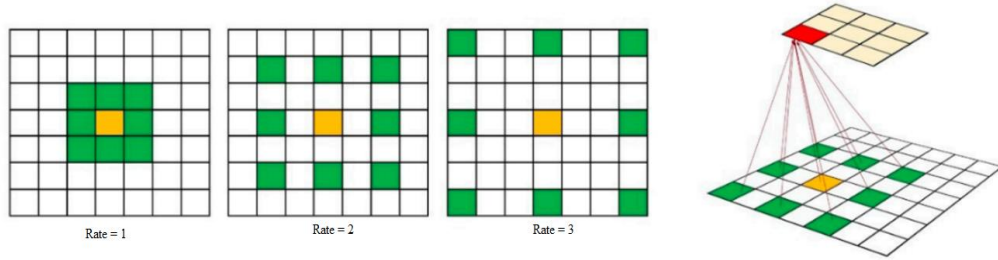


Figure 4. Dilated convolutions with different dilation rates. A dilation rate of one is normal convolution.

and can be applied to every input scale to strengthen the consistency of the features that emphasize the result.

We implemented an attention block in the proposed method for medical image segmentation. We placed the attention block before the up-sampling layer at each level of the U-Net decoder path. Specifically, the model encodes various semantic feature maps at various stages. The attention mechanism is used to enhance the flow of the spatial feature map to the next level of the decoding side; to generate relevant feature maps, up-sampling information is fused with the corresponding encoder-side information. Thus, attention blocks at various stages allow the proposed network to encode low-level to high-level information at different scales and provide only relevant regions to the next layer.

## Dilated Convolution

The concept of dilated convolution comes from wavelet decomposition [38]. It is also called “atrous convolution” and “algorithm à trous.” Dilated convolution enables the model to arbitrarily expand the filter field of view at every DCNN (Deep Convolutional Neural Network) layer. In order to hold both the calculation and the number of parameters contained, CNNs usually use small convolutional kernels (typically 3×3). Dilated convolution with a rate of  $r$  adds  $r-1$  zeroes between the consequent filter values, as shown in Figure 4. It thus provides an efficient field of view control mechanism and finds the optimal trade-off between detailed localization (small field of view) and assimilation of context (large field of view).

The results section demonstrates that one of the keys to the success of our model is the use of dilated kernels, as they allow the network to increase the receptive field without adding computational complexity or increasing the network information capability. We replace normal convolutions from Equation (4) by atrous convolution in Equation (5) with a dilation rate of 2 for every layer of the network. In Equation (5),  $s + lt = p$  indicates that some points have been skipped during the convolution.

$$(F * k)(p) = \sum_{s+t=p} F(s)k(t) \quad (4)$$

$$(F * k_l)(p) = \sum_{s+lt=p} F(s)k(t) \quad (5)$$

## Experiments

We evaluated our proposed A-DenseUNet architecture on two public segmentation datasets, Kvasir-SEG [39] and CVC-612 [40].

### Datasets

*Kvasir-SEG*: We utilized the Kvasir-SEG dataset [39], which has 1,000 polyp images and their

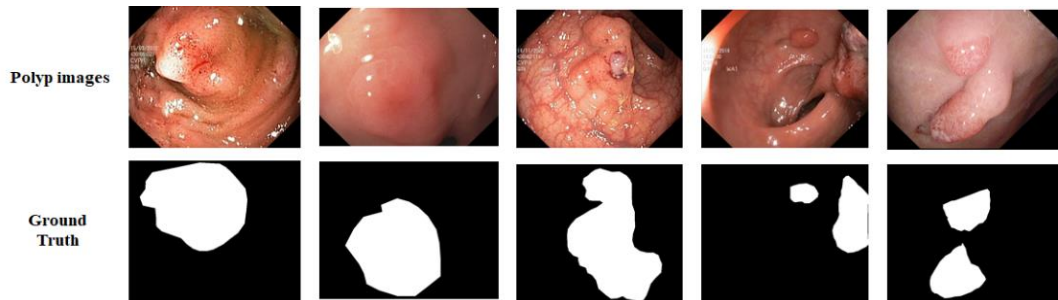


Figure 5. Example of data from Kvasir-SEG dataset. The first row shows original images and the second row presents their respective ground truth.

corresponding ground truth masks annotated by professional gastroenterologists from Vestre Viken Health Trust in Norway, as shown in Figure 5. The images have sizes ranging from  $332 \times 487$  to  $1920 \times 1072$  pixels, but training and testing were performed with an image resolution of  $224 \times 224$  pixels. The images were randomly split into 80% for training, 10% for validation, and 10% for testing.

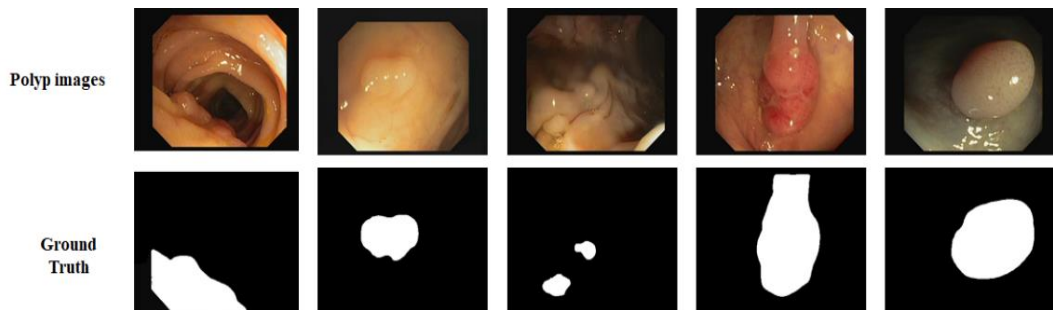


Figure 6. Images and ground truth masks from the CVC-612 dataset.

*CVC-612*: In addition, we used the CVC-612 [40] dataset, which has 612 images with a size of  $384 \times 288$  pixels from 31 colonoscopy series. The images were split into training, validation, and testing sets in the ratio of 80:10:10. All training, validation, and testing were performed with an image size of  $224 \times 224$  pixels. Figure 6 shows some example images and corresponding masks from the CVC-612 dataset.

### Data Augmentation

The effectiveness of deep learning networks depends significantly on the size of the training dataset. It is clear that in the case of polyp segmentation, the training dataset is limited, at least with respect to typical training images employed in the context of deep learning. Furthermore, certain polyp forms are not represented in the dataset, and for other types only a few examples are available.

Hence, it is important to extend the training dataset by data augmentation. Data augmentation is conducted to provide additional polyp images for training deep neural networks. Even though this approach cannot produce new polyp forms, it can provide extra data samples based on various image acquisition conditions, such as colon deformations, camera position, and illumination.

All training samples were resized to 224×224 pixels in a manner such that the image aspect ratio was retained. This process included random cropping augmentation. All images were augmented using four augmentation techniques: (1) rotation, with the angle of rotation randomly chosen from the range 0° to 90°; (2) reflection, horizontally and vertically; (3) elastic deformation with a fixed 10×10 grid; and (4) color adjustment by random gamma augmentation. After augmentation, the Kvasir-SEG training dataset contained a total of 8,000 images.

## Evaluation Metrics

To evaluate the polyp segmentation, we used the following well-known segmentation evaluation metrics: recall, precision, intersection over union (IoU), and Dice coefficient. We calculated these metrics using well-known parameters such as true positive (TP), false positive (FP), and false negative (FN).

$$\text{Recall} = \frac{TP}{TP + FN} \quad (6)$$

$$\text{Precision} = \frac{TP}{TP + FP} \quad (7)$$

Intersection over Union: The IoU is a standard metric for evaluating segmentation models. The equation presents the similarity between the predicted pixels ( $Y'$ ) and the true mask ( $Y$ ).

$$\text{IoU} = \frac{|Y' \cap Y|}{|Y' + Y|} = \frac{TP}{TP + FP + FN} \quad (8)$$

Dice similarity coefficient: The Dice similarity coefficient is a standard metric for comparing the pixel-wise results between the ground truth and predicted segmentation. The formula of the Dice coefficient is defined as

$$\text{Dice coefficient} = \frac{2 * |Y' \cap Y|}{|Y'| + |Y|} = \frac{2 * TP}{TP + FP + FN} \quad (9)$$

where  $Y'$  is the predicted set of pixels, and  $Y$  signifies the ground truth of the item.

