

# **Breast Ultrasound Images Segmentation Using Deep Neural Networks**

by

Khang Trong Nguyen

B.Sc., University of Texas at Dallas, 2020

A Report Submitted in Partial Fulfillment of the  
Requirements for The Degree of

**MASTER OF ENGINEERING**

in the Department of Electrical and Computer Engineering

©Khang Trong Nguyen,

2024

University of Victoria

All rights reserved. This report may not be reproduced in whole or in part,  
by photocopy or other means, without the permission of the author

# **Breast Ultrasound Images Segmentation Using Deep Neural Networks**

by

Khang Trong Nguyen

B.Sc., University of Texas at Dallas, 2020

## Supervisory Committee

Dr. Daler Rakhmatov, Supervisor

Department of Electrical and Computer Engineering

Dr. Kin F. Li, Departmental Member

Department of Electrical and Computer Engineering

## Abstract

Breast Ultrasound (US) imaging has emerged as an important diagnostic technique for detecting and characterizing breast tumors. Accurate segmentation of breast US images plays an essential role in enhancing the efficiency and precision of clinical assessments. This report explores the application of several well-known deep neural networks to the breast US image segmentation task. Specifically, we train and evaluate the following five models: SegNet, U-Net, and DeepLab V3+ with three different backbones (ResNet-18, ResNet-50, and Xception).

The presented results are based on two labeled datasets. One is Breast US Images (BUSI) dataset, which was used for training, validation, and testing. The other is Breast US Lesions (BUL) dataset, which was used exclusively for testing. Data augmentation was applied to increase the number and diversity of the data samples by randomly varying the contrast, brightness, and gamma of US images.

The performance of each model was evaluated based on Global Accuracy, Mean Accuracy, Mean Intersection-over-Union (IoU), Weighted IoU, Mean Boundary-F1 (BF) score, Average Dice score of Background, Average Dice score of Tumor, Mean Dice score of Tumor, and the model's cost. Overall, our results showed that Xception-based DeepLab V3+ and U-Net outperformed the other models under consideration when segmenting breast US images.

# Content

|   |     |
|---|-----|
| <b>Abstract</b> .....                                       | iii |
| <b>Content</b> .....  | iv  |
| <b>List of Tables</b> .....                                 | vi  |
| <b>List of Figures</b> .....                                | vii |
| <b>Glossary</b> .....                                       | ix  |
| <b>Acknowledgements</b> .....                               | x   |
| <b>Chapter 1. Introduction</b> .....                        | 1   |
| <b>1.1 Ultrasound imaging</b> .....                         | 1   |
| <b>1.2 Segmentation in breast ultrasound images</b> .....   | 4   |
| 1.2.1 Thresholding-based segmentation. ....                 | 4   |
| 1.2.2 Clustering-based segmentation .....                   | 5   |
| 1.2.3 Segmentation based on active contour model.....       | 6   |
| 1.2.4 Segmentation using deep neural networks .....         | 7   |
| <b>1.3 Deep Neural Networks</b> .....                       | 8   |
| 1.3.1 Basic Structure of DNNs.....                          | 8   |
| 1.3.2 Activation functions .....                            | 9   |
| 1.3.3 Parameter learning .....                              | 10  |
| <b>1.4 Report contribution and organization</b> .....       | 10  |
| <b>Chapter 2. US Datasets and Evaluation Criteria</b> ..... | 12  |
| <b>2.1 Breast US Images (BUSI) dataset</b> .....            | 12  |
| <b>2.2 Breast US Lesions (BUL) dataset</b> .....            | 12  |
| <b>2.3 Segmentation quality criteria</b> .....              | 13  |
| 2.3.1 Accuracy .....  | 14  |
| 2.3.2 BF score.....   | 14  |
| 2.3.3 Intersection over union (IoU).....                    | 15  |
| 2.3.4 Dice score.....                                       | 15  |
| <b>Chapter 3. Segmentation Networks</b> .....               | 16  |
| <b>3.1 DeepLab V3+ architecture</b> .....                   | 16  |
| 3.1.1 ResNet-18.....  | 18  |
| 3.1.2 ResNet-50.....  | 20  |
| 3.1.3 Xception .....  | 21  |

|  |    |
|--|----|
| <b>3.2 U-Net architecture</b> .....                                | 24 |
| <b>3.3 SegNet architecture</b> .....                               | 26 |
| <b>Chapter 4. Network Training</b> .....                           | 28 |
| <b>4.1 Dataset preparation</b> .....                               | 28 |
| 4.1.1 Data splitting, grayscale conversion, and resizing.....      | 28 |
| 4.1.2 Data augmentation .....                                      | 28 |
| <b>4.2. Training settings</b> .....                                | 30 |
| 4.2.1. Experimental parameters .....                               | 30 |
| 4.2.2. Learning option parameters .....                            | 31 |
| <b>4.3 Learning Curves</b> .....                                   | 35 |
| 4.3.1 DeepLab V3+ with ResNet-18 bondnet.....                      | 35 |
| 4.3.2 DeepLab V3+ with ResNet-50 bondnet.....                      | 36 |
| 4.3.3 DeepLab V3+ with Xception bondnet .....                      | 37 |
| 4.3.4 SegNet architecture .....                                    | 37 |
| 4.3.5 U-Net architecture .....                                     | 38 |
| <b>Chapter 5. Evaluation Results</b> .....                         | 40 |
| <b>5.1 Performance of DeepLab V3+ with ResNet-18 bondnet</b> ..... | 40 |
| <b>5.2 Performance of DeepLab V3+ with ResNet-50 bondnet</b> ..... | 42 |
| <b>5.3 DeepLab V3+ with Xception bondnet</b> .....                 | 44 |
| <b>5.4 Performance of SegNet</b> .....                             | 45 |
| <b>5.5 Performance of U-Net</b> .....                              | 47 |
| <b>5.6 Comparative summary</b> .....                               | 49 |
| <b>Chapter 6. Conclusion and Future Work</b> .....                 | 53 |
| <b>6.1. Conclusion</b> .....                                       | 53 |
| <b>6.2. Future work</b> .....                                      | 53 |
| <b>Bibliography</b> .....  | 55 |

## List of Tables

|  |    |
|--|----|
| Table 4.1 Experimental hardware and software specifications.....             | 30 |
| Table 4.2 DeepLab V3+ training option hyperparameters.....                   | 31 |
| Table 4.3 U-Net training option hyperparameters. ....                        | 32 |
| Table 4.4 SegNet training option hyperparameters.....                        | 33 |
| Table 5.1 Evaluation results for DeepLab V3+ with ResNet-18 bondnet.....     | 41 |
| Table 5.2 Evaluation results for DeepLab V3+ with ResNet-50 bondnet.....     | 43 |
| Table 5.3 Evaluation results for DeepLab V3+ with Xception bondnet.....      | 44 |
| Table 5.4 Evaluation results for SegNet. ....                                | 46 |
| Table 5.5 Evaluation results for U-Net. ....                                 | 48 |
| Table 5.6 Latency, size, and number of MACs associated with each model. .... | 51 |

## List of Figures

|   |    |
|---|----|
| Figure 1.1 A-mode sonogram of a renal calculus [33].....  | 2  |
| Figure 1.2 B-mode ultrasound image of breast tumor from BUSI dataset [20].....  | 2  |
| Figure 1.3 C-mode ultrasound image of carotid artery [43]. .....  | 3  |
| Figure 1.4 M-Mode imaging provides a visual representation of how the tissues move over time, allowing for detailed analysis and monitoring of their dynamics [42]. ..... | 3  |
| Figure 1.5 Breast US image segmentation using thresholding [6]. .....   | 5  |
| Figure 1.6 Breast US image segmentation using K-means clustering [4]. .....   | 6  |
| Figure 1.7 Breast US image segmentation using active contour model [7]. .....   | 7  |
| Figure 1.8 The structure of DL model [16].....  | 8  |
| Figure 1.9 Basics of (Deep) Neural Network [17].....  | 10 |
| Figure 2.1 Breast US Images (BUSI) dataset samples [20].....  | 12 |
| Figure 2.2 Example images with a cyst (a), a fibroadenoma lesion (b), and an invasive ductal carcinoma lesion (c) shown [21].....   | 13 |
| Figure 2.3 FN, FP, and TP [32].....   | 14 |
| Figure 3.1 Encoder-Decoder with atrous convolution and spatial pyramid pooling [22].....  | 16 |
| Figure 3.2 DeepLab V3+ architecture [22].....   | 18 |
| Figure 3.3 Atrous separable convolution illustrated (c) where atrous convolution is adopted for depthwise convolution (a) with rate = 2 [22].....                         | 19 |
| Figure 3.4 Architecture of typical ResNet model [30].....   | 19 |
| Figure 3.5 DeepLab V3+ with ResNet-18 bondnet as implemented in MATLAB.....   | 20 |
| Figure 3.6 Left: regular residual block. Right: bottleneck residual block. [24].....  | 20 |
| Figure 3.7 DeepLab V3+ with ResNet-50 bondnet as implemented in MATLAB.....   | 21 |
| Figure 3.8 A canonical Inception module (Inception V3) [29]. .....  | 22 |
| Figure 3.9 Architecture of Xception model [30].....   | 23 |
| Figure 3.10 DeepLab V3+ with Xception bondnet as implemented in MATLAB. ....  | 23 |
| Figure 3.11 Architecture of U-Net model [30]. .....   | 24 |
| Figure 3.12 U-Net architecture as implemented in MATLAB.....  | 25 |
| Figure 3.13 Architecture of SegNet model [31].....  | 26 |
| Figure 3.14 SegNet architecture as implemented in MATLAB.....   | 27 |
| Figure 4.1 The left top corner is the original image. The others are augmented images using randomized brightness adjustment in the range [-0.5, 0.5]......               | 29 |
| Figure 4.2 The left top corner is the original image. The others are augmented images using randomized contrast adjustment in the range [0.5, 1.5]......                  | 29 |

|  |    |
|--|----|
| Figure 4.3 The left corner is original image. The others are augmented images using randomized gamma adjustment in the range [0.5, 1.5]. | 30 |
| Figure 4.4 Learning curves for DeepLab V3+ based on ResNet-18.   | 36 |
| Figure 4.5 Learning curves for DeepLab V3+ based on ResNet-50 during training in MATLAB  | 37 |
| Figure 4.6 Learning curve for DeepLab V3+ based on Xception during training in MATLAB.   | 38 |
| Figure 4.7 Learning curve for SegNet during training in MATLAB   | 39 |
| Figure 4.8 Learning curves for U-Net during training in MATLAB.  | 40 |
| Figure 5.1 Evaluation results for DeepLab V3+ with ResNet-18 bondnet.  | 41 |
| Figure 5.2 Evaluation results for DeepLab V3+ with ResNet-50 bondnet.  | 43 |
| Figure 5.3 Evaluation results for DeepLab V3+ with Xception bondnet.   | 46 |
| Figure 5.4 Evaluation results for SegNet.  | 47 |
| Figure 5.5 Evaluation results for U-Net.   | 49 |
| Figure 5.6 Network comparison results: Unseen BUL dataset.   | 49 |
| Figure 5.7 Network comparison results: Test portion of BUSI dataset.   | 50 |

## Glossary

|      |                              |
|------|------------------------------|
| Adam | Adaptive moment estimation   |
| AI   | Artificial intelligence      |
| ANN  | Artificial Neural Network    |
| BF   | Boundary F1                  |
| BUL  | Breast US Lesions            |
| BUS  | Breast Ultrasound            |
| BUSI | Breast US Image              |
| CAD  | Computer-aided diagnosis     |
| CNN  | Convolutional neural network |
| DL   | Deep Learning                |
| DNNs | Deep Neural Networks         |
| FN   | False Negative               |
| FP   | False Positive               |
| IoU  | Intersection over union      |
| LR   | Learning rate                |
| NN   | Neural networks              |
| ReLU | Rectified linear unit        |
| TP   | True Positive                |
| US   | Ultrasound                   |

## **Acknowledgements**

First, I very much thank God for his countless blessings. I am also thankful to Uyen Nguyen, my spouse, who gave me strength, encouragement, and capability to get things done, as well as my parents, my sister, my mother-in-law and sister-in-law for their unconditional support, wishes, prayers, and encouragement.

I would like to thank my supervisor, Dr. Daler Rakhmatov, for taking me as a student and all his support, guidance, mentorship, flexibility, patience, and encouragement throughout this program.

I am also grateful to the University of Victoria for providing me with a professional environment to learn and grow with top-notch facilities.

# Chapter 1. Introduction

## 1.1 Ultrasound imaging

Ultrasound (US) is a safe and affordable medical imaging modality that allows for the examination of internal organs without invasive surgery. This non-invasive technique has garnered considerable attention from the research community due to its immense diagnostic value. Creating accurate and easily interpretable images is the key to the early detection and treatment of conditions such as cancer, tumors, and other critical diseases.

To better understand how US imaging works, we need to introduce the basic principles of US imaging. US refers to sound waves that are not detectable by the human ear with frequencies greater than 20,000 cycles/sec (Hz). Diagnostic US commonly uses frequencies between 2 and 15 MHz. Intravascular transducers commonly use frequencies up to 30 MHz, and US bio-microscopy transducers frequencies up to 100 MHz [2]. The basic principles of pulse-echo US imaging remain largely unchanged today compared to several decades ago. This process entails transmitting US pulses from a transducer into the body. As these US waves penetrate body tissues of varying acoustic impedances along their transmission path, some are reflected back to the transducer as echo signals, while others continue to penetrate deeper. The echo signals are then processed and combined to generate an image [10].

There are multiple modes and techniques that can capture the interaction between US waves and tissues. Four of the most frequently used ones are briefly described below.

- A-mode, also known as amplitude mode, is the simplest US technique that uses a single transducer to scan one line through the body and producing a plot of echo amplitudes as a function of depth. It is commonly used in therapeutic US procedures targeting specific tumors or calculi, allowing for pinpoint accurate focusing of destructive wave energy [14]. Fig. 1.1 is an example of the A-mode US image showing a single spike representing a reflection from the renal calculus [33].

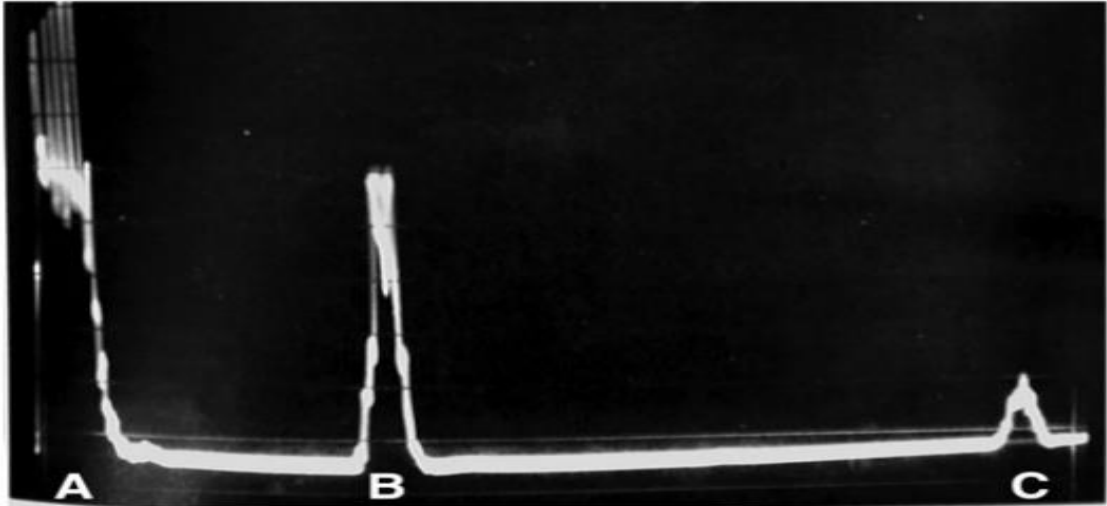


Figure 1.1 A-mode sonogram of a renal calculus [33].

- B-mode, also known as brightness mode, is the most common mode of US imaging, employing a linear array of transducers to scan a plane of tissue throughout the body, thereby providing a two-dimensional image [14]. The datasets used for training, testing, and evaluation are typically obtained from B-mode images. One such image is shown in Fig. 1.2.

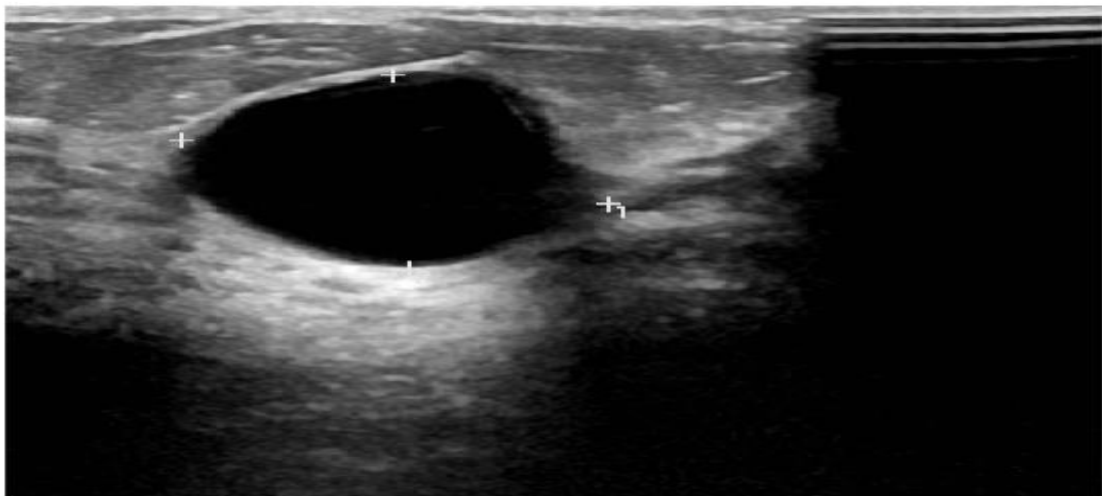


Figure 1.2 B-mode ultrasound image of breast tumor from BUSI dataset [20].