## Implementation Details

We trained all the methods in the Keras framework [41] with TensorFlow [42] as a backend. We trained our model with 224×224-pixel images. We set the batch size to 10, and we trained the model for 100 epochs. We used the Adam optimizer with reduced learning rate callback; the learning rate starts from 0.01 and is divided by 10 when the patience level exceeds 5. We used an early-stop mechanism on the validation set to avoid overfitting. We chose ReLU as the non-linear activation and binary cross-entropy as the loss function. To convert the predicted pixels to the background or foreground, we used a threshold value of 0.5. All the models were implemented using two NVIDIA

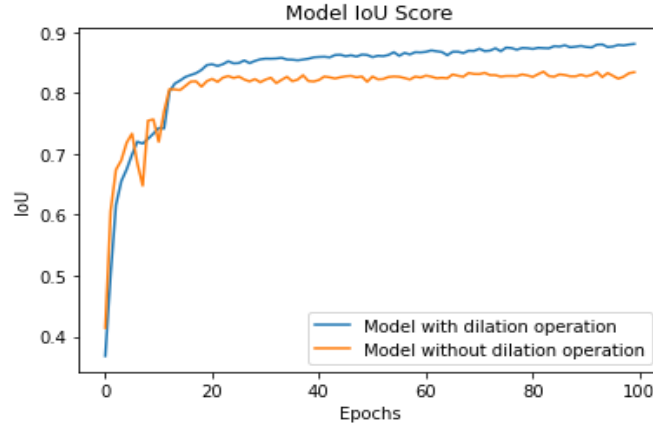


Figure 7. IoU score with and without dilated convolution.

GTX 1080 GPUs, each with 8 GB of memory. It took five hours to complete the training of our proposed model.

## Results

We performed comprehensive experiments to assess the effectiveness of our proposed A-DenseUNet architecture. Four state-of-the-art deep learning models, U-Net [7], wide U-Net [35], ResUNet [8], and UNet++ [35], were selected for comparative analysis of the proposed method.

To test the effectiveness of the dilated convolution, we trained our model with standard and dilated convolution with a 3×3 filter size and dilation rate of two. Figure 7 presents the training IoU score with and without dilated convolution, showing that the model with dilated convolution performed better than the model with standard convolution. In addition, to show the efficacy of the attention mechanism, we trained our proposed model with and without attention blocks. A comparison between the models trained with and without attention blocks shows that the model with an attention mechanism demonstrated strong attention ability by emphasizing the discriminative region

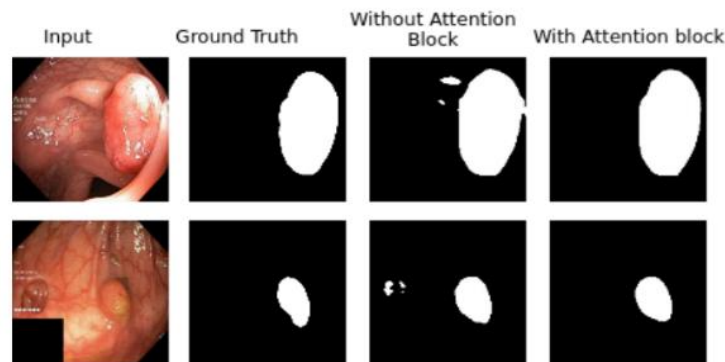


Figure 8. Effect of attention block in the network. By adding this, we were able to suppress the irrelevant regions.

of interest rather than concentrating on specular reflection and the normal area. Figure 8 depicts the qualitative difference between them.

*Kvasir-SEG dataset results:* We improved our proposed A-DenseUNet architecture with various sets of hyperparameters. During the model training, we manually tuned the hyperparameters

Table 2. Kvasir-SEG evaluation results of all methods.

Method	Params	Dice	IoU	Recall	Precision
<b>A-DenseUNet</b>	11.0M	<b>0.9085±0.77</b>	<b>0.8615±0.43</b>	<b>0.9448±0.31</b>	<b>0.9766±0.91</b>
UNet++ [35]	9.0M	0.8021±1.9	0.7215±0.34	0.7914±1.36	0.9321±0.61
ResUNet [8]	8.5M	0.7864±1.35	0.5421±0.58	0.7861±0.94	0.8912±1.45
Wide U-Net [35]	9.1M	0.7645±0.62	0.7112±0.91	0.7684±0.55	0.9231±0.57
U-Net [7]	7.8M	0.7062±0.46	0.5628±0.13	0.7768±0.61	0.9022±0.63

with various hyperparameter sets and evaluated the results. Table 2 presents the results of A-DenseUNet, ResUNet [8], UNet++ [35], wide U-Net, and the original U-Net [7] on the Kvasir-SEG [39] dataset. The data indicate that the proposed architecture outperformed all current methods.

*CVC-612 dataset results:* After achieving good results on the Kvasir-SEG dataset, we tested our method on the CVC-612 dataset. Table 3 presents the performance of all the models on the CVC-

Table 3. CVC-612 evaluation results of all methods.

Method	Params	Dice	IoU	Recall	Precision
<b>A-DenseUNet</b>	11.0M	<b>0.8912±0.57</b>	<b>0.8553±0.35</b>	<b>0.9448±0.35</b>	0.9266±0.75
UNet++	9.0M	0.7815±0.15	0.7241±1.67	0.8064±0.73	0.9076±0.20
ResUNet	8.5M	0.7397±0.59	0.5597±0.93	0.7643±0.36	0.8627±1.02
Wide U-Net	9.1M	0.7754±1.31	0.7078±1.31	0.7831±0.43	0.9113±0.59
U-Net	7.8M	0.6943±0.94	0.5798±0.97	0.7648±0.57	<b>0.9418±1.24</b>

612 dataset. The proposed method achieved the largest Dice coefficient, IoU, and recall. U-Net obtained the highest precision score, but its other important metric scores for segmentation were not competitive.

Figures 9 and 10 present the qualitative results for all deep learning methods. Tables 2 and 3 and the qualitative results show the dominance of A-DenseUNet over the baseline methods such as UNet++ [35], ResUNet [8], wide U-Net, and the original U-Net [7]. On the Kvasir-SEG dataset, the proposed architecture achieved mean improvements of 10.64%, 12.21%, 14.4%, and 20.23%, as measured by the Dice coefficient, and 14.12%, 32.21%, 15.03%, and 29.86%, as measured by the IoU score. The large margin of difference between the proposed architecture and the existing methods could be interpreted as indicating that the combination of dilated convolution, attention mechanism, and multiscale features plays a crucial role in optimizing segmentation efficiency. The proposed model encodes multiscale semantic information at every stage, which allows the conservation of more fine-grained feature maps at the decoder block, unlike U-Net and ResUNet, which use the same-scale feature map concatenation. Furthermore, the attention mechanism enhances the network by focusing on important information that boosts the segmentation efficiency.