- C-mode, similar to B-mode, produces a two-dimensional image. It combines aspects of both A-mode and B-mode by gathering data from two orthogonal B-mode slices, which are then visualized simultaneously parallel to the face of the transducer at any chosen depth as shown in Fig. 1.3. This enables the plotting of a two-dimensional image slice at the selected depth [14], [15].

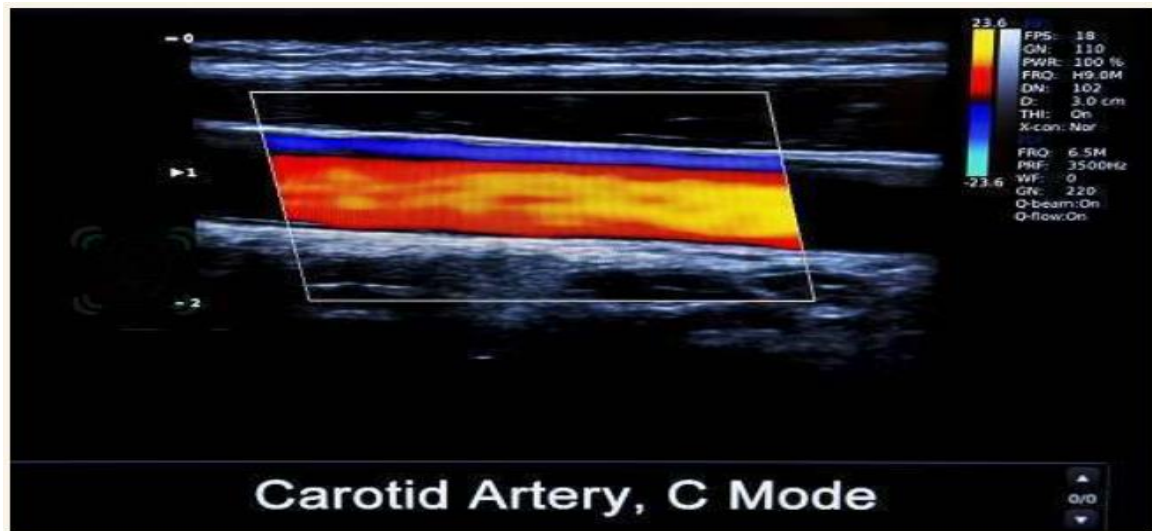


Figure 1.3 C-mode ultrasound image of carotid artery [43].

- M-mode, also known as motion mode, operates continuously over time. The transducer emits pulses continuously, capturing a series of either A-mode or B-mode images, thereby generating a 'video' of the ultrasound. This mode is often employed for velocity estimation [11], [15]. In M-mode, images are rapidly displayed in succession on a monitor, enabling real-time observations. Fig. 1.4 shows an example image snapshot. By isolating anatomy along a single line, changes in depth can be tracked over time, facilitating precise and accurate monitoring of the movement of structures or instruments by healthcare professionals.

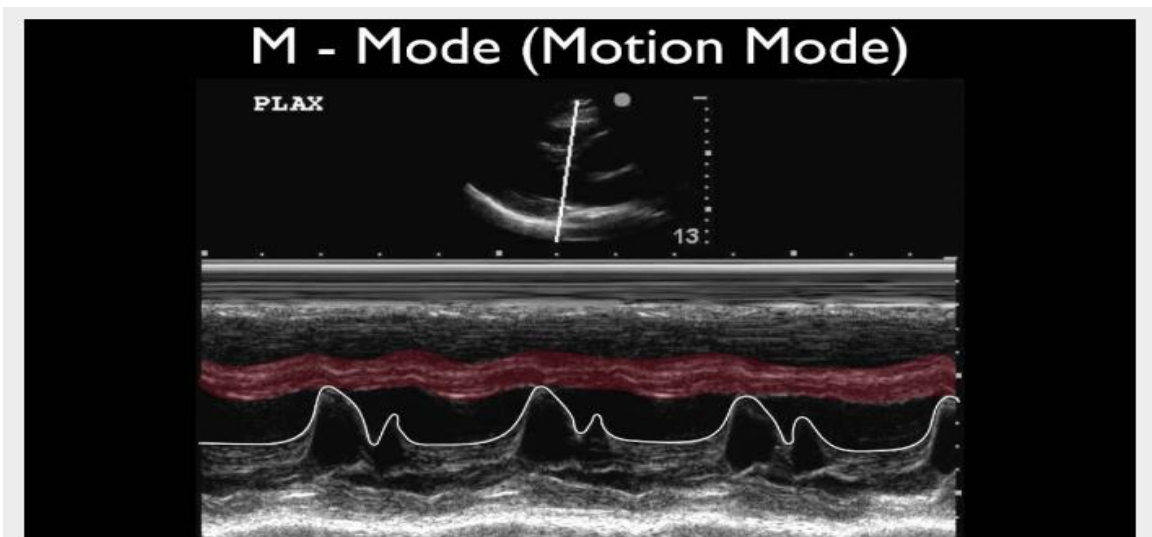


Figure 1.4 M-Mode imaging provides a visual representation of how the tissues move over time, allowing for detailed analysis and monitoring of their dynamics [42].

## 1.2 Segmentation in breast ultrasound images

Breast cancer is the most common form of cancer among women worldwide, with more than 8% of women experiencing this disease during their lifetime [3]. Since the causes of breast cancer remain unknown, early detection is key to reducing the death rate (40% or more) [3]. Therefore, early detection of breast cancer increases the chances of patient survivability and treatment options. One of the most popular diagnostic techniques for detecting and identifying breast abnormalities is US imaging. Recently, computer-aided diagnosis (CAD) systems using US images, and their segmentations have been developed to assist radiologists in improving diagnostic precision. However, due to the variations in ultrasonic elements, effective segmentation of US images remains a challenging task. A breast cancer CAD system based on US imaging typically involves the following four steps: preprocessing, segmentation, feature extraction, and classification [4].

Image segmentation involves the task of assigning a specific label to every pixel within an image, aiming to group pixels with similar characteristics into coherent regions or objects. The primary aim is to divide the image into discrete segments that exhibit uniform properties, such as color, texture, intensity, or other relevant visual attributes, facilitating subsequent analysis and interpretation of the image data. In this work, our segmentation task translates into a binary classification task, where each pixel in a grayscale ultrasound image is classified as either tumor or background. In this section, we briefly review several popular methods for segmenting images.

### 1.2.1 Thresholding-based segmentation

Thresholding is one of the popular techniques for monochrome image segmentation and has been widely applied to segmenting BUS images [5]. It is a straightforward yet powerful technique that assigns a label or value to each pixel within an image based on a predetermined threshold value or a set of threshold values. Pixels surpassing or falling below the specified threshold(s) are categorized as foreground or background, respectively. Therefore, the selection of an appropriate threshold value (or a set of threshold values) is crucial. Various methodologies have been developed to address this requirement in thresholding-based segmentation, including the maximum entropy method and Otsu's method (maximum variance).

Fig. 1.5 presents an example of thresholding-based segmentation as proposed in [6]. On the left, the original breast US image is displayed, providing the input data for analysis. In the

middle, the segmentation result obtained through the thresholding-based segmentation method is showcased, where pixels are categorized into foreground and background based on specified intensity thresholds. Finally, on the right, the region of interest (ROI) is depicted after applying the morphological chain method. The latter involves the application of morphological operations such as erosion, dilation, and closing to refine and enhance the segmented regions, ultimately producing a more accurate representation of the desired anatomical structures within the image.

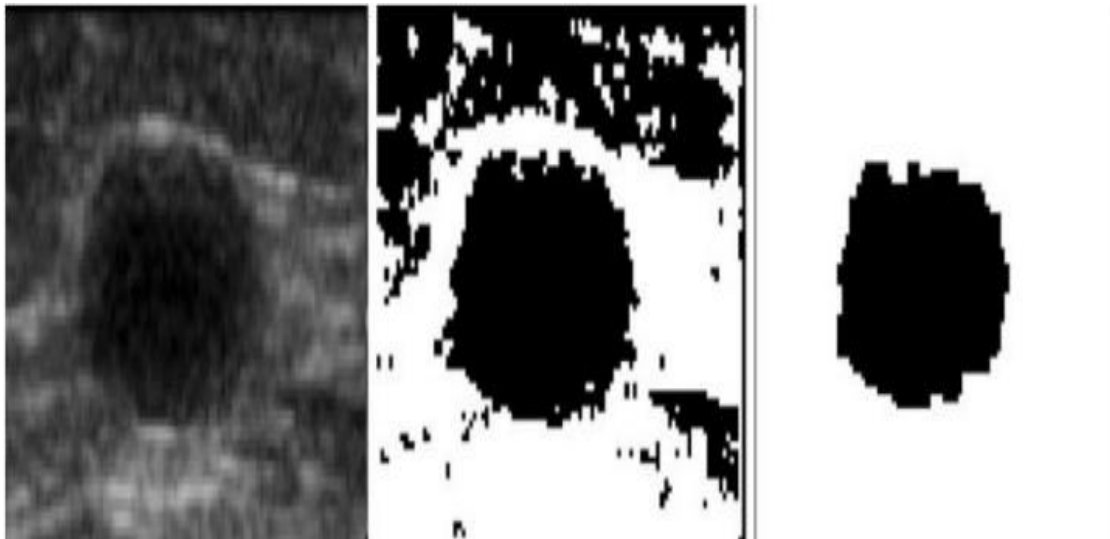


Figure 1.5 Breast US image segmentation using thresholding [6].

### 1.2.2 Clustering-based segmentation

Clustering is a technique used for image segmentation that groups pixels or regions together based on their similarity in terms of certain features or characteristics. It aims to identify natural clusters or groups within an image without explicitly defining specific thresholds. Pixels or regions are grouped into clusters based on the similarity of their feature values.

One of the popular methods used in clustering-based segmentation is the K-means algorithm. Its goal is to partition an image into K clusters. There are several key steps. Initially, K cluster centers are chosen, a process that can either be randomized or guided by specific heuristics. Subsequently, each pixel within the image is assigned to the cluster that minimizes the distance between the pixel and the cluster center. After this step, the cluster centers are updated by recalculating them based on the average of all pixels within the respective cluster. This iterative process continues until convergence is reached, signified by the absence of any

pixel changing clusters. The distance between a pixel and the cluster center is typically computed as the squared or absolute difference, based on factors such as pixel color, intensity, texture, and spatial location, or a weighted combination of these attributes. The selection of  $K$ , which represents the number of clusters, can be done manually, randomly, or through the application of a heuristic approach. The effectiveness of the solution is contingent upon both the initial set of clusters and the chosen value of  $K$ , highlighting the significance of these parameters in determining the quality of the clustering outcome [4].

Fig. 1.6 shows an example of  $K$ -means clustering application. The left side of the figure is the original breast US image, and on the right side, the segmentation result obtained using the  $K$ -means algorithm.

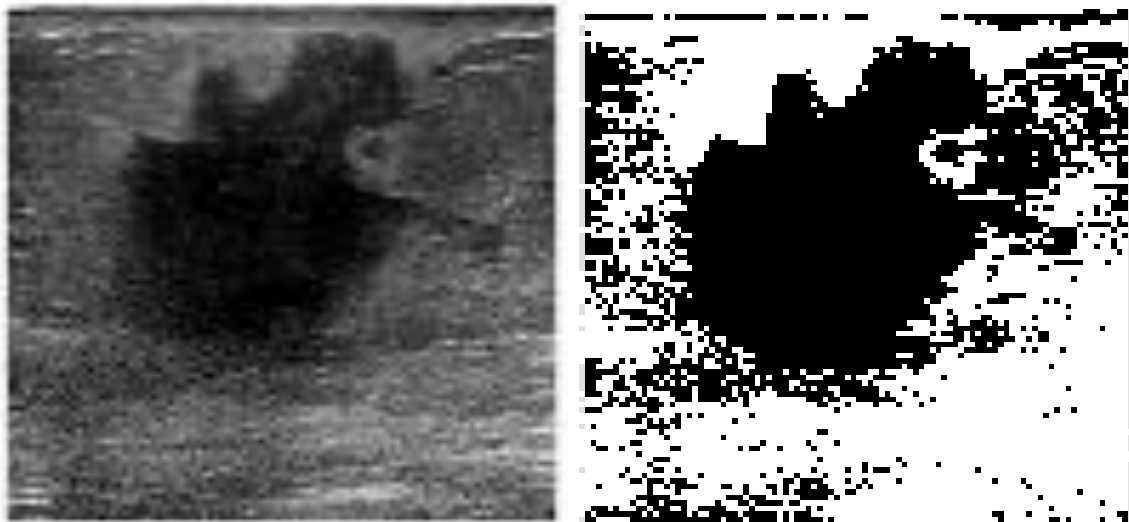


Figure 1.6 Breast US image segmentation using  $K$ -means clustering [4].

### 1.2.3 Segmentation based on active contour model

The active contour model, commonly referred to as a snake, is a popular segmentation technique for US images. It provides a framework for delineating an object outline from a possibly noisy 2D image and has been widely utilized as an edge-based segmentation method [4]. Typically, the contour is represented either as a parametric curve or as a series of interconnected points, enabling precise delineation of object boundaries within the image.

The active contour model leverages the energy minimization principle to iteratively deform the initial contour towards the desired object boundary. It defines an energy function comprising two key components: internal energy and external energy. The internal energy, derived from the contour model, controls the shape and regularity of the contour, while the

external energy, derived from image features, extracts the contour of the desired object [4]. This model is achieved by applying optimization techniques, such as gradient descent, to deform the contour towards total energy minimization, incorporating both internal and external energy.

The active contour model method offers the advantage of ensuring that the detected edge is closed and continuous. However, its disadvantage lies in the need for accurately initialized points. Additionally, obtaining a reasonable energy function can be challenging due to the quality of breast US images [7]. Fig. 1.7 illustrates some segmentation results based on the active contour model. The first row displays the original images overlaid with initial contours, serving as the starting boundaries for segmentation. In the second row, manual segmentation results are shown, providing boundaries delineated manually as a reference. Finally, the third row demonstrates the segmentation outcomes achieved through the active contour model approach.

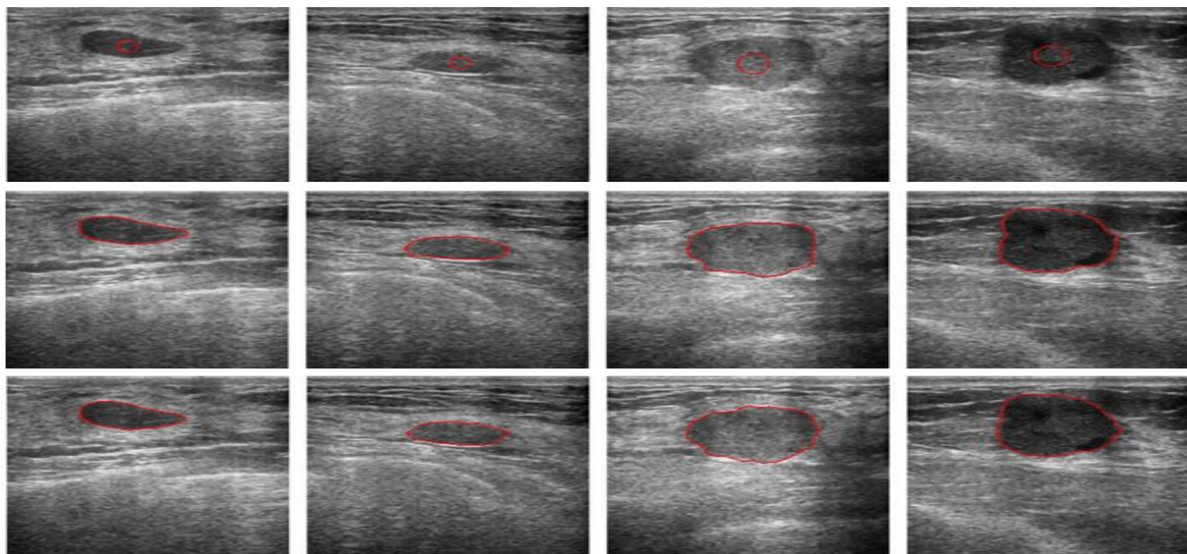


Figure 1.7 Breast US image segmentation using active contour model [7].

#### 1.2.4 Segmentation using deep neural networks

Learning-based segmentation methods, transforming the segmentation problem into a classification decision based on a set of input features, have proven to be highly accurate [8]. In 1943, Warren Mcculloch and Walter Pitts first investigated the concept of the Artificial Neural Network (ANN) and proposed a mathematical model of an artificial neuron [11]. In 2006, Hinton et al. [12] proposed the concept of Deep Learning (DL). They argued that an

ANN with multiple layers was in possession of extraordinary feature learning ability, allowing the learned feature data to represent the essence of the original data.

In this report, we applied three common types of neural network-based segmentation methods for breast US images. The corresponding details are provided in section 1.4 and Chapter 3, while some basic background on deep neural networks (DNNs) is given in the next section.

### 1.3 Deep Neural Networks

#### 1.3.1 Basic Structure of DNNs

Machine learning is a form of mathematical optimization that seeks to find optimally tuned parameters for a decision-making model based on available data. Machine learning algorithms can be roughly divided into three groups: supervised, unsupervised, and semi-supervised algorithms.

Deep learning refers to a subset of machine learning that allows computational models composed of multiple layers to learn data representations with multiple levels of abstraction [13]. The goal of DL is to construct neural networks suitable for various machine learning tasks, including regression, clustering, classification, segmentation, and generation problems. The basic structure of DNNs is presented in Fig. 1.8.