## Discussion

The proposed A-DenseUNet model achieved adequate results on both the CVC-612 and Kvasir-SEG datasets. From the qualitative results, it is obvious that the proposed model’s segmentation mask performed better than other methods to capture the shape of information on the Kvasir-SEG dataset. The results show that the predicted segmentation mask in A-DenseUNet is closer to the

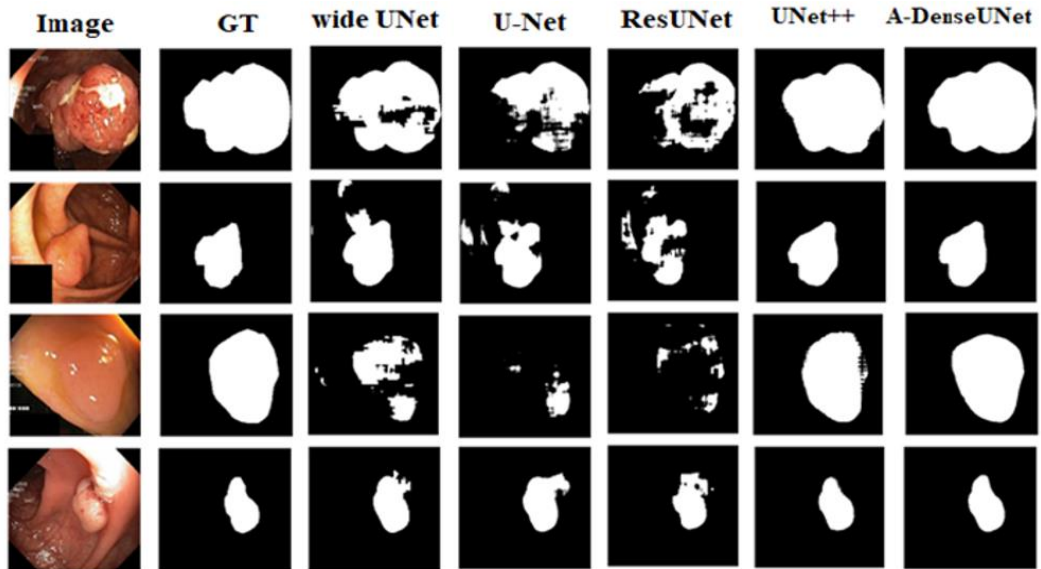


Fig. 9. Qualitative segmentation results of various models on the Kvasir-SEG dataset. Experimental results show that A-DenseUNet produces better segmentation masks than other state-of-the-art networks.

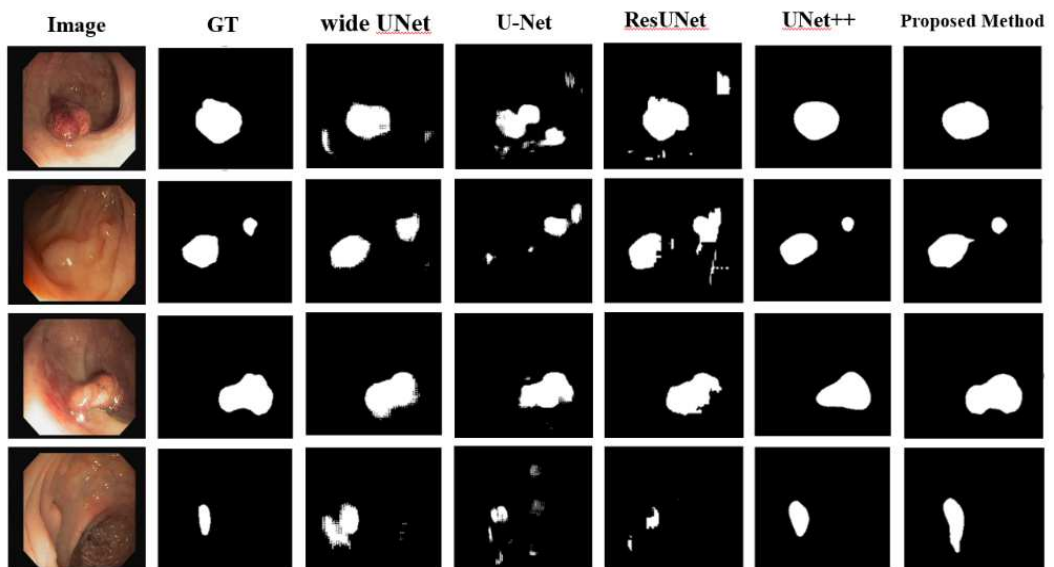


Fig. 10. Qualitative segmentation results of various models on the CVC-612 dataset.

ground truth mask than that in other state-of-the-art architectures. However, segmentation masks from UNet++ and wide U-Net are also competitive.

During the training process, we used various loss functions to improve our results, such as Jaccard loss, Dice loss, mean square loss, and binary cross-entropy loss. According to our experiments, the method achieved a better Dice coefficient value with all loss functions, whereas the IoU score was higher with a binary cross-entropy loss function. We chose the binary cross-entropy loss function based on our analytical assessment. In addition, we found that the number of kernels, batch size, optimizer, loss function, and depth of the model may affect the result.

We speculate that the efficiency of the model could be further improved by enlarging the dataset size, using various augmentations, and adding certain post-processing steps. We designed a very deep neural network architecture to achieve significant performance, although it can increase the

number of parameters. Our proposed A-DenseUNet model not only achieved better results on biomedical image segmentation but also achieved good results in pixel-wise image classification and natural image segmentation tasks. We used all our experience and knowledge to optimize the model, but there might be further optimizations, which could affect the performance of the method.

## Conclusion

In this paper, we have presented an end-to-end biomedical image segmentation architecture, A-DenseUNet, to achieve more accurate segmentation results. The proposed architecture takes advantage of dense blocks, atrous convolution, residual blocks, attention blocks, and restrictive skip connections. Experiments on the CVC-612 and Kvasir-SEG datasets demonstrate that the proposed method outperforms the state-of-the-art UNet++, ResUNet, and U-Net architectures in predicting accurate segmentation masks. Our model achieved the best Dice coefficient and IoU score among the models. Future studies need to evaluate our model on various medical and natural image segmentation datasets.

## Abbreviations.

ReLU – Rectified Linear Unit;  
IoU – Intersection over Union;  
FCN – Fully Convolutional Network;  
CAD – Computer-Aided Detection;  
CNN – Convolutional Neural Network;  
ML – Machine Learning;

## Availability of data and materials

Data sets are openly available.

## Competing interests

Not applicable.

## Funding

Not applicable.

## Authors' Contribution

S.Sirojbek designed the model and the computational framework and analyzed the data. Both S.Sirojbek and T.K.Whangbo authors contributed to the final version of the manuscript. T.K.Whangbo supervised the project.

## Acknowledgments

Not applicable.

## Authors' information

Sirojbek Safarov received the B.S. degree from Tashkent University of Information Technologies, Tashkent, Uzbekistan in 2018. Currently, he is M.S. at Gachon University, in the faculty of IT Convergence Engineering, Korea. His research interests include Computer Vision, Image Processing, Machine/Deep Learning, and Artificial Intelligence. He is working on machine learning and deep learning area.

Taeg Keun Whangbo received the M.S. degree from the City University of New York in 1988 and the Ph.D. degree both in Computer Science from Stevens Institute of Technology in 1995. Currently, he is a professor in the Department of Computer Science and vice president of Gachon University, Korea. Before he joined Gachon University, he was the software developer in Q-Systems which is located in New Jersey from 1988 to 1993. He was also a researcher in Samsung Electronics from 2005 to 2007. From 2006 to 2008, he was the president of the Association of Korea Cultural Technology. His research areas include Computer Graphics, HCI, and VR/AR.