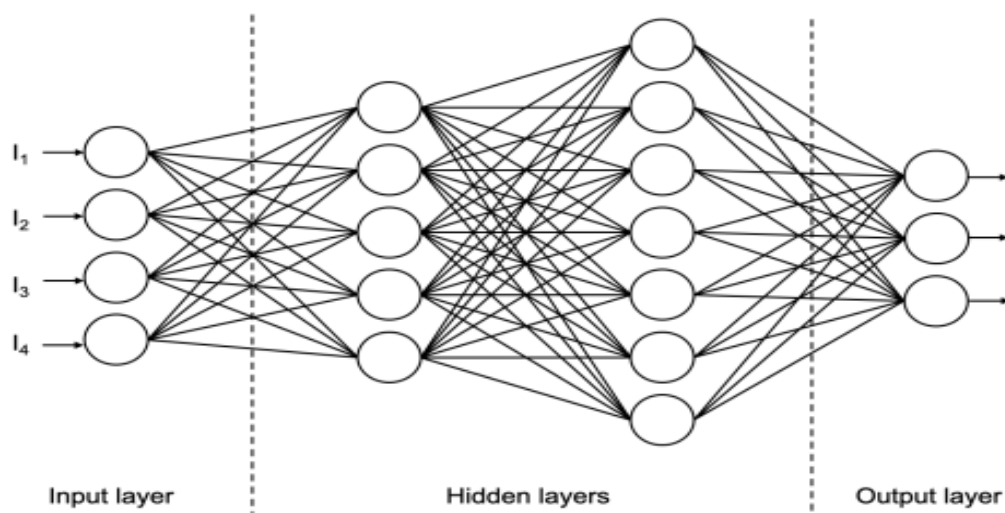


Figure 1.8 The structure of DL model [16].

The layers in DL networks can be divided into three groups: the input layer, the hidden layers, and the output layer, based on their position and functions. The first layer is the input

layer, where pre-processed data is fed into the model. The last layer is the output layer, where the model's result is obtained. The middle layers between the input and output layers are called hidden layers, and DL models typically consist of multiple hidden layers to improve their expressiveness and learn multi-level feature representations [16].

In DL models, the connections between layers vary. In Fig. 1.8, adjacent layers are fully connected, meaning that any neuron in the one layer is connected to any neuron in the next layer. Furthermore, in some complex DL structures, not only fully connected layers are used as hidden layers, but also layers composed of other types of neurons, such as convolutional layers, pooling layers, etc., can be used as hidden layers [16].

### 1.3.2 Activation functions

Another important factor in a DNN is activation functions, which are inspired by human neural activities. The activation function's role is to transform tightly clustered data into a more classifiable plane, thereby facilitating observation of the effects of different dimensions in the given problem [17].

For example, in Fig. 1.9 A, logistic regression can be conceptualized as a single neural unit, representing the probability of the decision made by that unit, with an output range from 0 to 1 determined by multiple inputs ( $X_1, X_2, X_3$ ). To facilitate logistic regression processing, activation functions  $a(Z)$ , such as sigmoid, hyperbolic tangent, and ReLU (see Fig. 1.9 B), are utilized. The sigmoid activation function, for instance, is a good fit for generating an output within the range of 0 and 1. Fig. 1.9 C illustrates a shallow network comprising one hidden layer containing three neural units. The final output represents the probability of decision derived from the output of each of the three neural units. Employing the sigmoid activation function, each of the three neural units and the final output generates outputs within the range  $[0,1]$ , where a value approaching 0 signifies a false outcome and a value nearing 1 denotes a true outcome. Fig. 1.9 D depicts the architecture of a deep neural network consisting of multiple hidden layers, each containing more than one neural unit [17]. This structure exemplifies how a deep neural network is essentially formed by stacking multiple shallow neural networks together.

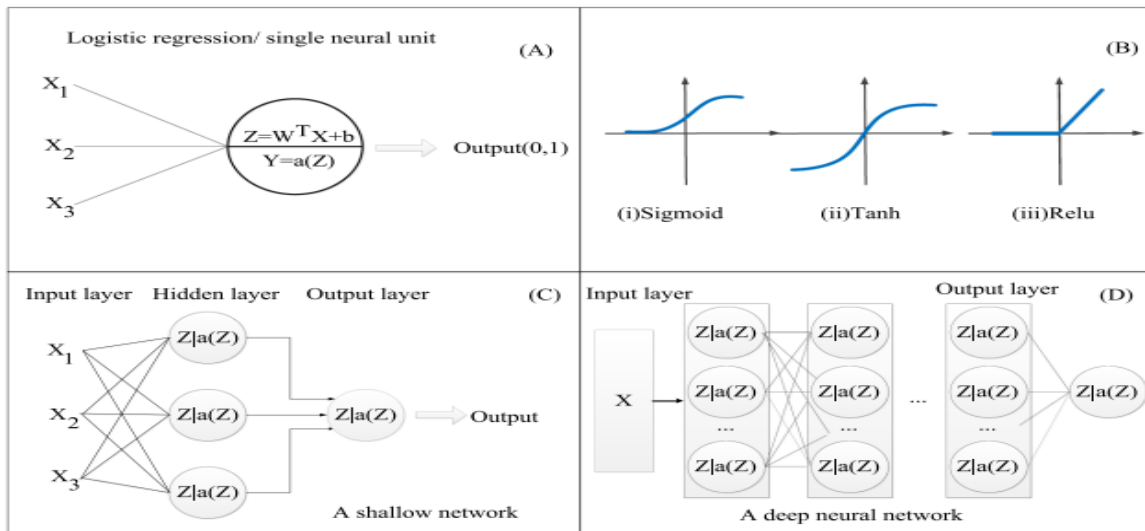


Figure 1.9 Basics of (Deep) Neural Network [17].

### 1.3.3 Parameter learning

DL models require parameter learning using mathematical algorithms such as gradient descent. The primary objective of gradient descent is to minimize the loss function by converging towards the global optimum. While gradient descent is suitable for convex functions, which possess single absolute minimum or maximum, non-convex functions present challenges due to multiple local optima. The learning parameter may get stuck in a locally optimal point and miss the global optimum [17]. Therefore, neural network optimization may or may not find the global maximum or minimum when trying to minimize the error between predicted and actual values.

### 1.4 Report contribution and organization

The key motivation for this work stemmed from the paper [22] by L. Chen et al., where the authors addressed the problem of semantic segmentation tasks by introducing the DeepLab V3+ architecture. This DNN architecture utilized an encoder-decoder structure with atrous separable convolution. Our project aimed to apply the architecture proposed in [22], as well as the U-Net [25] and SegNet [31] architectures, to the task of segmenting B-mode breast US images from the BUSI and BUS datasets. The objective was to evaluate the performance of these three architectures (U-Net, SegNet, and DeepLab V3+) under similar training settings and dataset conditions, thereby comparing their efficacy relative to each other.

U-Net is a widely used architecture for biomedical image segmentation. It features a contracting path and an expansive path. The contracting path follows the typical architecture

of a convolutional network. In the expansive path, each step involves upsampling the feature map followed by a 2x2 convolution, often referred to as "up-convolution," which reduces the number of feature channels by half [25]. U-Net incorporates skip connections, which involve copying and cropping feature maps from the contracting path to corresponding locations in the expansive path. These skip connections enable the fusion of low-level and high-level features, thereby facilitating precise localization and segmentation of objects.

SegNet is another popular architecture that focuses on semantic pixel-wise segmentation using convolutional neural networks (CNNs). It consists of two parts: the encoder and the decoder. The encoder network consists of 13 convolutional layers, which correspond to the first 13 convolutional layers in the VGG16 network originally designed for object classification [27]. These layers are primarily responsible for analyzing object information. Each encoder layer has a corresponding decoder layer. The decoder network is tasked with reconstructing the segmented image based on the features extracted by the encoder.

DeepLab is a state-of-the-art DNN architecture targeting semantic image segmentation. DeepLab V3+ introduces several key improvements aimed at enhancing the accuracy and efficiency of segmentation tasks. One significant enhancement is the adoption of an encoder-decoder architecture, which combines the powerful encoder module from DeepLab V3 with a simple yet effective decoder module.

The DeepLab V3+ architecture has previously demonstrated success in conventional image segmentation tasks. This report explores the merits of DeepLab V3+ further by applying it to B-Mode US images and comparing its performance with other architectures such as U-Net and SegNet. By evaluating the effectiveness of DeepLab V3+ in this context, the project aims to provide insights into its applicability and potential benefits in US image analysis.

The rest of this report is organized as follows. Chapter 2 introduces BUSI and BUL datasets used in this project. The BUSI dataset is used for both training and testing of the DNN models under consideration. The BUL dataset is used exclusively for testing, serving as a source of truly unseen data samples. Furthermore, Chapter 2 describes the segmentation quality criteria used for evaluation purposes. Chapter 3 gives technical details about U-Net, SegNet, and DeepLab V3+ architectures and their specific configurations used in this project. Chapter 4 discusses our training process including dataset preparation, training settings, and learning curves. Chapter 5 presents the performance results for the DNNs models under consideration. Finally, Chapter 6 concludes this report and provides some suggestions for future work.

## Chapter 2. US Datasets and Evaluation Criteria

### 2.1 Breast US Images (BUSI) dataset

The training process in this project was based on the Breast US Images (BUSI) dataset [20], compiled and stored in a DICOM format at Baheya Hospital [20]. The consumed time used to collect and annotate the images was about one year. The total size of the dataset is 197 MB. It contains 133 normal scans, 487 scans with benign tumors, and 210 scans with malignant tumors. Each US image has a corresponding tumor mask label that has been reviewed by clinical radiologists.

The collected data samples include breast US images among 600 female patients aged between 25 and 75 years old, with an average image size of 500 x 500 pixels. A few example images are illustrated in Fig. 2.1.

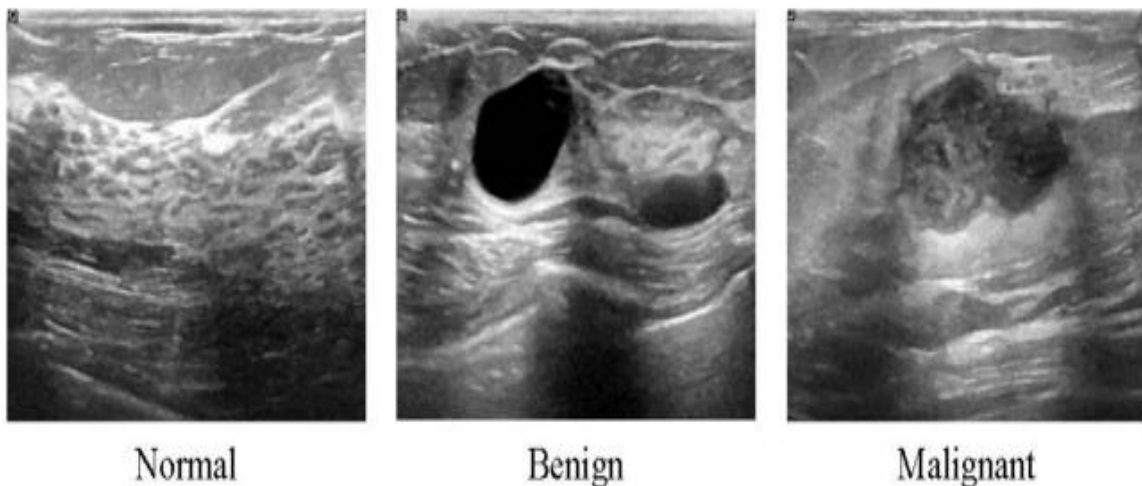


Figure 2.1 Breast US Images (BUSI) dataset samples [20].

For this project, only the images from benign and malignant categories are used. Among these images, 80% were used for training, 10% for validation, and the remaining 10% for testing.

### 2.2 Breast US Lesions (BUL) dataset

The BUL dataset was collected at the UDIAT Diagnostic Centre of the Parc Tauli' Corporation, Sabadell (Spain) with a Siemens ACUSON Sequoia C512 system 17L5 HD linear array transducer (8.5 MHz). The dataset consists of 163 images from different female patients with a mean image size of  $760 \times 570$  pixels, where each of the images presented one or more

lesions. Among these images, 53 were with cancerous masses and 110 with benign lesions. From the malignant images, 40 were with invasive ductal carcinomas, 4 with ductal carcinomas in situ, 2 with invasive lobular carcinomas, and 7 with other unspecified malignant lesions. From the benign images, 65 were with unspecified cysts, 39 with fibroadenomas, and 6 with some other types of benign lesions. All the lesions were delineated by experienced radiologists. A few example images are illustrated in Fig. 2.2 [21]. In this project, the BUL dataset was used for additional (unseen dataset) testing to evaluate each architecture objectively.

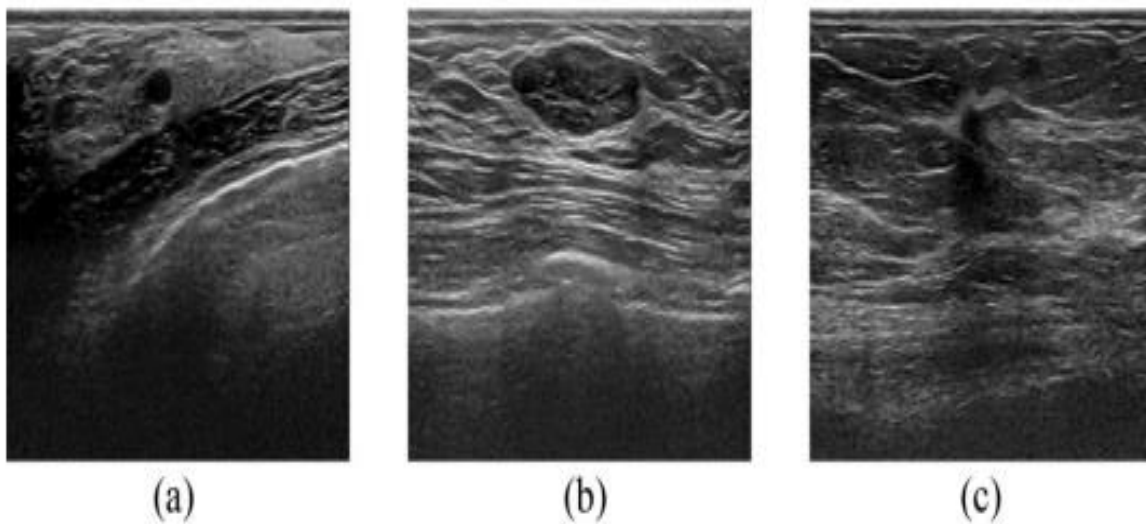


Figure 2.2 Example images with a cyst (a), a fibroadenoma lesion (b), and an invasive ductal carcinoma lesion (c) shown [21].

### 2.3 Segmentation quality criteria

Following the completion of the training phase, the performance of each model on the test datasets was assessed and compared using the following five performance metrics: accuracy, Dice score, IoU, and BF-score [18] [32]. These metrics are described in the rest of this section.

The performance metrics rely on statistics such as False Negative (FN), False Positive (FP), and True Positive (TP). They are illustrated in Fig. 2.3 and defined below.

- True Positive (TP): the number of pixels that the prediction model classified in class and were actually in class.
- False Positive (FP): the number of pixels that the prediction model misclassified in class, but they were not in class.
- False Negative (FN): the number of pixels the prediction model misclassified not in class, but they were in class.

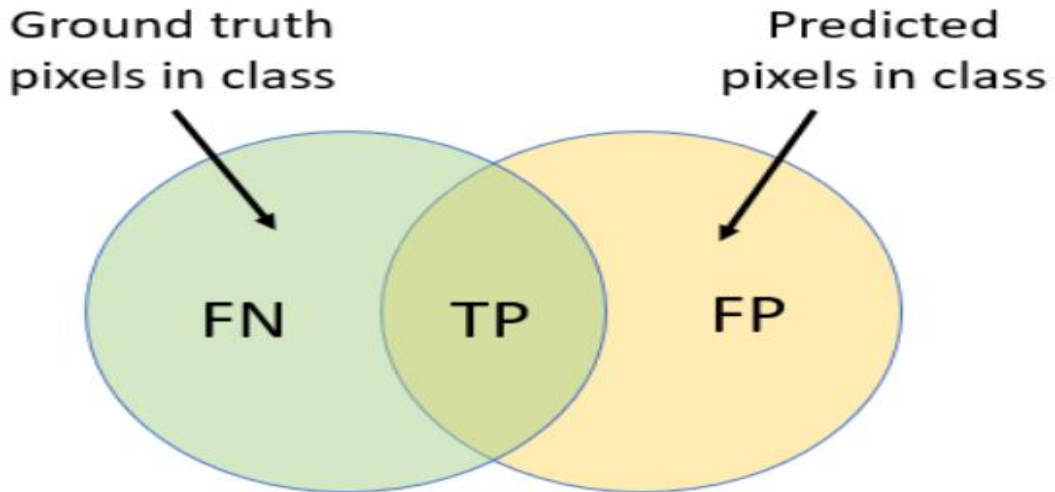


Figure 2.3 FN, FP, and TP [32].

### 2.3.1 Accuracy

Accuracy is the ratio of correctly classified pixels (TP) to the total number of pixels of the ground truth in that class. Generally, it is given by the following equation:

$$\text{Accuracy score} = \frac{\text{TP}}{(\text{TP} + \text{FN})} \quad (2.1)$$

Global accuracy is the ratio of correctly classified pixels to the total number of pixels for all classes (tumor and background). Mean accuracy is the average accuracy over all classes (tumor and background) in all images.

### 2.3.2 BF score

The BF score, also known as the Boundary F1 contour matching score, measures how close the predicted boundary of an object matches the ground truth boundary [32]. It is the ratio of the doubled product of precision and recall to the sum of them, as shown below:

$$\text{BF score} = \frac{2 * \text{precision} * \text{recall}}{(\text{recall} + \text{precision})} \quad (2.2)$$

In equation 2.2, precision is the ratio of the number of correctly classified pixels on the predicted boundary (TP) to the total number of pixels of the predicted boundary. It can be computed as follows:

$$\text{Precision} = \frac{\text{TP}}{(\text{TP} + \text{FP})} \quad (2.3)$$

In equation 2.2, recall is the ratio of the number of correctly classified pixels on the predicted boundary (TP) to the total number of pixels of the ground truth boundary. It can be computed as follows:

$$\text{Recall} = \frac{\text{TP}}{(\text{TP} + \text{FN})} \quad (2.4)$$

Precision is very important when false positives could result in unnecessary medical interventions or treatments. Recall becomes critical in scenarios where missing tumor regions could lead to delayed diagnosis or treatment, potentially impacting patient outcomes.

### 2.3.3 Intersection over union (IoU)

IoU, also known as the Jaccard similarity coefficient, measures the intersection (TP) over the union of the prediction and ground truth for each class and reports the average. IoU thus takes into account both the false alarms and the missed values for each class [18]. It is computed as shown below.

$$\text{IoU} = \frac{\text{TP}}{(\text{TP} + \text{FP} + \text{FN})} \quad (2.5)$$

There are two popular IoU versions: weighted IoU and mean IoU. Mean IoU is the average IoU score of all classes in all images. Weighted IoU is the average IoU of each class, weighted by the number of pixels in that class [32].

### 2.3.4 Dice score

Dice score measures the doubled intersection (TP) over the sum of the prediction and ground truth for each class and reports the average. It has a value in the range [0, 1]. A value of 1 means that the segmentation in the two images is a perfect match. It is given by the following equation:

$$\text{Dice score} = \frac{2 * \text{TP}}{(2 * \text{TP} + \text{FP} + \text{FN})} \quad (2.6)$$

In this project, we have used two versions of the Dice score: the average Dice score of background, and the average Dice score of tumors. Additionally, we also evaluate the median Dice score of tumors, which is simply the median value of all individual Dice score of tumor calculated across a set of samples.

## Chapter 3. Segmentation Networks

### 3.1 DeepLab V3+ architecture

DeepLab V3+ represents a significant advancement over its predecessor, DeepLab V3, in the field of semantic image segmentation. Developed as an extension of the DeepLab series by Google Research, DeepLab V3+ incorporates several key improvements aimed at enhancing the accuracy and efficiency of semantic segmentation.

One notable enhancement is the adoption of an encoder-decoder architecture, leveraging DeepLab V3 as a powerful encoder module combined with a simple yet effective decoder module that utilizes atrous convolutions with spatial pyramid pooling to recover object boundaries, as depicted in Fig. 3.1. This enhancement enables the model to discern fine details and capture complex structures in images, resulting in more precise segmentation results. Additionally, DeepLab V3+ integrates a feature pyramid network, facilitating the integration of information from multiple scales of the input image.

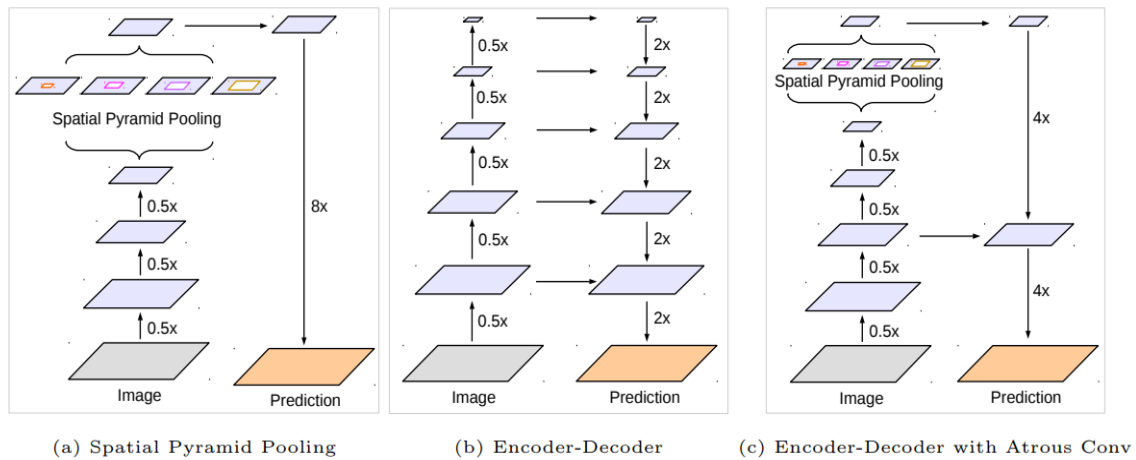


Figure 3.1 Encoder-Decoder with atrous convolution and spatial pyramid pooling [22].

Initially, the encoder features undergo bilinear upsampling, increasing their scale by a factor of 4. These upsampled features are then concatenated with the corresponding low-level features, which are obtained between the  $1 \times 1$  convolution and the encoder backbone. Following this concatenation, the features are refined by applying several  $3 \times 3$  convolutions. After this refinement process, another straightforward bilinear upsampling is executed, increasing the scale by a factor of 4, subsequently yielding the final output [22], as illustrated

in Fig. 3.2. The decoding structure gradually reconstructs spatial information to better capture object boundaries [26].

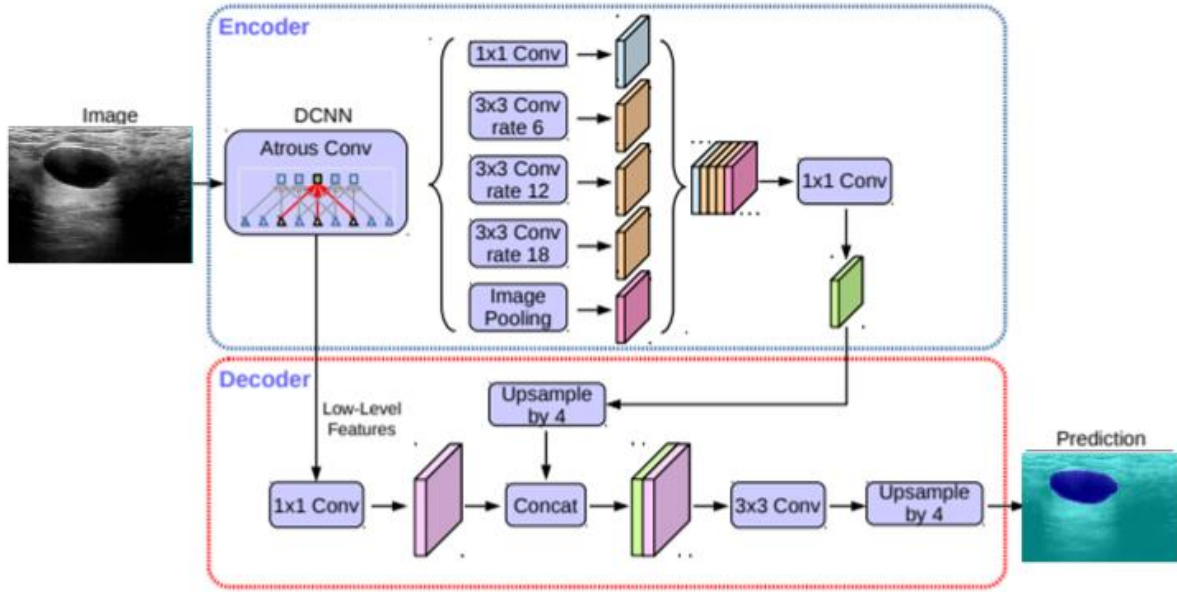


Figure 3.2 DeepLab V3+ architecture [22].

The purpose of aforementioned atrous convolution is to regulate the feature resolution computed by deep convolutional neural networks and to tune the filter's field-of-view to collect multi-scale information [22]. In the context of two-dimensional input, for every position  $i$  on the output feature map  $y$  and a convolution filter  $w$ , atrous convolution operates on the input feature map  $x$  in the following manner:

$$y[i] = \sum_k x[i + r \cdot k]w[k] \quad (3.1)$$

where the atrous rate  $r$  determines the stride with which we sample the input signal.

Atrous convolution can be combined with separable convolution, as depicted in Fig. 3.3. Depthwise separable convolution decomposes a standard convolution into two steps: depthwise convolution followed by pointwise convolution (i.e.,  $1 \times 1$  convolution). In this process, the depthwise convolution performs spatial convolutions independently for each input channel, while the pointwise convolution merges the outputs from the depthwise convolution [22].

## DeepLabv3+: Encoder-Decoder with Atrous Separable Convolution

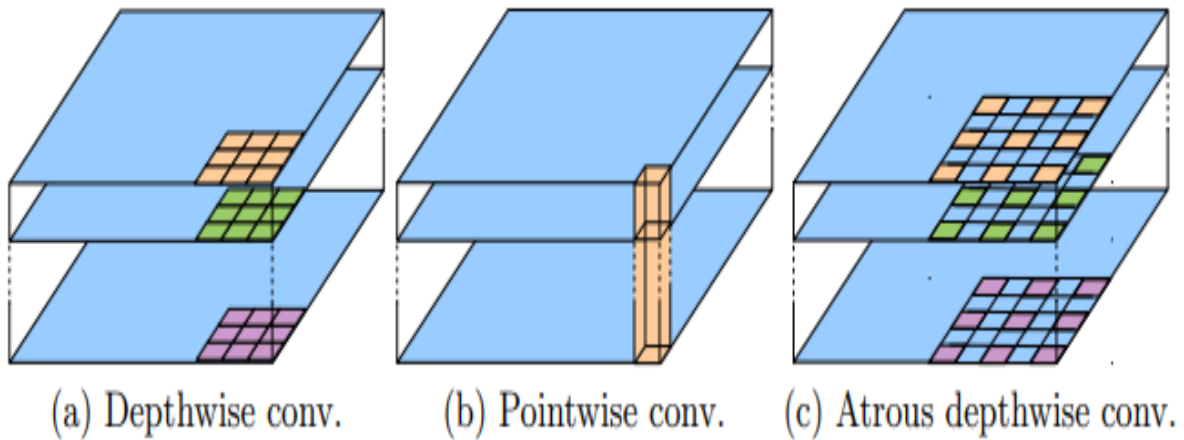


Figure 3.3 Atrous separable convolution illustrated (c) where atrous convolution is adopted for depthwise convolution (a) with rate = 2 [22].

In this project, we evaluate the DeepLab V3+ architecture with three different backbone networks used inside its encoder module: ResNet-18, ResNet-50, and Xception. Further discussion on these networks is provided in the following sections.

### 3.1.1 ResNet-18

Inspired by the VGG network [34], ResNet (Residual Network) basic architecture utilizes convolutional layers with  $3 \times 3$  filters. The architecture incorporates two key concepts for model optimization: maintaining the same number of filters for identical output feature maps across layers and doubling the number of filters when the output feature map's size is halved to preserve each layer's time complexity [30]. Fig. 3.4 illustrates the architecture of a ResNet model.

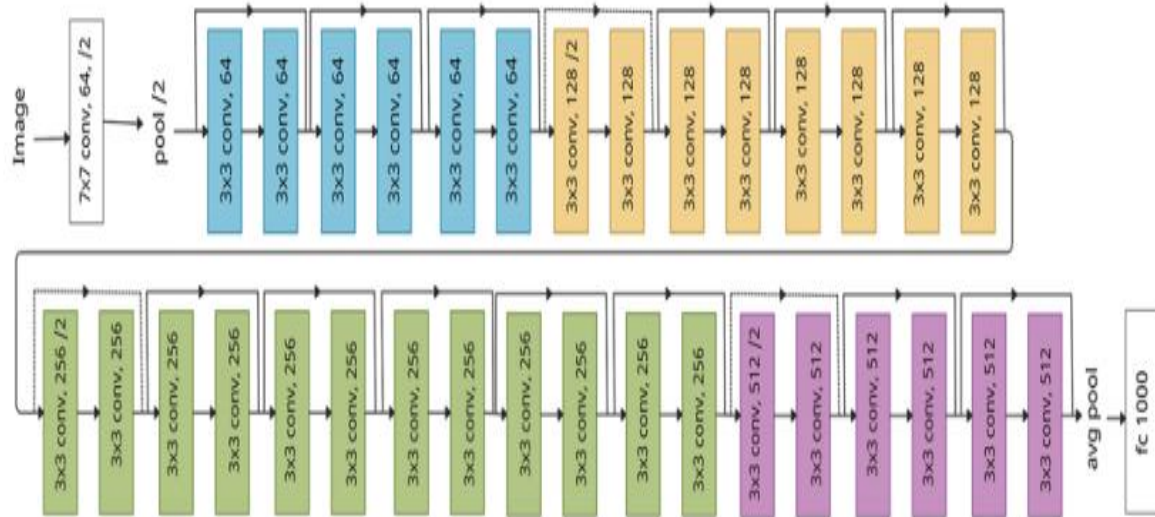


Figure 3.4 Architecture of typical ResNet model [30].

ResNet-18 belongs to the ResNet family [24]. This architecture tackles the difficulties associated with training extremely deep neural networks by introducing residual learning blocks. These blocks facilitate the learning of residual functions relative to the layer inputs, thereby addressing the vanishing gradient problem. A key innovation is the inclusion of skip connections within the residual blocks, which enhances the flow of gradient information throughout the network.

A typical residual block in ResNet-18 comprises a 3x3 convolutional layer with ReLU activation, followed by another 3x3 convolutional layer. The shortcut connection skips one convolutional layer at a time, and the input is added element-wise to the output from the residual block. Unlike larger ResNet variants, ResNet-18 does not utilize more complex bottleneck blocks. Fig. 3.5 illustrates the DeepLab V3+ architecture using the ResNet-18 backbone. This architecture consists of a total of 100 layers, including 29 convolutional layers, with 20.6 million learnable parameters, and it requires 62.5352 MB of memory.

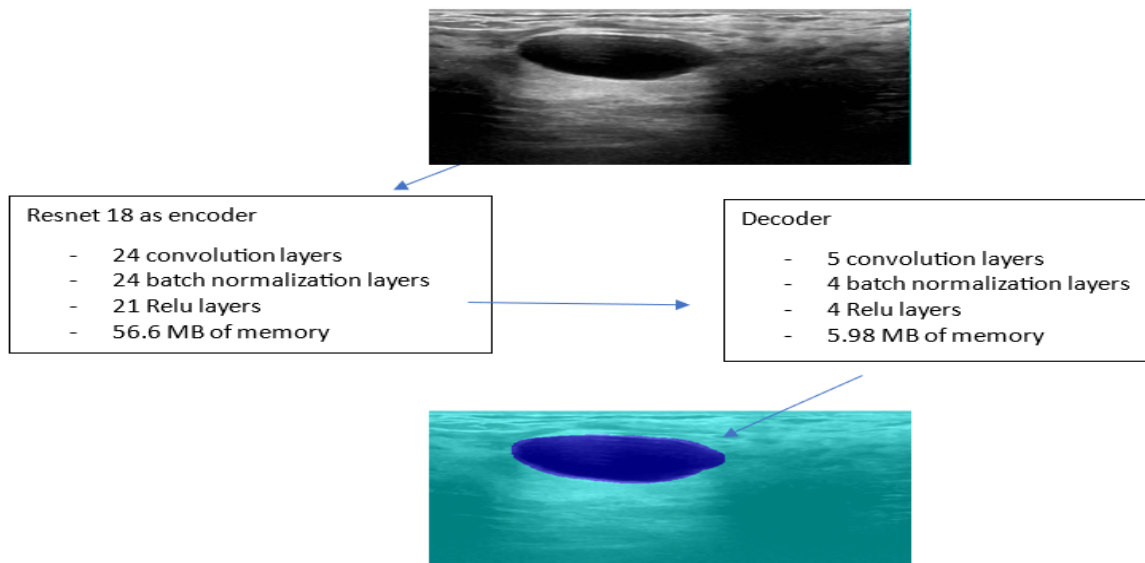


Figure 3.5 DeepLab V3+ with ResNet-18 bondnet as implemented in MATLAB.

### 3.1.2 ResNet-50

The architecture of ResNet-50 incorporates a more sophisticated structure known as the bottleneck residual block. Each bottleneck residual block consists of three convolutional layers: 1x1, 3x3, and 1x1, complemented by batch normalization and ReLU activations. Fig. 3.6 illustrates the difference between regular residual blocks and bottleneck residual blocks.

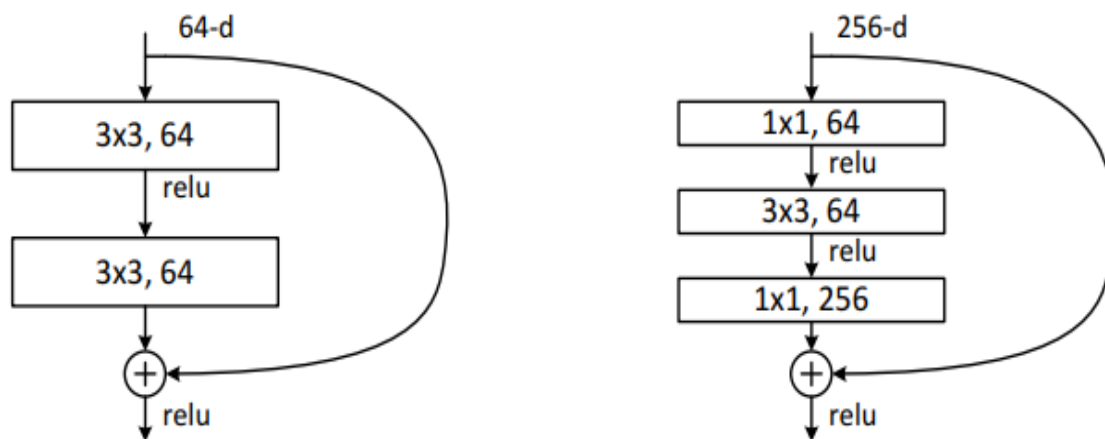


Figure 3.6 Left: regular residual block. Right: bottleneck residual block [24].

The bottleneck residual blocks are designed to capture richer feature representations, contributing to the model's overall efficacy. Similar to ResNet-18, ResNet-50 adopts global

average pooling to reduce spatial dimensions. The final layer is a fully connected layer with SoftMax activation, facilitating classification tasks. Skip connection, a key part of ResNet architectures, are present in ResNet-50 as well. These connections skip one or more layers, enhancing the flow of gradients during training and reducing the vanishing gradient problem [24].