## References

- [1] J. Silva, A. Histace, O. Romain, X. Dray, and B. Granado. "Toward embedded detection of polyps in wce images for early diagnostics of colorectal cancer." *International Journal of Computer Assisted Radiology and Surgery* 9(2), 2014; 283-293.
- [2] M. Arnold, M.S. Sierra, M. Laversanne, I. Soerjomataram, A. Jamel, and F. Bray: "Global patterns and trends in colorectal cancer incidence and mortality", *Gut*, 2016; pages gutjnl 2015.
- [3] A. M. Leufkens, M. G. H. van Oijen, F. P. Vleggaar, and P.D. Siersema. "Factors influencing the miss rate of polyps in a back-to-back colonoscopy study", *Endoscopy*, 2012; 44(05), 470-475.
- [4] L. Rabeneck, H. B. El-Serag, J. A. Davila, and R. S. Sandler, "Outcomes of colorectal cancer in the united states: No change in survival (1986-1997)", *The American journal of gastroenterology*, 2003; vol.98, no. 2, pp. 471-477.
- [5] Y. Mori and S.-e. Kudo, "Detecting colorectal polyps via machine learning", *Nature biomedical engineering*, 2018, vol. 2, no. 10, p.713.
- [6] L. Yu, H. Chen, Q. Dou, J. Qin, and P. A. Heng, "Integrating online and offline three-dimensional deep learning for automated polyp detection in colonoscopy videos", *IEEE JBHI* 21(1), 2016; 65-75.
- [7] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation", in Proceedings of International Conference on Medical image computing and computer-assisted intervention. *Springer*, 2015; pp 234-241.
- [8] Z. Zhang, Q. Liu, and Y. Wang, "Road extraction by deep residual unet," *IEEE Geoscience and Remote Sensing Letters*, 2018;vol. 15, no. 5, pp. 749-753.
- [9] Z.Zhou, M.M.R. Siddiquee, N. Tajbakhsh, and J. Liang.: "Unet++: A nested u-net architecture for medical image segmentation". *IEEE TMI*, 2019; pp. 3-11.
- [10] S. Karkanis, D. K. Iakovidis, D. E. Maroulis, D. Karras, M. Tzivras *et al.*, "Computer-aided tumor detection in endoscopic video using color wavelet features," *Information Technology in Biomedicine, IEEE Transactions on*, 2003; vol. 7, no. 3, pp. 141-152.
- [11] D. K. Iakovidis, D. E. Maroulis, S. A. Karkanis, and A. Brokos, "A comparative study of texture features for the discrimination of gastric polyps in endoscopic video," in *18th IEEE Symposium on ComputerBased Medical Systems (CBMS'05)*. IEEE, 2005; pp. 575-580.
- [12] L. A. Alexandre, N. Nobre, and J. Casteleiro, "Color and position versus texture features for endoscopic polyp detection," in *2008 International Conference on BioMedical Engineering and Informatics*, 2008; vol. 2. *IEEE*, pp. 38-42.
- [13] S. Ameling, S. Wirth, D. Paulus, G. Lacey, and F. Vilarino, "Texturebased polyp detection in colonoscopy," in *Bildverarbeitung fur die Medizin 2009*. Springer, 2009; pp. 346-350



- [14] S. Hwang, J. Oh, W. Tavanapong, J. Wong, and P. C. De Groen, "Polyp detection in colonoscopy video using elliptical shape feature," in *Image Processing, 2007. ICIP 2007. IEEE International Conference on*, 2017; vol. 2. IEEE, pp. II-465.
- [15] Y. Wang, W. Tavanapong, J. Wong, J. Oh, and P. C. de Groen, "Partbased multiderivative edge cross-sectional profiles for polyp detection in colonoscopy," *Biomedical and Health Informatics, IEEE Journal of*, 2014; vol. 18, no. 4, pp. 1379-1389.
- [16] Y. Wang, W. Tavanapong, J. Wong, J. H. Oh, and P. C. De Groen, "Polyp-alert: Near real-time feedback during colonoscopy," *Computer methods and programs in biomedicine*, 2015; vol. 120, no. 3, pp. 164-179.
- [17] J. Silva, A. Histace, O. Romain, X. Dray, and B. Granado, "Toward embedded detection of polyps in wce images for early diagnosis of colorectal cancer," *International Journal of Computer Assisted Radiology and Surgery*, 2014; vol. 9, no. 2, pp. 283-293.
- [18] N. Tajbakhsh, S. R. Gurudu, and J. Liang, "Automated polyp detection in colonoscopy videos using shape and context information," *IEEE Transactions on Medical Imaging*, 2016; vol. 35, no. 2, pp. 630-644.
- [19] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of IEEE conference on computer vision and pattern recognition (CVPR)*, 2015; pp. 3431-3440.
- [20] F. Milletari, N. Navab, and S.-A. Ahmadi, "V-net: Fully convolutional neural networks for volumetric medical image segmentation," in *Proceeding of International Conference on 3D Vision (3DV)*, 2016; IEEE, pp. 565-571.
- [21] P. Brandao, O. Zisimopoulos, E. Mazomenos, G. Ciuti, J. Bernal, M. Visentini-Scarzanella, A. Mencias, P. Dario, A. Koulaouzidis, A. Arezzo *et al.*, "Towards a computed-aided diagnosis system in colonoscopy: automatic polyp segmentation using convolution neural networks," *Journal of Medical Robotics Research*, 2018; vol. 3, no. 2, p. 1840002.
- [22] Y. Wang, W. Tavanapong, J. Wong, J. Oh, and P. C. De Groen, "Partbased multiderivative edge cross-sectional profiles for polyp detection in colonoscopy," *IEEE Journal of Biomedical and Health Informatics*, 2013; vol. 18, no. 4, pp. 1379-1389.
- [23] A. A. A. Setio, F. Ciampi, G. Litjens, P. Gerke, C. Jacobs, S. J. van Riel, M. M. W. Wille, M. Naqibullah, C. I. Sanchez, and B. van Ginneken, "Pulmonary nodule detection in ct images: false positive reduction using multi-view convolutional networks," *IEEE transactions on medical imaging*, 2016; vol. 35, no. 5, pp. 1160-1169.
- [24] H. Chen, Q. Dou, D. Ni, J.-Z. Cheng, J. Qin, S. Li, and P.-A. Heng, "Automatic fetal ultrasound standard plane detection using knowledge transferred recurrent neural networks," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015; pp. 507-514.
- [25] M. Anthimopoulos, S. Christodoulidis, L. Ebner, A. Christe, and S. Mougiakakou, "Lung pattern classification for interstitial lung diseases using a deep convolutional neural network," *IEEE transactions on medical imaging*, 2016; vol. 35, no. 5, pp. 1207-1216.
- [26] H.-C. Shin, H. R. Roth, M. Gao, L. Lu, Z. Xu, I. Nogues, J. Yao, D. Mollura, and R. M. Summers, "Deep convolutional neural networks for computer-aided detection: Cnn architectures, dataset characteristics and transfer learning," *IEEE transactions on medical imaging*, 2016; vol. 35, no. 5, pp. 1285-1298.
- [27] N. Tajbakhsh, J. Y. Shin, S. R. Gurudu, R. T. Hurst, C. B. Kendall, M. B. Gotway, and J. Liang, "Convolutional neural networks for medical image analysis: Full training or fine tuning?" *IEEE transactions on medical imaging*, 2016; vol. 35, no. 5, pp. 1299-1312.
- [28] H. R. Roth, L. Lu, A. Farag, H.-C. Shin, J. Liu, E. B. Turkbey, and R. M. Summers, "Deeporgan: Multi-level deep convolutional networks for automated pancreas segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015; pp. 556-564.
- [29] H. Chen, Q. Dou, X. Wang, J. Qin, and P. A. Heng, "Mitosis detection in breast cancer histology images via deep cascaded networks," in *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- [30] H. R. Roth, L. Lu, A. Farag, H.-C. Shin, J. Liu, E. B. Turkbey, and R. M. Summers, "Deeporgan: Multi-level deep convolutional networks for automated pancreas segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015; pp. 556-564.
- [31] H. Chen, X. Qi, L. Yu, and P.-A. Heng, "Dcan: Deep contour-aware networks for accurate gland segmentation," *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2016;

- [32] P. Moeskops, M. A. Viergever, A. M. Mendrik, L. S. de Vries, M. J. Benders, and I. Išgum, "Automatic segmentation of mr brain images with a convolutional neural network," *IEEE transactions on medical imaging*, 2016; vol. 35, no. 5, pp. 1252–1261.
- [33] N. Tajbakhsh, S. R. Gurudu, and J. Liang, "Automatic polyp detection in colonoscopy videos using an ensemble of convolutional neural networks," in *Biomedical Imaging (ISBI), 2015 IEEE 12th International Symposium on*. IEEE, 2015; pp. 79–83.
- [34] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017;
- [35] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang.: "Unet++: Redesigning Skip Connections to Exploit Multiscale Features in Image Segmentation". *IEEE TMI*, 2020;
- [36] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017; pp. 5998–6008.
- [37] H. Li, P. Xiong, J. An, and L. Wang, "Pyramid attention network for semantic segmentation," *arXiv preprint arXiv:1805.10180*, 2018;
- [38] M. Holschneider, R. Kronland-Martinet, J. Morlet, and P. Tchamitchian, "A real-time algorithm for signal analysis with the help of the wavelet transform," in *Wavelets: Time-Frequency Methods and Phase Space*, 1989; pp. 289–297.
- [39] D. Jha, P. H. Smedsrud, M. Riegler, P. Halvorsen, T. de Lange, D. Johansen, and H. Johansen, "Kvasir-seg: A segmented polyp dataset," in *International Conference on Multimedia Modeling*. Springer, 2020; [Online]. Available: <https://datasets.simula.no/kvasir-seg/>.
- [40] J. Bernal, F. J. Sanchez, G. Fernandez-Esparrach, D. Gil, C. Rodriguez, and F. Vilarino, "Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians," *Computerized Medical Imaging and Graphics*, 2015, vol. 43, pp. 99–111.
- [41] F. Chollet et al., "Keras," 2015;
- [42] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard et al., "Tensorflow: A system for large-scale machine learning," in *Proceeding of {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI})*, 2016; pp. 265–283
- [43] D. Jha, P. H. Smedsrud, M. A. Riegler, D. Johansen, T. De Lange, P. Halvorsen, and H. D. Johansen, "ResUNet++: An advanced architecture for medical image segmentation," in *Proceeding of IEEE International Symposium on Multimedia (ISM)*, 2019; pp. 225–2255.
- [44] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet large scale visual recognition challenge," *Int. J. Comput. Vis.*, 2015; vol. 115, no. 3, pp. 211–252.
- [45] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2014; arXiv:1409.1556. [Online]. Available: <http://arxiv.org/abs/1409.1556>
- [46] K. Yun, J. Park, and J. Cho, "Robust human pose estimation for rotation via self-supervised learning," *IEEE Access*, 2020; vol. 8, pp. 32502–32517.
- [47] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, 2017; vol. 39, no. 12, pp. 2481–2495.
- [48] L.-C. Chen, Y. Yang, J. Wang, W. Xu, and A. L. Yuille, "Attention to scale: Scale-aware semantic image segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016; pp. 3640–3649.
- [49] 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)*, 2018; pp. 267–283.
- [50] J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, and H. Lu, "Dual attention network for scene segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019; pp. 3146–3154.
- [51] P. Zhang, W. Liu, H. Wang, Y. Lei, and H. Lu, "Deep gated attention networks for large-scale street-level scene segmentation," *Pattern Recognition*, 2019; vol. 88, pp. 702–714.
- [52] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang, "Residual attention network for image classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017; pp. 3156–3164.

- [53] K. Li, Z. Wu, K.-C. Peng, J. Ernst, and Y. Fu, "Tell me where to look: Guided attention inference network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018; pp. 9215–9223.
- [54] W. Li, K. Liu, L. Zhang, and F. Cheng, "Object detection based on an adaptive attention mechanism," *Scientific Reports*, 2020; vol. 10, no. 1, pp. 1–13.

# Figures

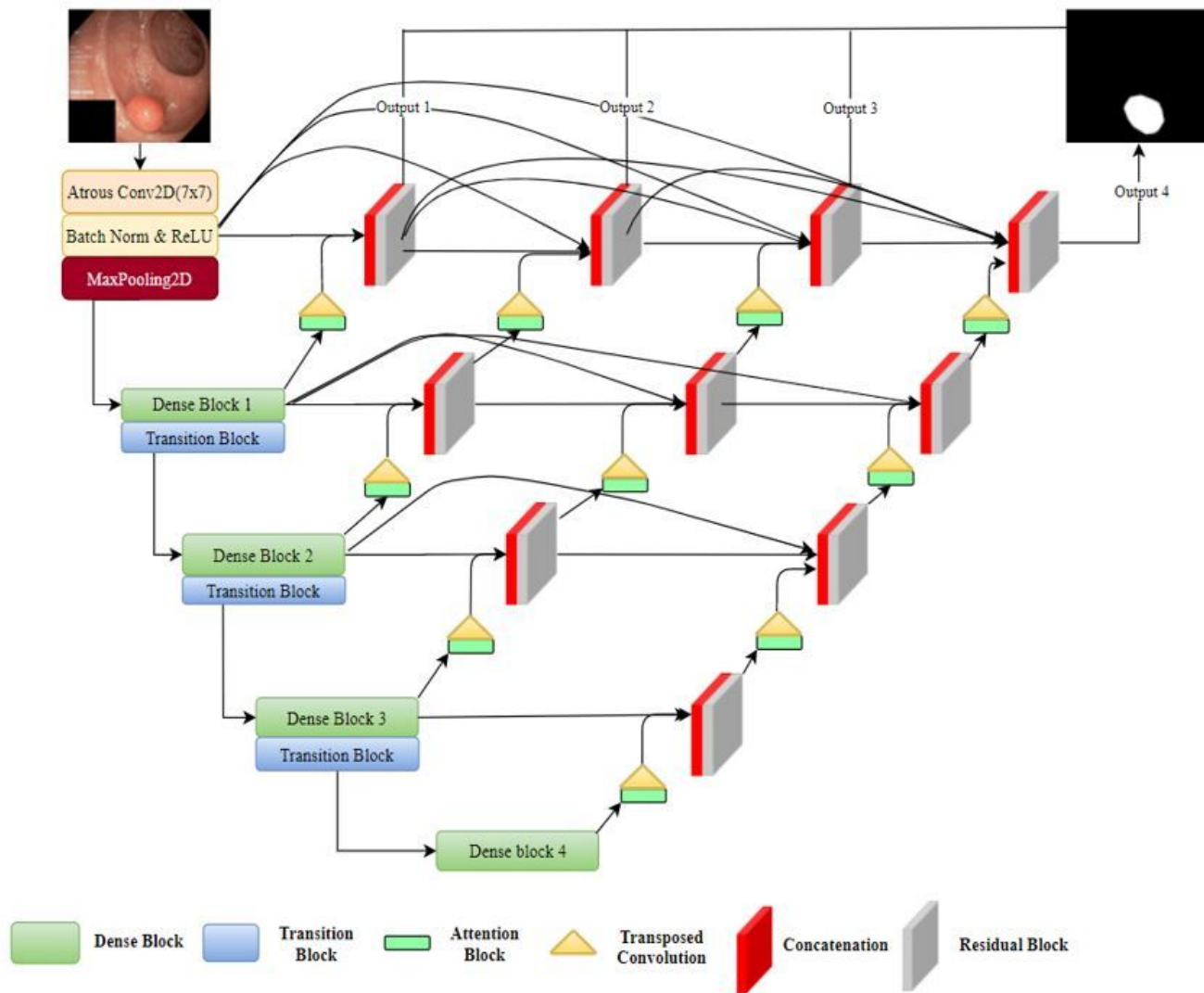


Figure 1

Block diagram of the proposed A-DenseUNet architecture: DenseNet is used as an encoder, Transposed convolution is performed for up-sampling between levels.

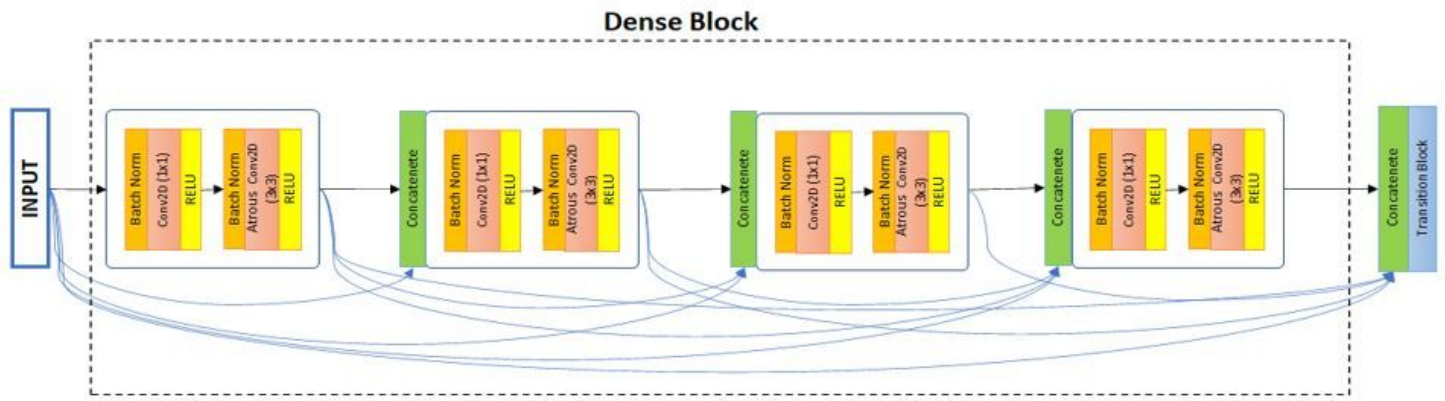


Figure 2

Five-layer dense block with grow rate  $n = 4$ . Each layer takes all previous information and includes batch normalization, atrous convolution, and ReLU activation.