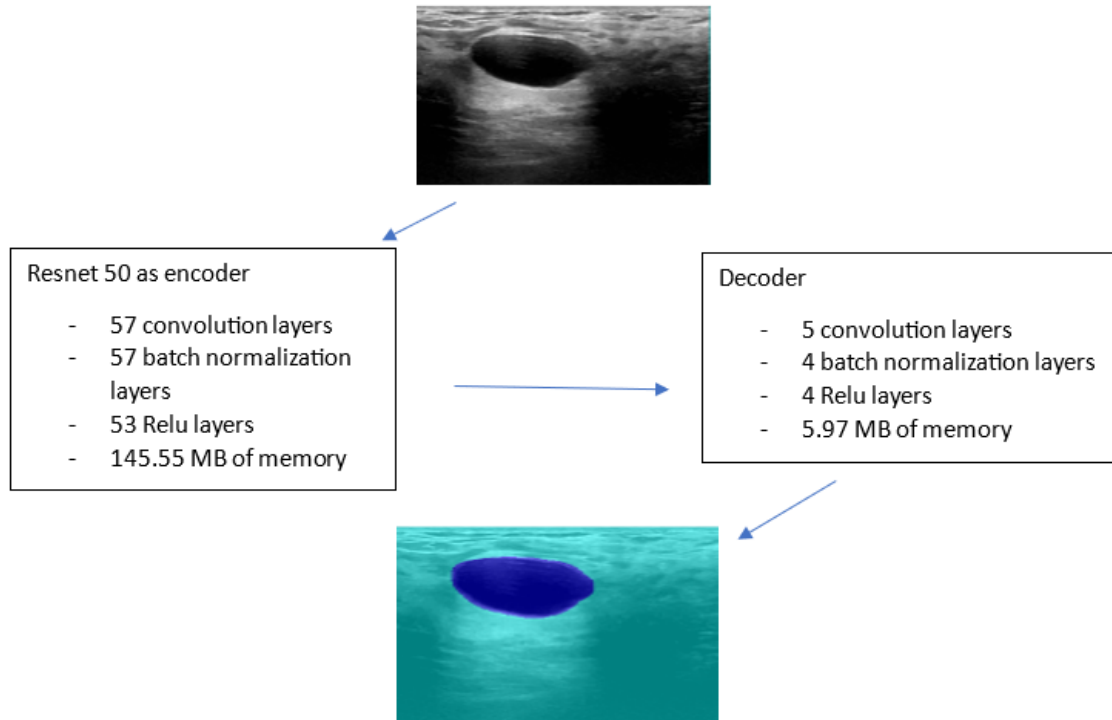


Figure 3.7 DeepLab V3+ with ResNet-50 bondnet as implemented in MATLAB.

Fig 3.7 depicts the DeepLab V3+ architecture utilizing the ResNet-50 backbone. This architecture consists of a total of 206 layers, including 62 convolutional layers, with 43.9 million learnable parameters. It requires approximately 151.5281 MB of memory.

### 3.1.3 Xception

Xception, short for 'Extreme Inception,' is a convolutional neural network architecture developed by François Chollet, the creator of the Keras deep learning library [29]. The concept behind the Inception module is to simplify and optimize data processing by breaking it down into a sequence of operations that independently examine cross-channel correlations and spatial correlations. A standard Inception module initially addresses cross-channel correlations through a series of 1x1 convolutions. These convolutions reduce the dimensionality of the input data into 3 or 4 separate spaces, which are smaller than the original input space. Subsequently,

the module utilizes regular 3x3 or 5x5 convolutions to capture all correlations within these smaller 3D spaces. This approach enables efficient feature extraction by separately handling cross-channel and spatial correlations, thereby enhancing the overall effectiveness of the convolutional neural network [29]. This is illustrated in Fig. 3.8.

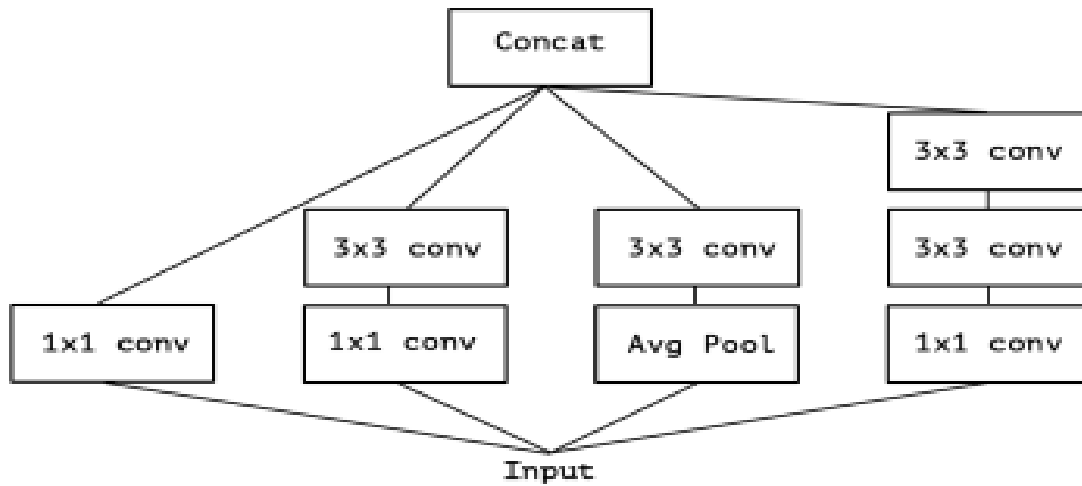


Figure 3.8 A canonical Inception module (Inception V3) [29].

Unlike Inception, an Xception module has two levels. The first level, containing a single 1x1 convolutional filter, splits the input into multiple segments (output channels). The second level contains multiple 3x3 convolutional filters (one per segment), whose outputs are subsequently concatenated. Another important aspect of Xception model is the use of separable convolutions [23–25]. The architecture of the Xception model is illustrated in Fig. 3.9.

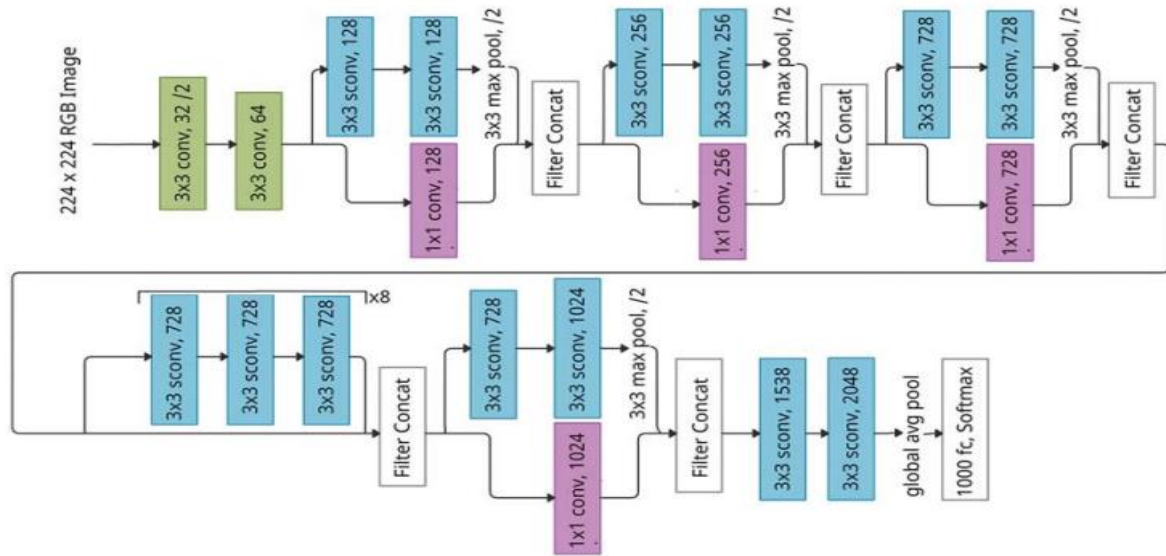


Figure 3.9 Architecture of Xception model [30].

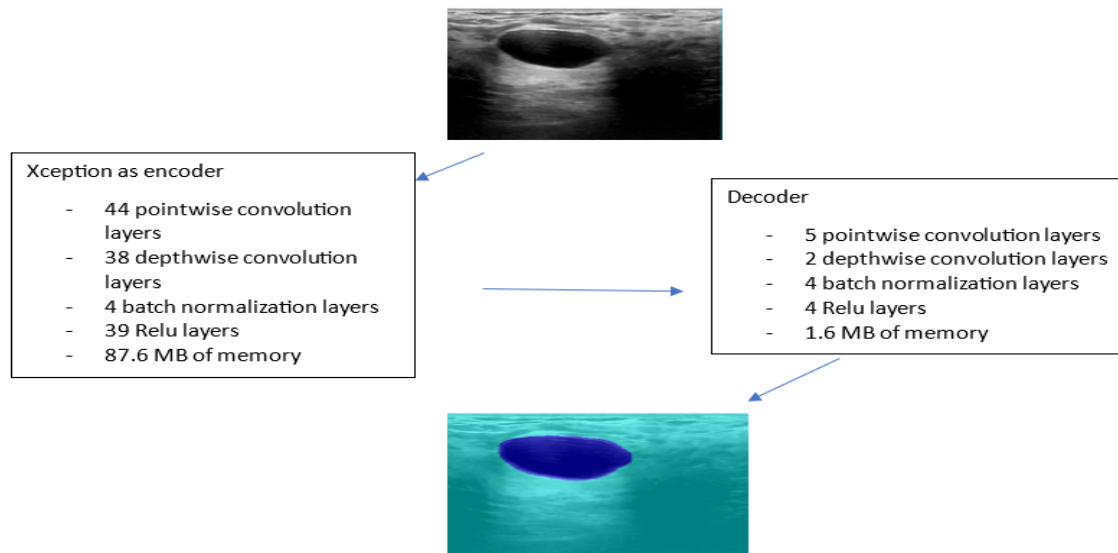


Figure 3.10 DeepLab V3+ with Xception bondnet as implemented in MATLAB.

Fig. 3.10 illustrates the DeepLab V3+ architecture utilizing the Xception backbone. This architecture comprises a total of 205 layers, including 49 pointwise convolution layers and 40 depthwise convolution layers. It has 27.6 million learnable parameters and requires approximately 89.2042 MB of memory.

### 3.2 U-Net architecture

U-Net is a convolutional neural network architecture designed specifically for image segmentation tasks, particularly in biomedical image analysis. It was introduced by Olaf Ronneberger, Philipp Fischer, and Thomas Brox in 2015. The name 'U-Net' is derived from the architecture's shape, which resembles the letter 'U'.

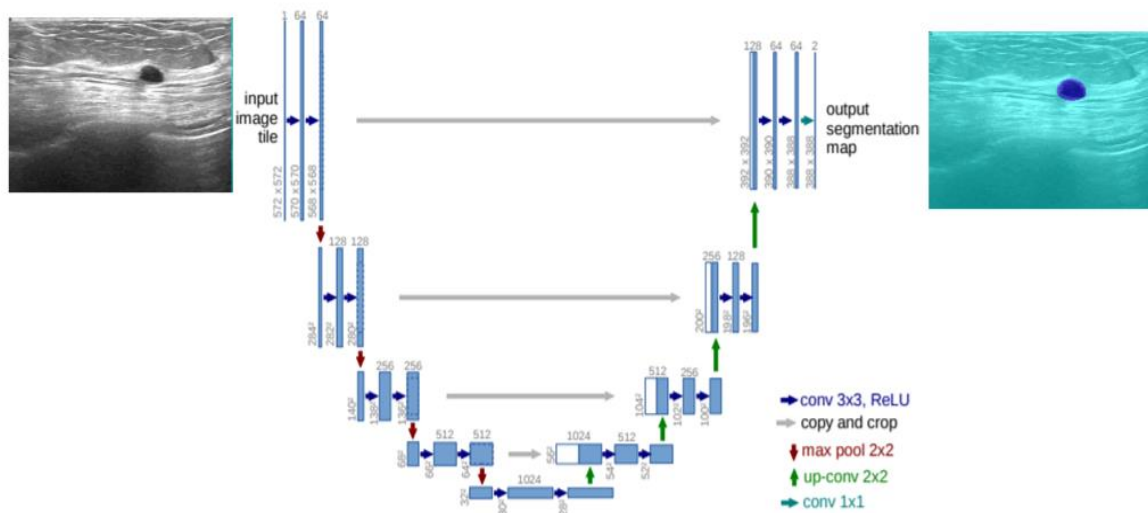


Figure 3.11 Architecture of U-Net model [30].

The network architecture is illustrated in Figure 3.11. It is comprised of two main paths: a contracting path (on the left side) and an expansive path (on the right side). The contracting path adheres to the typical design of a convolutional network, involving the repeated application of two 3x3 convolutions (without padding), each followed by a rectified linear unit (ReLU) and a 2x2 max pooling operation with a stride of 2 to downsample the feature maps. At each downsampling step, the number of feature channels is doubled [25].

In contrast, each step in the expansive path involves upsampling the feature map, followed by a 2x2 convolution (referred to as 'up-convolution') that halves the number of feature channels. The upsampled feature map is then concatenated with the correspondingly cropped feature map from the contracting path. Following the concatenation, two 3x3 convolutions, each followed by a ReLU, are applied. At the final layer, a 1x1 convolution is employed to transform each 64-component feature vector into the desired number of classes. This architectural design, incorporating both contraction and expansion paths, proves particularly effective for tasks such as semantic segmentation [25].

An important aspect of U-Net is its utilization of skip connections, which directly connect feature maps from the contracting path to corresponding feature maps in the expansive path. These skip connections play a crucial role in preserving spatial information and gradients during training, facilitating more effective segmentation. Additionally, the concatenation of up-sampled feature maps with those from the contracting path enables the leveraging of high-resolution information.

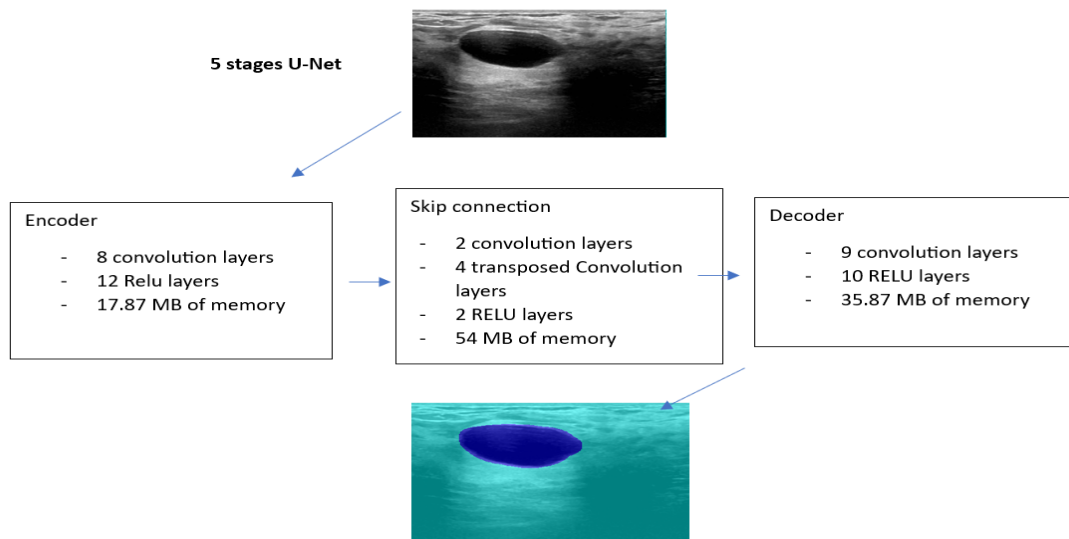


Figure 3.12 U-Net architecture as implemented in MATLAB.

Fig. 3.12 depicts the U-Net architecture implemented in MATLAB. It consists of 58 layers in total, has 28.3 million learnable parameters, and requires approximately 107.74 MB of memory.

### 3.3 SegNet architecture

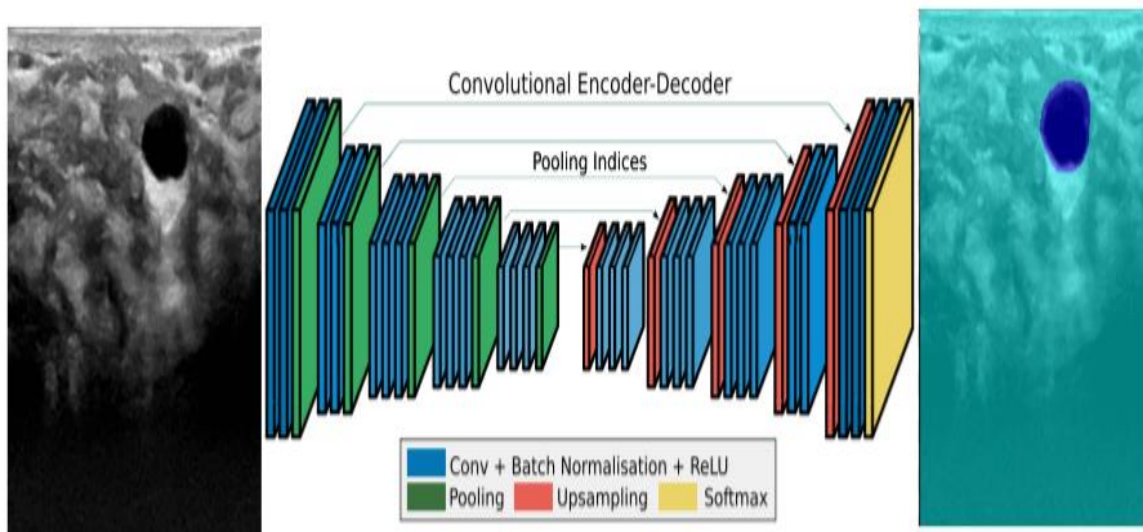


Figure 3.13 Architecture of SegNet model [31].

SegNet is a convolutional neural network designed for semantic segmentation tasks, focusing on pixel-wise classification. It was introduced in 2015 by Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla [31]. The encoder network consists of 13 convolutional layers, which correspond to the first 13 convolutional layers in the VGG16 network designed for object classification [27], primarily for analyzing object information. The fully connected layers are excluded to retain higher resolution feature maps at the deepest encoder output, significantly reducing the number of parameters in the SegNet encoder network. Each encoder layer has a corresponding decoder layer, resulting in a decoder network with 13 layers. The final decoder output is fed to a multi-class softmax classifier to produce class probabilities for each pixel independently. The network structure of SegNet is illustrated in Fig. 3.13 [31].

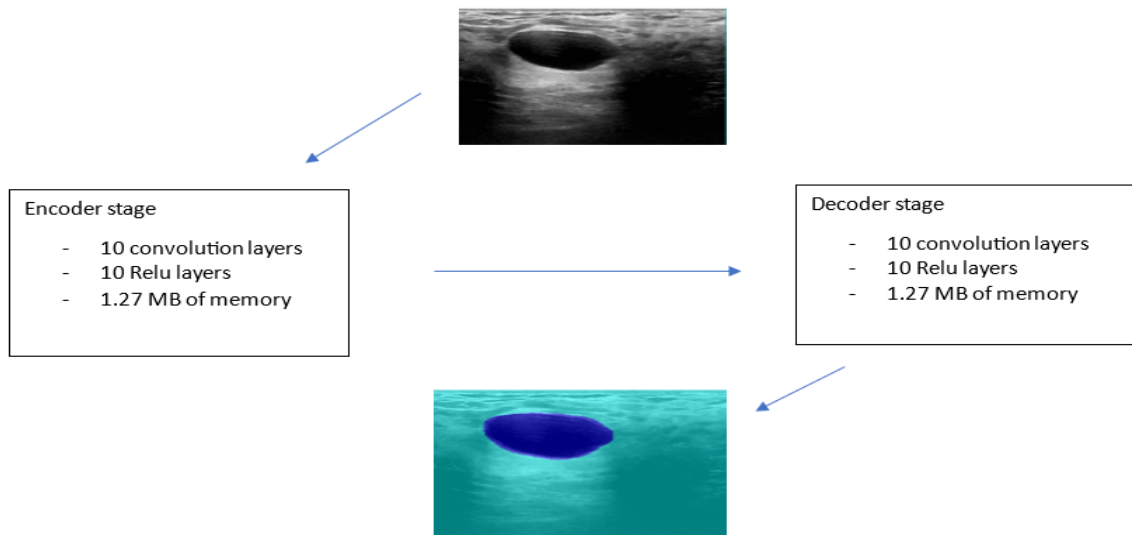


Figure 3.14 SegNet architecture as implemented in MATLAB.

Fig. 3.14 illustrates the SegNet architecture implemented in MATLAB. It has a total of 73 layers and contains 124.3 million learnable parameters. It requires approximately 2.54 MB of memory.

## Chapter 4. Network Training

### 4.1 Dataset preparation

#### 4.1.1 Data splitting, grayscale conversion, and resizing

**Data splitting.** The choice of how to split the available training data into different subsets is an important element in machine learning. The two most often used split ratios are 80:10:10 and 70:20:10. In this experiment, 80% of the BUSI dataset is used for training, 10% for validation, and the remaining 10% for testing. For the BUL dataset, 100% is used as "unseen data" for testing only.

**Grayscale conversion.** The BUSI dataset consists of B-mode ultrasound images having three channels. To simplify processing, the datasets are converted into grayscale, resulting in images with a single channel representing intensity. This conversion reduces computational complexity by simplifying the data, thus facilitating easier processing and analysis.

**Resizing.** Resizing the US images is one of the key steps in data preprocessing since various network architectures require different image sizes. To reduce processing time and adhere to the recommended image sizes for network architectures, all our images were adjusted to either 224x224 pixels or 299x299 pixels to fit each model's recommended input size.

#### 4.1.2 Data augmentation

Datasets are not always easily accessible, because it is challenging to collect a large number of high-quality samples. In order to increase the diversity of our dataset, the data augmentation technique was employed. There are many ways to perform augmentation such as horizontal or vertical flipping, rotation, shifting, zooming, cropping, blurring, adding noise, contrast and brightness adjustments, etc. In our case, each image in the training set was subjected to the following randomized methods:

**Brightness:** shifting the intensity of the image by a randomly selected value from the range  $[-0.5, 0.5]$ , as shown in Fig. 4.1.

**Contrast:** scaling the intensity of the image by a randomly selected value from the range  $[0.5, 1.5]$ , as shown in Fig. 4.2.

**Gamma:** correcting the intensity of the image with a randomly selected gamma value from the range  $[0.5, 1.5]$ , as shown in Fig. 4.3.

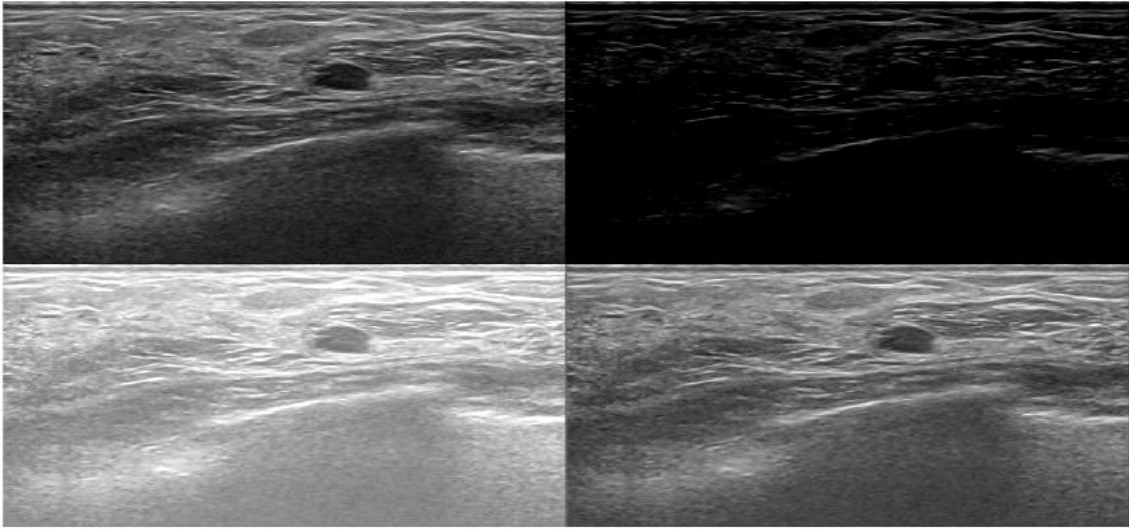


Figure 4.1 The left top corner is the original image. The others are augmented images using randomized brightness adjustment in the range  $[-0.5, 0.5]$ .

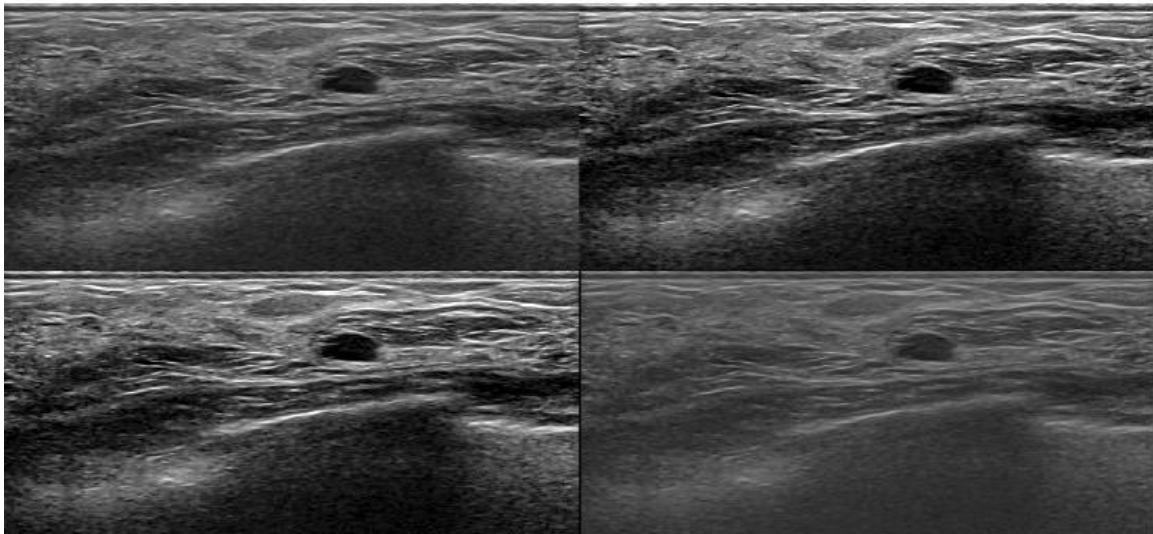


Figure 4.2 The left top corner is the original image. The others are augmented images using randomized contrast adjustment in the range  $[0.5, 1.5]$ .

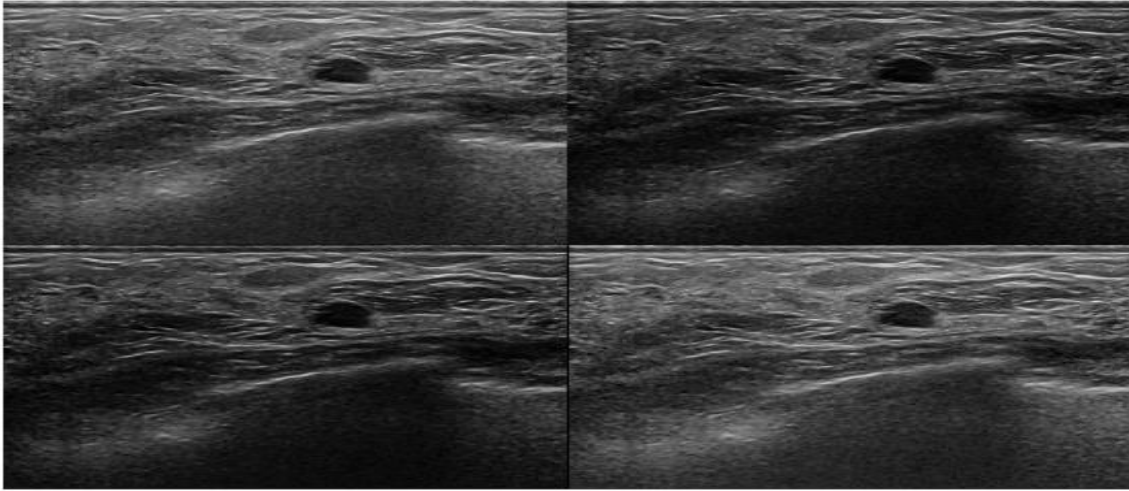


Figure 4.3 The left corner is original image. The others are augmented images using randomized gamma adjustment in the range [0.5, 1.5].

## 4.2. Training settings

### 4.2.1. Experimental parameters

MATLAB R2022b and related toolboxes, such as Deep Learning toolbox [35], Image Processing toolbox [37], Computer Vision toolbox [36], were used throughout this project to build and evaluate all the models. These were executed on a MATLAB cloud AWS version with a Nvidia A10G GPU processor having 24 GB of RAM. Table 4.1 provides the hardware and software specifications used to conduct our computational experiments.

| Item                 | Specification   |
|----------------------|---|
| Machine name         | Amazon AWS<br>cluster name: g5.2xlarge                                      |
| GPU                  | Nvidia A10G   |
| GPU memory           | 24 GB   |
| CPU model            | AMD EPYC 7R32 processor   |
| Cache size           | 19136 KB  |
| CPU MHz              | 2800.000  |
| Available RAM        | 32 GB   |
| Programming language | MATLAB R2022b   |
| MATLAB toolboxes     | Deep Learning toolbox, Image Processing toolbox,<br>Computer Vision toolbox |

Table 4.1 Experimental hardware and software specifications.

### 4.2.2. Learning option parameters

As stated earlier, the BUSI dataset split ratio was 80:10:10 for the three subsets: training, validation, and testing. The testing portion was put aside for model evaluation after the training and validation processes were done. The training option hyperparameters of all architectures are listed in Tables 4.2., 4.3, and 4.4 below.

| DeepLab V3+ architecture                                    |   |
|---|---|
| Hyperparameters   | Value   |
| Bondnets  | ResNet-18, ResNet-50, Xception                          |
| Input size  | 224 x 224 grayscale<br>299 x 299 grayscale for Xception |
| Optimizer   | Adam  |
| Execution environment                                       | Nvidia A10G GPU   |
| Learn rate schedule   | Piecewise   |
| Learn rate drop factor                                      | 0.5   |
| Learn rate drop period                                      | 100   |
| Initial learning rate ( $\alpha$ )                          | 1e-3  |
| Max number of epochs  | 300   |
| Minibatch size  | 16  |
| Verbose frequency   | 20  |
| Data shuffling  | Once  |
| Decay rate of gradient moving average ( $\beta_1$ )         | 0.9   |
| Decay rate of squared gradient moving average ( $\beta_2$ ) | 0.999   |
| Denominator offset ( $\epsilon$ )                           | 1e-8  |
| Factor for L <sub>2</sub> regularization ( $\lambda$ )      | 1e-4  |
| Statistics in batch normalization layers                    | Population statistics                                   |

Table 4.2 DeepLab V3+ training option hyperparameters.

| U-Net architecture  |                       |
|---|-----------------------|
| Hyperparameters   | Value                 |
| Encoder depth   | 5                     |
| Input size  | 224 x 224 grayscale   |
| Optimizer   | Adam                  |
| Execution environment                                       | Nvidia A10G GPU       |
| Learn rate schedule   | Constant              |
| Initial learning rate ( $\alpha$ )                          | 1e-3                  |
| Max number of epochs  | 100                   |
| Minibatch size  | 16                    |
| Verbose frequency   | 20                    |
| Data shuffling  | Once                  |
| Decay rate of gradient moving average ( $\beta_1$ )         | 0.9                   |
| Decay rate of squared gradient moving average ( $\beta_2$ ) | 0.999                 |
| Denominator offset ( $\epsilon$ )                           | 1e-8                  |
| Factor for $L_2$ regularization ( $\lambda$ )               | 1e-4                  |
| Statistics in batch normalization layers                    | Population statistics |

Table 4.3 U-Net training option hyperparameters.

| SegNet architecture   |                       |
|---|-----------------------|
| Hyperparameters   | Value                 |
| Encoder depth   | 5                     |
| Input size  | 224 x 224 grayscale   |
| Optimizer   | Adam                  |
| Execution environment                                       | Nvidia A10G GPU       |
| Learn rate schedule   | Constant              |
| Initial learning rate ( $\alpha$ )                          | 1e-3                  |
| Max number of epochs  | 100                   |
| Minibatch size  | 16                    |
| Verbose frequency   | 20                    |
| Data shuffling  | Once                  |
| Decay rate of gradient moving average ( $\beta_1$ )         | 0.9                   |
| Decay rate of squared gradient moving average ( $\beta_2$ ) | 0.999                 |
| Denominator offset ( $\epsilon$ )                           | 1e-8                  |
| Factor for $L_2$ regularization ( $\lambda$ )               | 1e-4                  |
| Statistics in batch normalization layers                    | Population statistics |

Table 4.4 SegNet training option hyperparameters.

The Adam (Adaptive moment estimation) optimizer, which is a popular optimization algorithm used in training machine learning models, includes the benefits of two other stochastic gradient descent extensions: AdaGrad and RMSProp [19]. AdaGrad is an adaptive gradient algorithm that keeps a per-parameter learning rate constant to improve performance when there are sparse gradients (e.g., computer vision problems and natural language processing). Root Mean Square Propagation (RMSProp) additionally maintains per-parameter learning rates based on the average of recent magnitudes of gradients for the weights (e.g., how quickly they are changing). This implies that the method performs well on both online and non-stationary data [19].

Adam uses a parameter update that is similar to RMSProp, but with an added momentum term. It keeps an element-wise moving average of both the parameter gradients and their squared values. Generally, it is given by the following equations:

$$\mathbf{m}_\ell = \beta_1 \mathbf{m}_{\ell-1} + (1 - \beta_1) \nabla \mathbf{E}(\boldsymbol{\theta}_\ell) \quad (4.1)$$

$$\mathbf{v}_\ell = \beta_2 \mathbf{v}_{\ell-1} + (1 - \beta_2) [\nabla \mathbf{E}(\boldsymbol{\theta}_\ell)]^2 \quad (4.2)$$

where  $\beta_1$  and  $\beta_2$  are the gradient decay (0.9) and squared gradient decay factors (0.999), respectively. Symbol  $\nabla \mathbf{E}(\boldsymbol{\theta}_\ell)$  is the gradient with respect to stochastic objective at timestep  $\ell$ . Symbol  $\mathbf{m}_\ell$  represents the quantity used to update biased first moment estimate. Symbol  $\mathbf{v}_\ell$  represents the quantity used to update biased second raw moment estimate. Adam relies on the moving averages to update the network parameters as follows:

$$\boldsymbol{\theta}_{\ell+1} = \boldsymbol{\theta}_\ell - \frac{\alpha \mathbf{m}_\ell}{\sqrt{\mathbf{v}_\ell + \varepsilon}} \quad (4.3)$$

where  $\varepsilon$  is denominator offset (1e-8) and  $\alpha$  is initial learning rate (1e-3).

The loss function is a critical component in machine learning algorithms, serving as a measure of how well the model performs on the given data. The primary objective of training is to minimize prediction errors. These prediction errors, or discrepancies between the model's predictions and the true values in the dataset, are quantified by the loss function to assess the model's accuracy. To prevent overfitting, adding a regularization term for the weights to the loss function  $\mathbf{E}(\boldsymbol{\theta})$  is a common trick:

$$\mathbf{E}_R(\boldsymbol{\theta}) = \mathbf{E}(\boldsymbol{\theta}) + \lambda \boldsymbol{\Omega}(\mathbf{w}) \quad (4.4)$$

where  $\mathbf{E}_R(\boldsymbol{\theta})$  is loss function after adding L2 regularization,  $\mathbf{w}$  is the weight vector,  $\lambda$  is the regularization factor (1e-4), and  $\boldsymbol{\Omega}(\mathbf{w})$  is the regularization function. In our case, we use

$$\boldsymbol{\Omega}(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \quad (4.5)$$

Learning Rate (LR) is a training parameter that influences how quickly the model weights are computed. A high LR might cause the model to converge too soon, whereas a low LR can result in more accurate weights (up to convergence) but requires more computing time. In our experiments, we used the learning rate of 0.001.