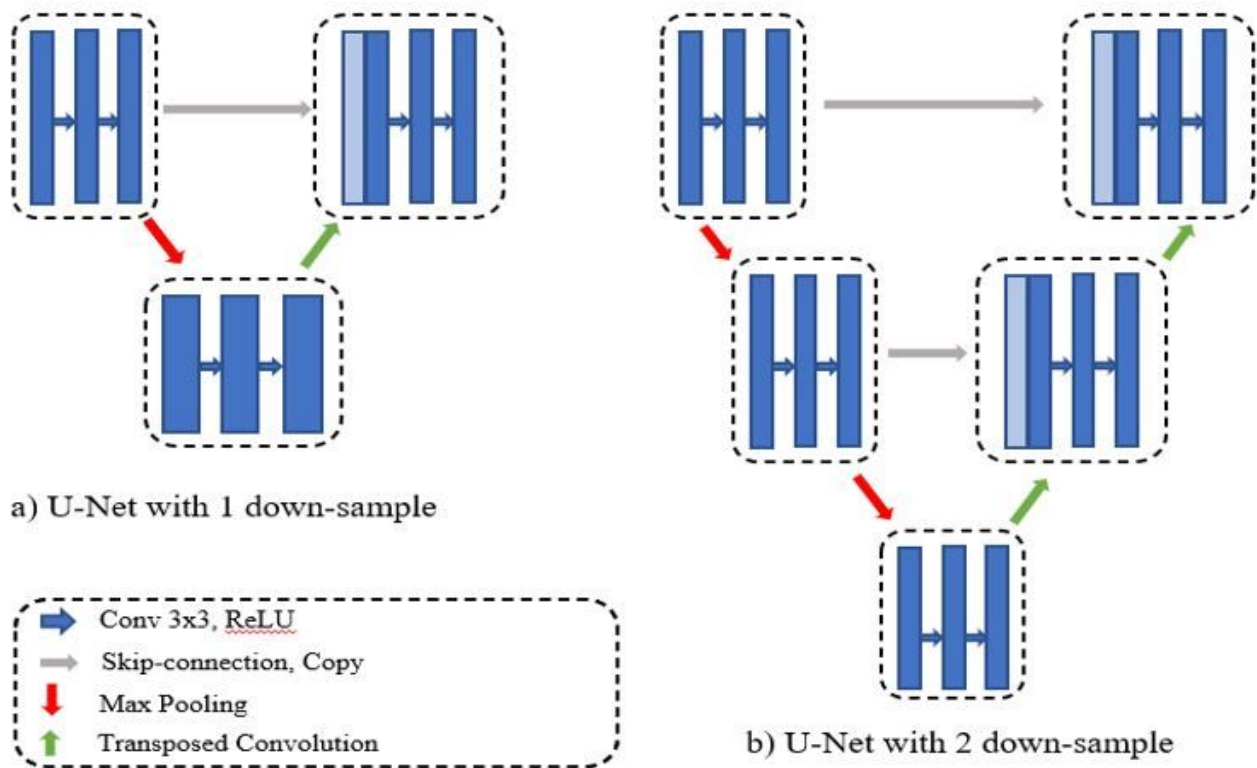
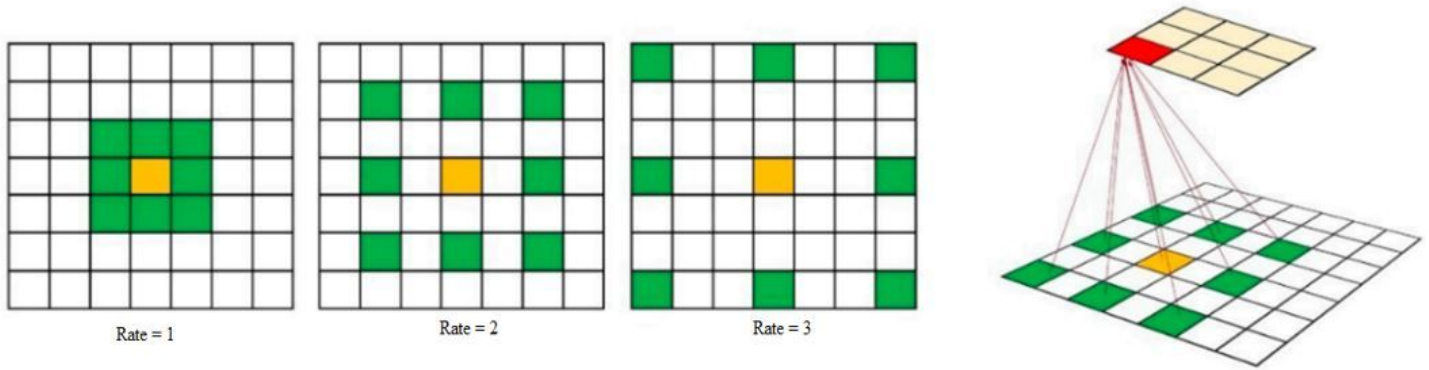


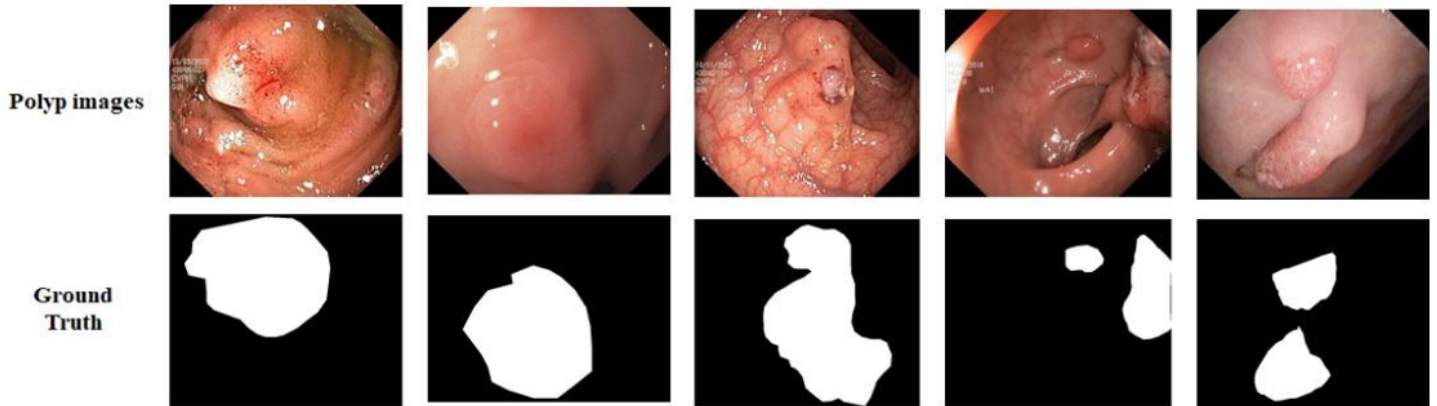
Figure 3

Multi-depth U-Net models.



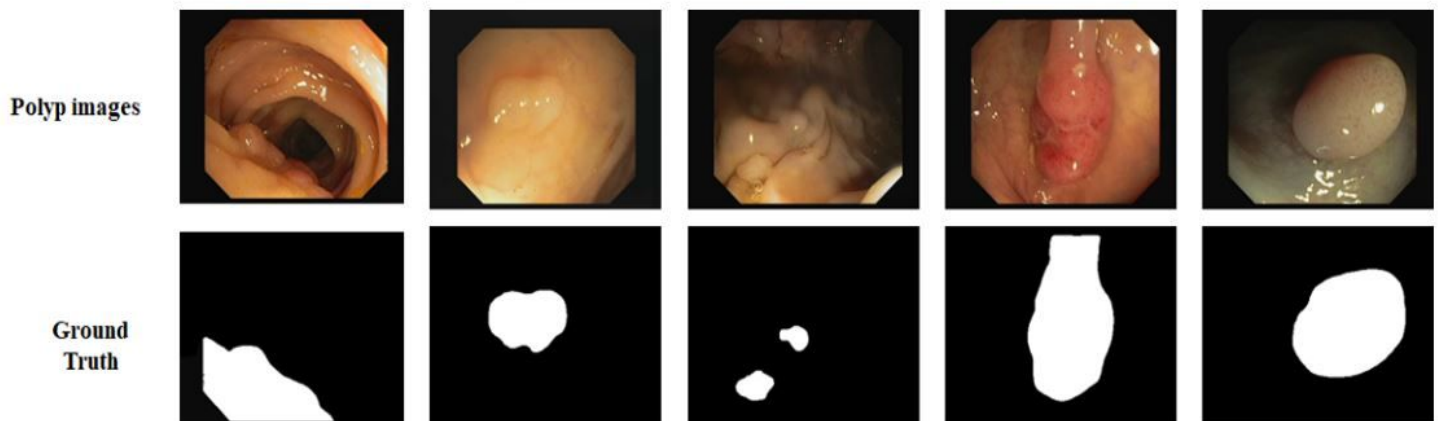
**Figure 4**

Dilated convolutions with different dilation rates. A dilation rate of one is normal convolution.



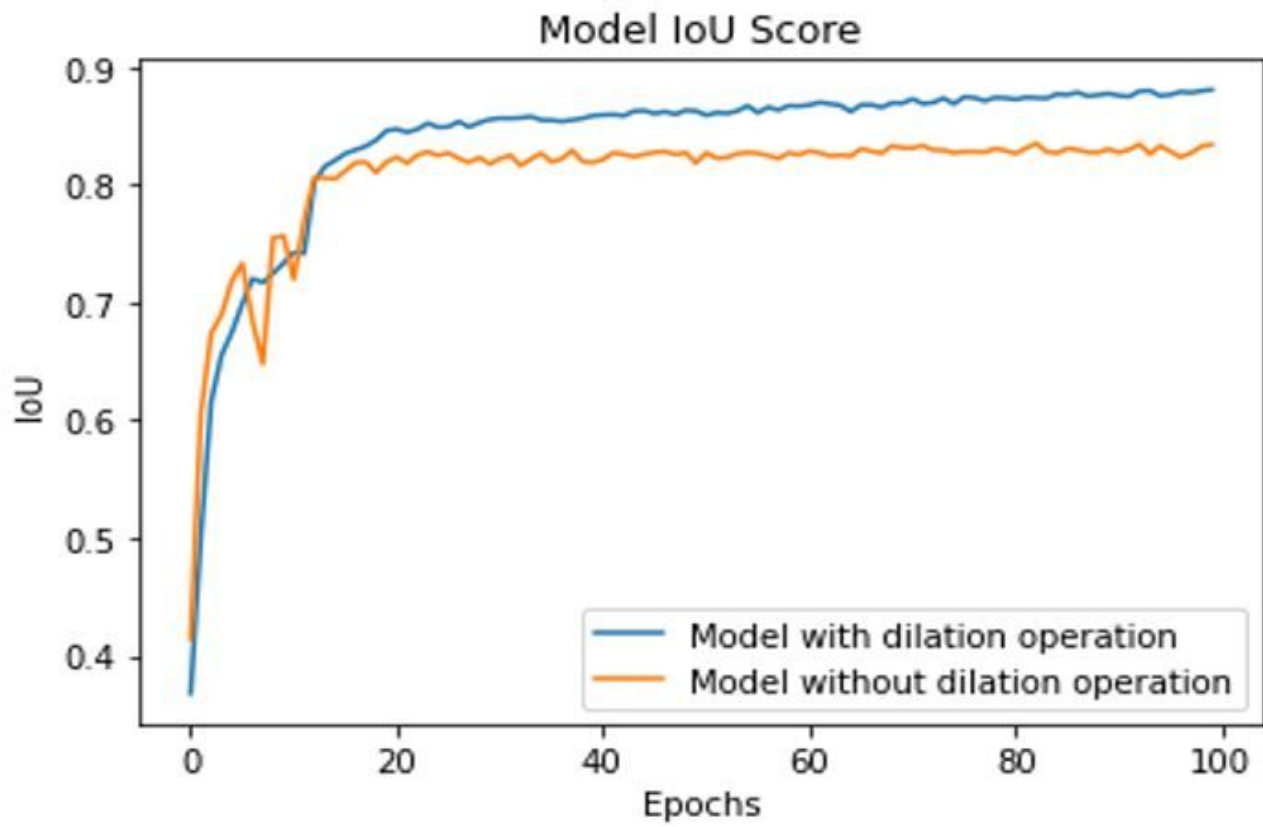
**Figure 5**

Example of data from Kvasir-SEG dataset. The first row shows original images and the second row presents their respective ground truth.



**Figure 6**

Images and ground truth masks from the CVC-612 dataset.



**Figure 7**

IoU score with and without dilated convolution.

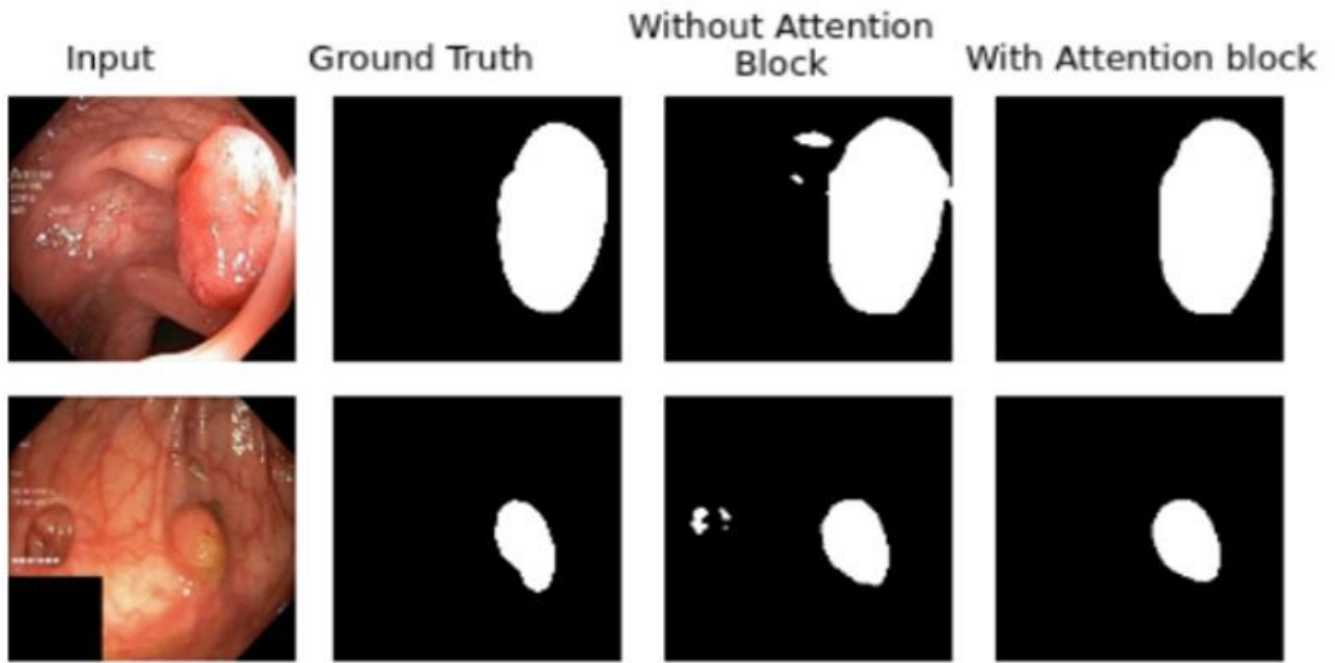


Figure 8

Effect of attention block in the network. By adding this, we were able to suppress the irrelevant regions.