Validation portion of a dataset helps evaluate a trained machine learning model's performance before deployment. It acts as a gauge to assess how well the training performance

generalizes to unseen data and helps tune hyperparameters. The training and validation sets were shuffled once before training.

- The number of training epochs for DeepLab V3+ was set to 300. Initially, we attempted to train SegNet and U-Net for 300 epochs as well, but this resulted in overfitting to the training dataset, leading to poor performance on the test set during segmentation. To address this issue, the number of epochs for SegNet and U-Net was reduced to 100, resulting in improved performance.

- Given that DeepLab V3+ training lasted 300 epochs, a piecewise learning rate schedule was implemented to prevent overfitting. This strategy involves halving the learning rate after 100 epochs, aiding the model in identifying optimal solutions as the training progresses slower. In contrast, SegNet and U-Net, being trained only for 100 epochs, did not require a piecewise learning rate schedule. Instead, a constant learning rate was maintained.

### 4.3 Learning Curves

In this section we show the learning curves of the loss function and accuracy metric for the training and validation data. In all experiments the batch size was set to 16 (see Tables 4.2-4.6) and each epoch took 32 iterations.

#### 4.3.1 DeepLab V3+ with ResNet-18 bondnet

The training and validation results for the DeepLab V3+ model with the ResNet-18 bondnet are shown in Fig. 4.4. Notably, the training process involving 300 epochs took approximately 56 minutes, or 11 seconds per epoch (i.e., per 32 iterations).

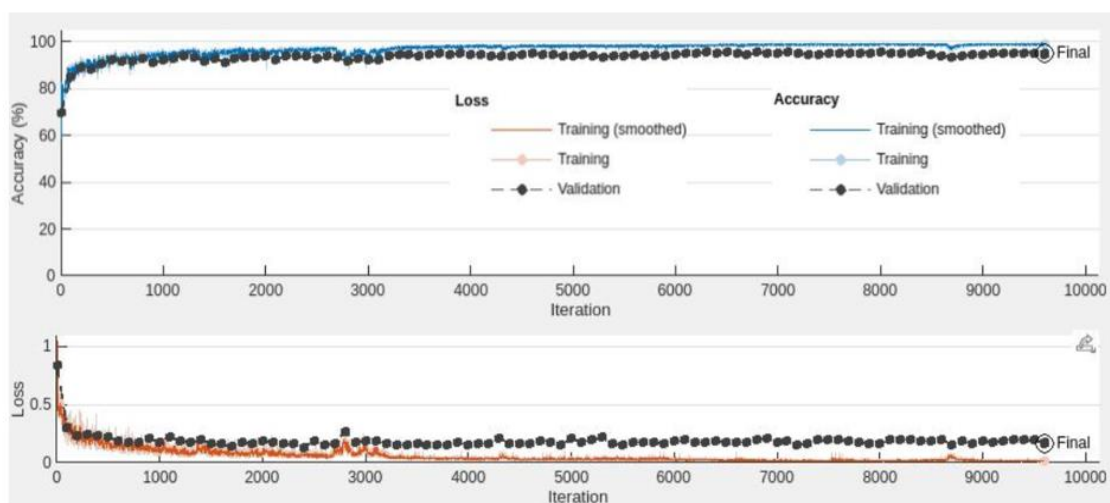


Figure 4.4 Learning curves for DeepLab V3+ based on ResNet-18

As one can see in Fig. 4.4, there is a noteworthy change for the training accuracy from approximately 69% at epoch 1 to 85.73% around epoch 4. It is followed by minor oscillations with a stable ascent between epoch 4 (iteration 128) and epoch 100 (iteration 3200), reaching a maximum of 98.84%. As for the validation accuracy, it shows an initial surge, increasing slowly afterwards to 95.23% at the end of epoch 300.

The values of the training data loss curve rapidly decreased from epoch 0 to epoch 5, and then reached 0.02 for train data (starting from 1.04 at the beginning). On the other hand, the validation data curve showed a noticeable drop from 0.8 to 0.17 at epoch 300.

### 4.3.2 DeepLab V3+ with ResNet-50 bondnet

The training and validation results for the DeepLab V3+ model with the ResNet-50 bondnet are shown in Fig. 4.5. Notably, the training process involving 300 epochs took approximately 89 minutes, or 18 seconds per epoch (i.e., per 32 iterations).

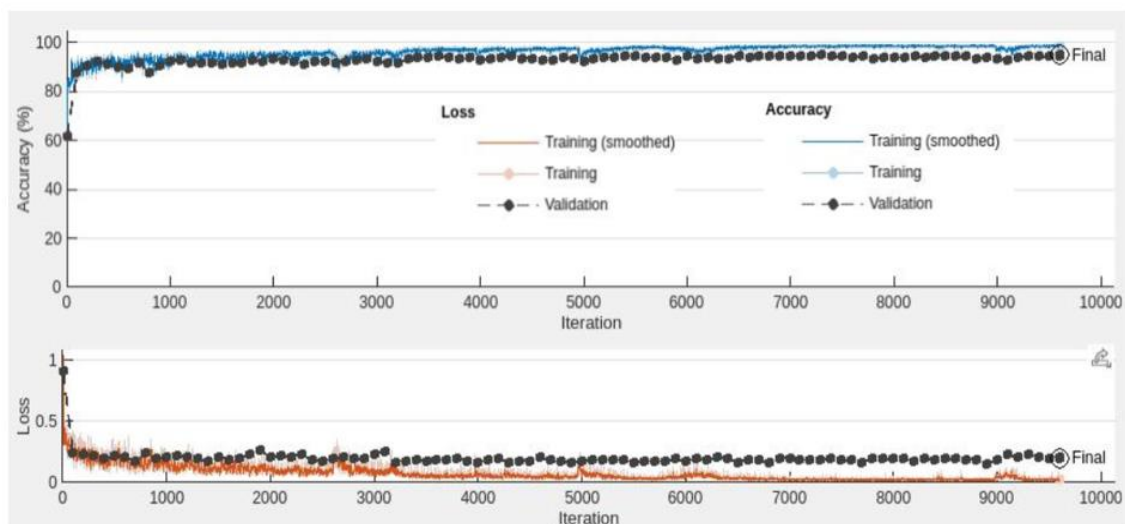


Figure 4.5 Learning curves for DeepLab V3+ based on ResNet-50 during training in MATLAB.

A noticeable increase in the training accuracy can be seen from approximately 61% at epoch 1 to 87.58% at epoch 5, reaching 99.36% at the end. For the validation data accuracy, a quick jump can be observed from epoch 0 to epoch 5, and then it oscillates while gradually increasing to 94.93% accuracy.

As for the loss curves, from epoch 0 to epoch 4, the loss is rapidly falling, and then begins to converge to 0.01 and 0.22 for the training and the validation data, respectively.

### 4.3.3 DeepLab V3+ with Xception bondnet

The training and validation results for the DeepLab V3+ model with the Xception bondnet are shown in Fig. 4.6. Notably, the training process involving 300 epochs took approximately 105 minutes, or 21 seconds per epoch (i.e., per 32 iterations).

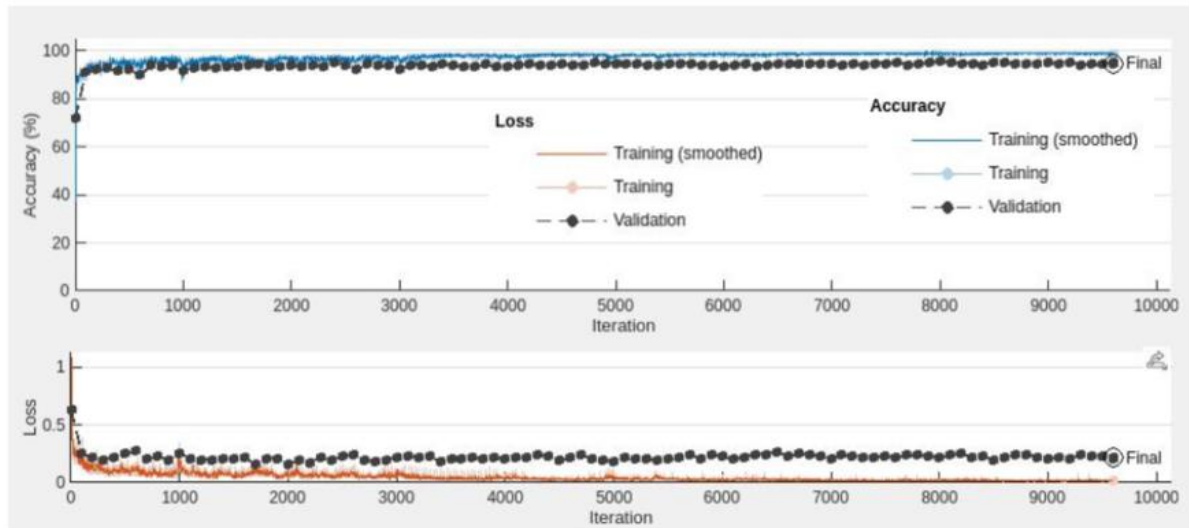


Figure 4.6 Learning curve for DeepLab V3+ based on Xception during training in MATLAB.

The resultant training data accuracy curve is growing quickly until it reaches 95% at the final epoch. After epoch 20, the accuracy steadily converges to 99% for the training data and 96% for validation data.

Meanwhile, the training loss is rapidly decreasing from epoch 0 to epoch 5, where it is equal to 0.2, and then begins to converge to 0.014. As for the validation loss curve, it falls to 0.2 and fluctuates around this value it until the last epoch.

### 4.3.4 SegNet architecture

The training and validation results for the SegNet model are shown in Fig. 4.7. Notably, the training process involving 100 epochs took approximately 26 minutes, or 16 seconds per epoch (i.e., per 32 iterations).

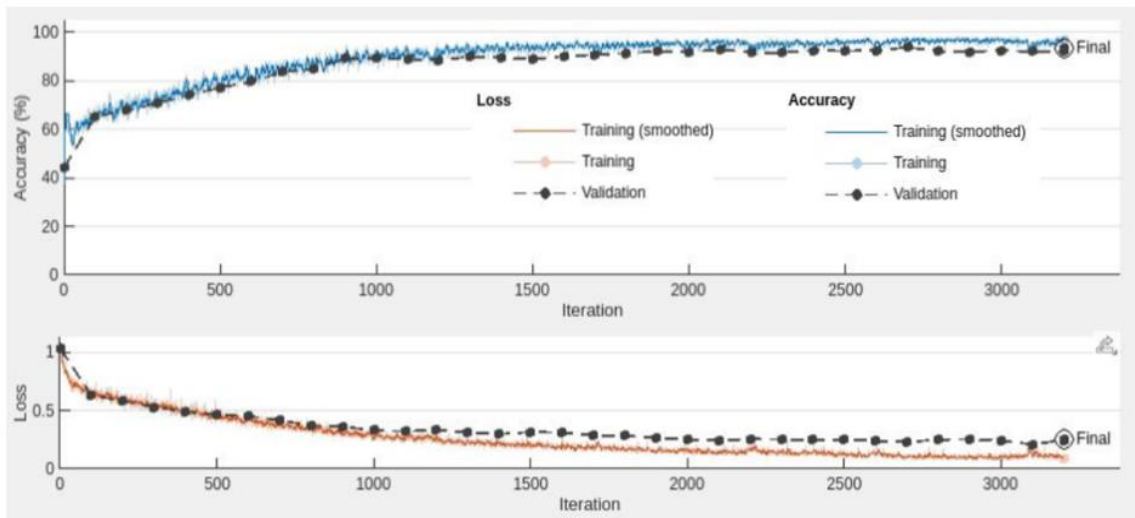


Figure 4.7 Learning curve for SegNet during training in MATLAB.

The resultant training data accuracy curve is growing quickly at first, from 42% at epoch 1 to about 89% at epoch 29 (iteration 928). After that, the accuracy slowly increases to 96.83% for the training data and 92.27% for the validation data.

Meanwhile, the training loss is rapidly decreasing from epoch 0 to epoch 19, where it is equal to 0.42, and then begins to converge to 0.0928. On the other hand, the validation loss curve initially falls to 0.4525 and then converges to 0.2478.

### 4.3.5 U-Net architecture

The training and validation results for the U-Net model are shown in Fig. 4.8. Notably, the training process involving 100 epochs took approximately 24 minutes, or 14 seconds per epoch (i.e., per 32 iterations).

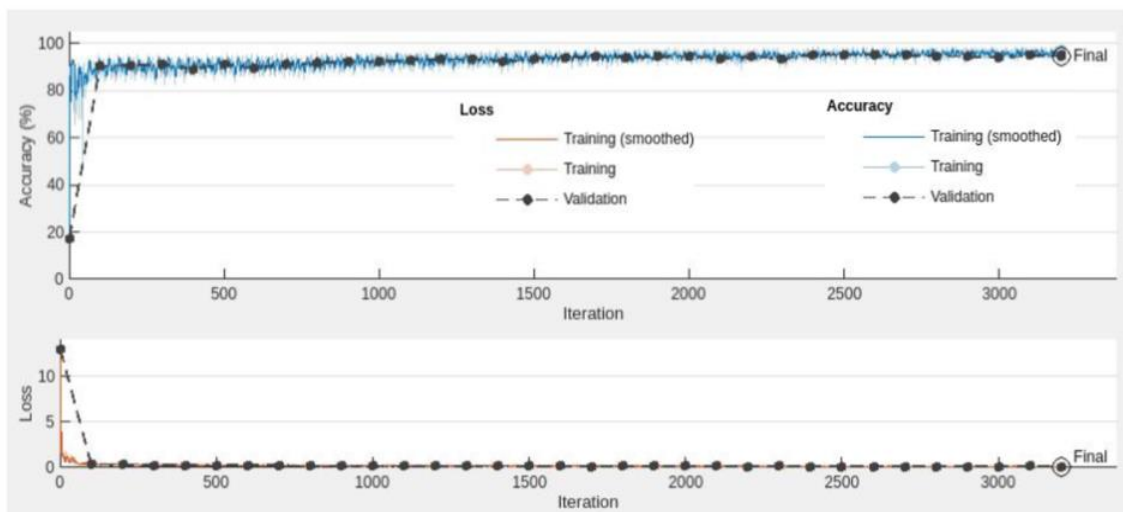


Figure 4.8 Learning curves for U-Net during training in MATLAB.

The resultant train data accuracy curve is growing quickly from 19% to 92% at epoch 5. Afterwards, the accuracy begins to converge to 98% for the training data and 96% for the validation data.

For the training data, the loss curve shows a sharp decline from epoch 0 to epoch 3, where the loss drops from 13 to 0.2. At epoch 100, the training loss is equal to 0.1. Similarly, the validation loss drops to 0.2 at the beginning, then decreased a loss of 0.15 at epoch 100.

## Chapter 5. Evaluation Results

In this chapter we present the inference-mode segmentation results obtained by the five networks models under consideration. Their performance on the training (80%) and validation (10%) portions of the BUSI dataset were discussed in Chapter 4. Here, we use the remaining 10% (65 images) of the BUSI dataset for testing purposes. To get additional evaluation results, we also utilize the entire BUL dataset (163 images) introduced in Chapter 2. The trained networks have not "seen" this dataset before, which is beneficial for assessing their generalization performance. Our evaluations are based on the segmentation quality metrics described in Chapter 2.

### 5.1 Performance of DeepLab V3+ with ResNet-18 bondnet

Table 5.1 and Fig. 5.1 summarize the segmentation performance of the DeepLab V3+ model based on ResNet-18. Notably, the model achieved high accuracy on both datasets, with global accuracy rates of approximately 94.7% for BUSI and 96.4% for BUL. Despite this high accuracy, there were also noticeable differences in other performance metrics between the two datasets.

Comparing the mean accuracy metric, ResNet-18 exhibited a slightly lower performance on the BUSI dataset (85.0%) compared to BUL (89.3%), implying potential challenges in accurately segmenting certain classes within the BUSI dataset. On the other hand, the mean IoU values were lower on the BUL dataset (72.1%), suggesting that segmentation performance was somewhat diminished on this unseen data compared to the BUSI dataset (74.3%). Similarly, the weighted IoU and mean BF score values were higher for the BUL dataset (94.3% and 64.6%) compared to BUSI (90.8% and 61.6%). However, the average and median Dice scores for tumors were notably higher for the BUSI dataset (66.1% and 74.0%) compared to BUL (52.9% and 62.6%). This discrepancy implies that the model's ability to achieve higher segmentation quality in terms of the tumor Dice scores is more pronounced on the BUSI dataset, indicating potential differences in tumor characteristics or dataset complexities. It should be mentioned that the average background Dice score yields the highest values among all metrics, reaching 96.9% and 98.0% for the BUSI and BUL datasets, respectively.

When comparing the testing accuracy of the ResNet-18 model for the BUSI dataset to its validation accuracy during training, the former remained comparably high (94.7% vs. 95.2%).

The consistency between testing and validation accuracies suggests that the model did not overfit BUSI training data.

| ResNet-18                               |                              |                    |
|---|------------------------------|--------------------|
|   | Test Portion of BUSI Dataset | Unseen BUL Dataset |
| <b>Global Accuracy</b>                  | 0.947                        | 0.964              |
| <b>Mean Accuracy</b>                    | 0.850                        | 0.893              |
| <b>Mean IoU</b>                         | 0.743                        | 0.721              |
| <b>Weighted IoU</b>                     | 0.908                        | 0.944              |
| <b>Mean BF score</b>                    | 0.620                        | 0.646              |
| <b>Average Dice score of Background</b> | 0.970                        | 0.980              |
| <b>Average Dice score of Tumor</b>      | 0.661                        | 0.530              |
| <b>Median Dice score of Tumor</b>       | 0.740                        | 0.630              |

Table 5.1 Evaluation results for DeepLab V3+ with ResNet-18 bondnet.