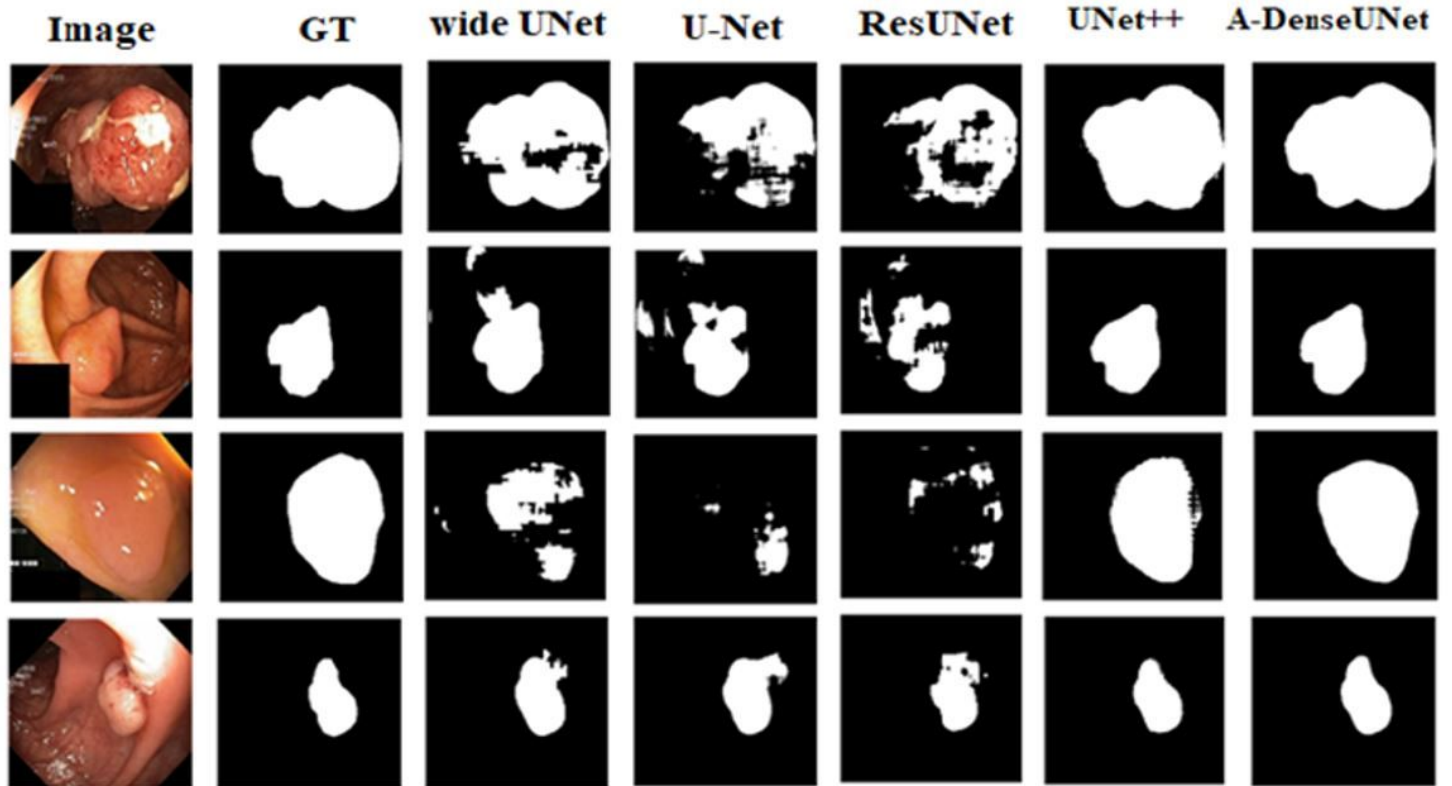


Figure 9



Qualitative segmentation results of various models on the Kvasir-SEG dataset. Experimental results show that A-DenseUNet produces better segmentation masks than other state-of-the-art networks.

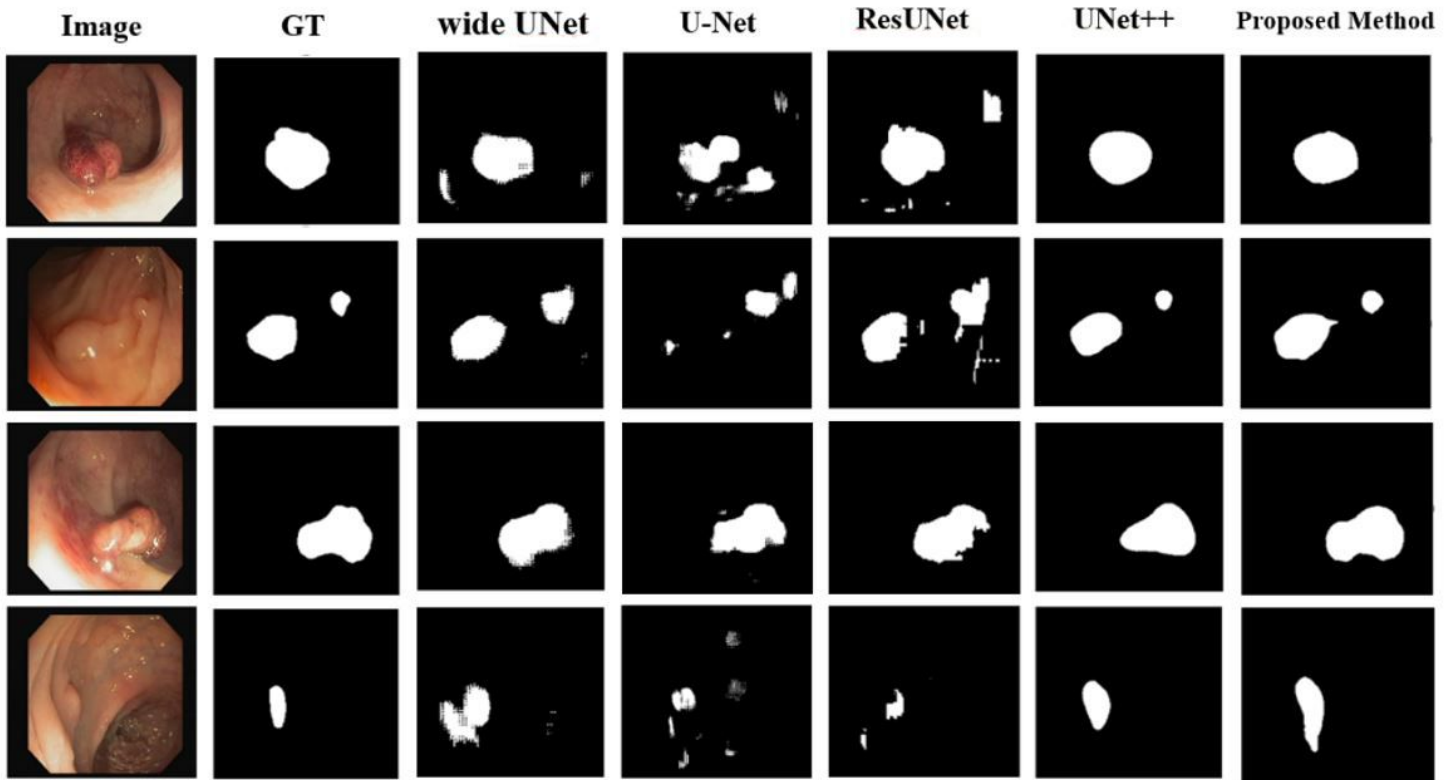


Figure 10

Qualitative segmentation results of various models on the CVC-612 dataset.