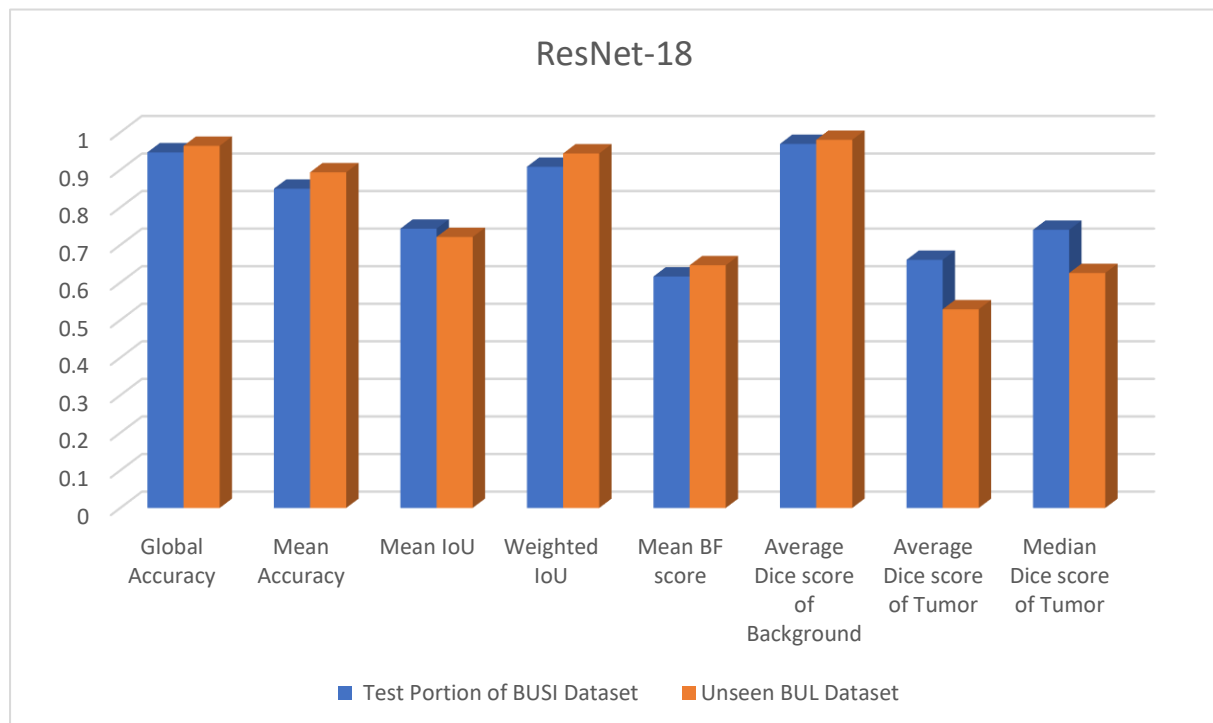


Figure 5.1 Evaluation results for DeepLab V3+ with ResNet-18 bondnet.

## 5.2 Performance of DeepLab V3+ with ResNet-50 bondnet

Table 5.2 and Fig. 5.2 summarize the segmentation performance of the DeepLab V3+ model based on ResNet-50. Notably, the model demonstrated high accuracy in pixel classification, with global accuracy reaching 94.9% and 96.6% for the BUSI and BUL datasets, respectively. In term of the mean accuracy metric, ResNet-50 achieved scores of 86.7% for the BUSI dataset and 88.8% for the BUL dataset. This suggests that the model performed slightly better when classifying pixels across different classes in the BUL dataset compared to the BUSI dataset. In terms of the mean and weighted IoU metrics, ResNet-50 achieved 75.7% and 91.3% on the BUSI dataset compared to 72.4% and 94.5% on the BUL dataset. Moreover, ResNet-50 demonstrated high average Dice score of background equal to 97.0% and 98.1% for the BUSI and BUL datasets, respectively. However, the average and median Dice score of tumor were significantly lower: 69.0% and 81.0% for the BUSI dataset compared to even worse 55.2% and 68.9% for the BUL dataset. In comparison to the ResNet-18 case, the mean BF score increased slightly to 65.2% and 67.5% for the BUSI and BUL datasets, respectively.

When comparing the testing accuracy of ResNet-50 for the BUSI dataset to its training and validation accuracies, we observed several key insights. The model achieved a very high accuracy of 99.4% on the training data. However, during validation, the model's accuracy decreased to 94.9%, which matched the testing accuracy on the unseen samples of the BUSI dataset. Such consistency between testing and validation accuracies indicates that the model did not overfit to the BUSI training data.

| ResNet-50                        |                              |                    |
|----------------------------------|------------------------------|--------------------|
|                                  | Test Portion of BUSI Dataset | Unseen BUL Dataset |
| Global Accuracy                  | 0.950                        | 0.966              |
| Mean Accuracy                    | 0.867                        | 0.888              |
| Mean IoU                         | 0.757                        | 0.724              |
| Weighted IoU                     | 0.913                        | 0.945              |
| Mean BF score                    | 0.652                        | 0.675              |
| Average Dice score of Background | 0.970                        | 0.981              |
| Average Dice score of Tumor      | 0.690                        | 0.552              |
| Median Dice score of Tumor       | 0.810                        | 0.689              |

Table 5.2 Evaluation results for DeepLab V3+ with ResNet-50 bondnet.



Figure 5.2 Evaluation results for DeepLab V3+ with ResNet-50 bondnet.

### 5.3 DeepLab V3+ with Xception bondnet

Table 5.3 and Fig. 5.3 summarize the segmentation performance of the DeepLab V3+ model based on Xception. Across various metrics, the model demonstrated similar performance to the ResNet-50 case on both datasets. The global accuracy of Xception was 95.4% for the BUSI dataset and 96.0% for the BUL dataset. The mean accuracy was 88.3% for the BUSI dataset and 90.6% for the BUL dataset, indicating a slightly better performance in pixel classification for the latter. In terms of the mean and weighted IoU metrics, Xception achieved 77.6% and 92.0% on the BUSI dataset compared to 70.9% and 93.8% on the BUL dataset. Furthermore, Xception exhibited high average Dice score of background equal to 97.3% and 97.8% for the BUSI and BUL datasets, respectively. However, the average and median Dice score of tumor were only 69.8% and 82.5% for the BUSI dataset compared to 56.0% and 70.5% for the BUL dataset. The mean BF score was 65.2% and 67.0% for the BUSI and BUL datasets, respectively, which is very close to the ResNet-50 case.

When comparing the testing accuracy of Xception for the BUSI dataset to its validation accuracy during training, the former remained comparably high (96% vs. 95.4%). This consistency suggests that the model did not overfit to the BUSI training data.

| <b>Xception</b>                         |   |                               |
|---|---|-------------------------------|
|   | <b>Test Portion of<br/>BUSI Dataset</b> | <b>Unseen<br/>BUL Dataset</b> |
| <b>Global Accuracy</b>                  | 0.954                                   | 0.968                         |
| <b>Mean Accuracy</b>                    | 0.883                                   | 0.906                         |
| <b>Mean IoU</b>                         | 0.776                                   | 0.709                         |
| <b>Weighted IoU</b>                     | 0.920                                   | 0.938                         |
| <b>Mean BF score</b>                    | 0.652                                   | 0.670                         |
| <b>Average Dice score of Background</b> | 0.974                                   | 0.978                         |
| <b>Average Dice score of Tumor</b>      | 0.698                                   | 0.560                         |
| <b>Median Dice score of Tumor</b>       | 0.825                                   | 0.705                         |

Table 5.3 Evaluation results for DeepLab V3+ with Xception bondnet.

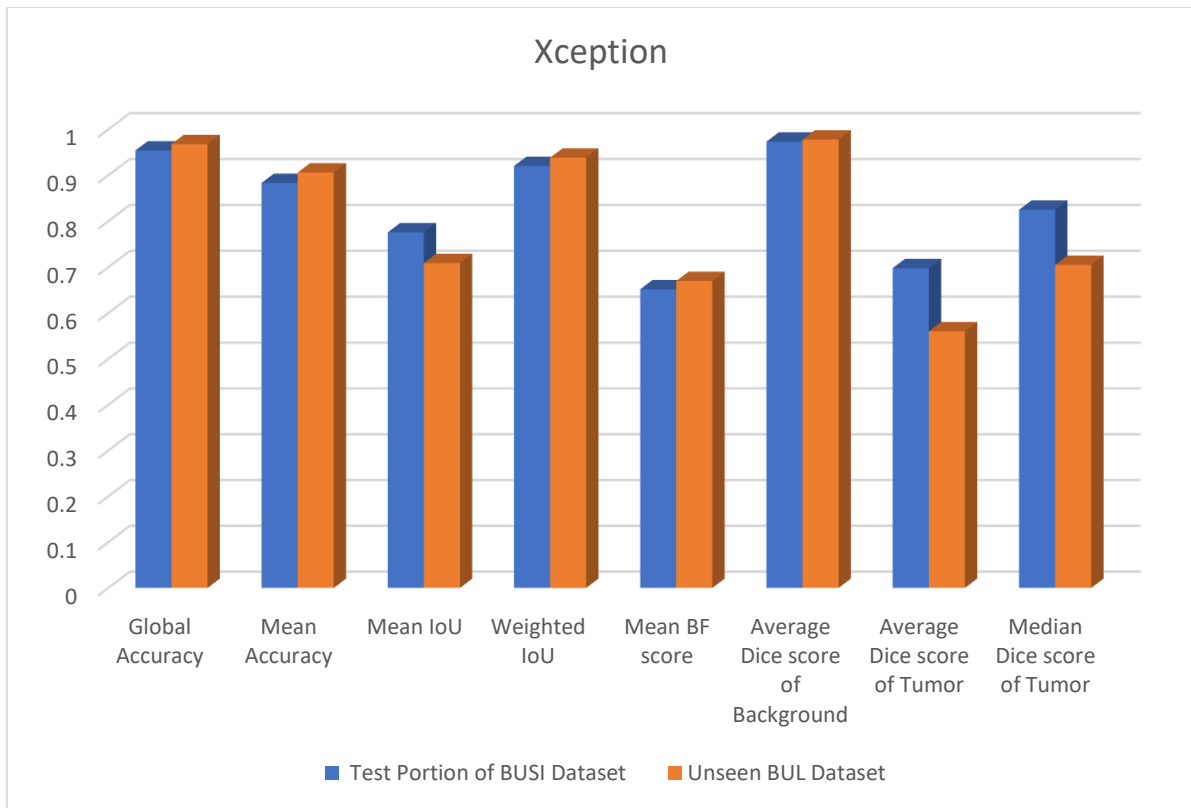


Figure 5.3 Evaluation results for DeepLab V3+ with Xception bondnet.

#### 5.4 Performance of SegNet

Table 5.4 and Fig. 5.4 summarize the segmentation performance of the SegNet model. The global accuracy of SegNet was 94.3% for the BUSI dataset and 95.6% for the BUL dataset. The mean accuracy was 85.2% and 89.9% for the BUSI and BUL datasets, respectively, suggesting a slightly better performance in pixel classification on the BUL dataset. In terms of the mean and weighted IoU metrics, SegNet achieved 73.6% and 90.4% on the BUSI dataset compared to 69.5% and 93.4% on the BUL dataset. As with the other models, SegNet yielded high average Dice score of background equal to 96.7% for the BUSI dataset and 97.5% for BUL. However, the average and median Dice score of tumor were 66.3% and 74.7% for the BUSI dataset, which is close to the ResNet-18 performance. The respective scores for the BUL dataset were 50.3% and 60.2%, which is worst among all five models. The mean BF score was 59.5% and 60.7% for the BUSI and BUL datasets, respectively, which also corresponds to the worst performance among all five models.

Recall that the model's validation accuracy for the BUSI dataset was 92.3%. Its testing accuracy for this dataset was slightly higher, reaching 94.3%, which suggests that SegNet did not overfit to the BUSI training data.

| SegNet                                  |                              |                    |
|---|------------------------------|--------------------|
|   | Test Portion of BUSI Dataset | Unseen BUL Dataset |
| <b>Global Accuracy</b>                  | 0.943                        | 0.956              |
| <b>Mean Accuracy</b>                    | 0.852                        | 0.899              |
| <b>Mean IoU</b>                         | 0.736                        | 0.695              |
| <b>Weighted IoU</b>                     | 0.904                        | 0.934              |
| <b>Mean BF score</b>                    | 0.595                        | 0.607              |
| <b>Average Dice score of Background</b> | 0.967                        | 0.975              |
| <b>Average Dice score of Tumor</b>      | 0.663                        | 0.503              |
| <b>Median Dice score of Tumor</b>       | 0.747                        | 0.601              |

Table 5.4 Evaluation results for SegNet.

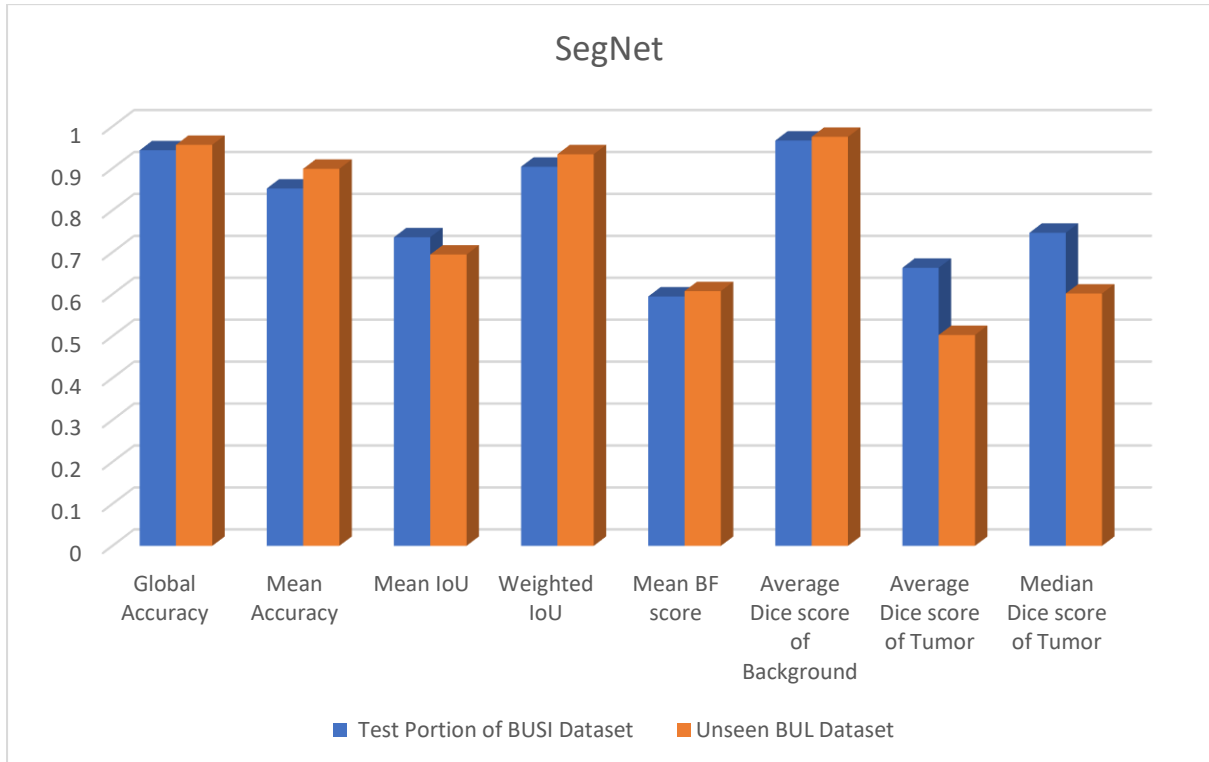


Figure 5.4 Evaluation results for SegNet.

## 5.5 Performance of U-Net

Table 5.5 and Fig. 5.5 summarize the segmentation performance of the U-Net model. The global and mean accuracy scores for the BUSI dataset were 95.6% and 85.4%, compared to 97.5% and 85.1% for the BUL dataset. Regarding the IoU metrics, U-Net achieved a mean IoU of 77.2% on the BUSI dataset and 73.2% on the BUL dataset. Their respective weighted IoU scores were 92.1% and 95.6%, which is best among all five models. U-Net exhibited high average Dice scores for background pixels, yielding 97.5% for the BUSI dataset and 98.5% for the BUL dataset. However, for tumor pixels, the model achieved much lower average and median Dice scores: 68.3% and 83.6% for the BUSI dataset compared to 52.8% and 66.9% for the BUL dataset, which is in line with the other models. The mean BF score was 65.2% and 67.5% for the BUSI and BUL datasets, respectively, which matches the ResNet-50 case.

As with the other models, the validation and testing accuracy scores for the BUSI dataset were numerically close, equal to 95.8% and 95.6%, respectively. This suggests that U-Net did not overfit to the BUSI training data.

| U-Net                                   |                              |                    |
|---|------------------------------|--------------------|
|   | Test Portion of BUSI Dataset | Unseen BUL Dataset |
| <b>Global Accuracy</b>                  | 0.956                        | 0.975              |
| <b>Mean Accuracy</b>                    | 0.854                        | 0.851              |
| <b>Mean IoU</b>                         | 0.772                        | 0.732              |
| <b>Weighted IoU</b>                     | 0.921                        | 0.956              |
| <b>Mean BF score</b>                    | 0.652                        | 0.675              |
| <b>Average Dice score of Background</b> | 0.975                        | 0.985              |
| <b>Average Dice score of Tumor</b>      | 0.682                        | 0.528              |
| <b>Median Dice score of Tumor</b>       | 0.836                        | 0.669              |

Table 5.5 Evaluation results for U-Net.

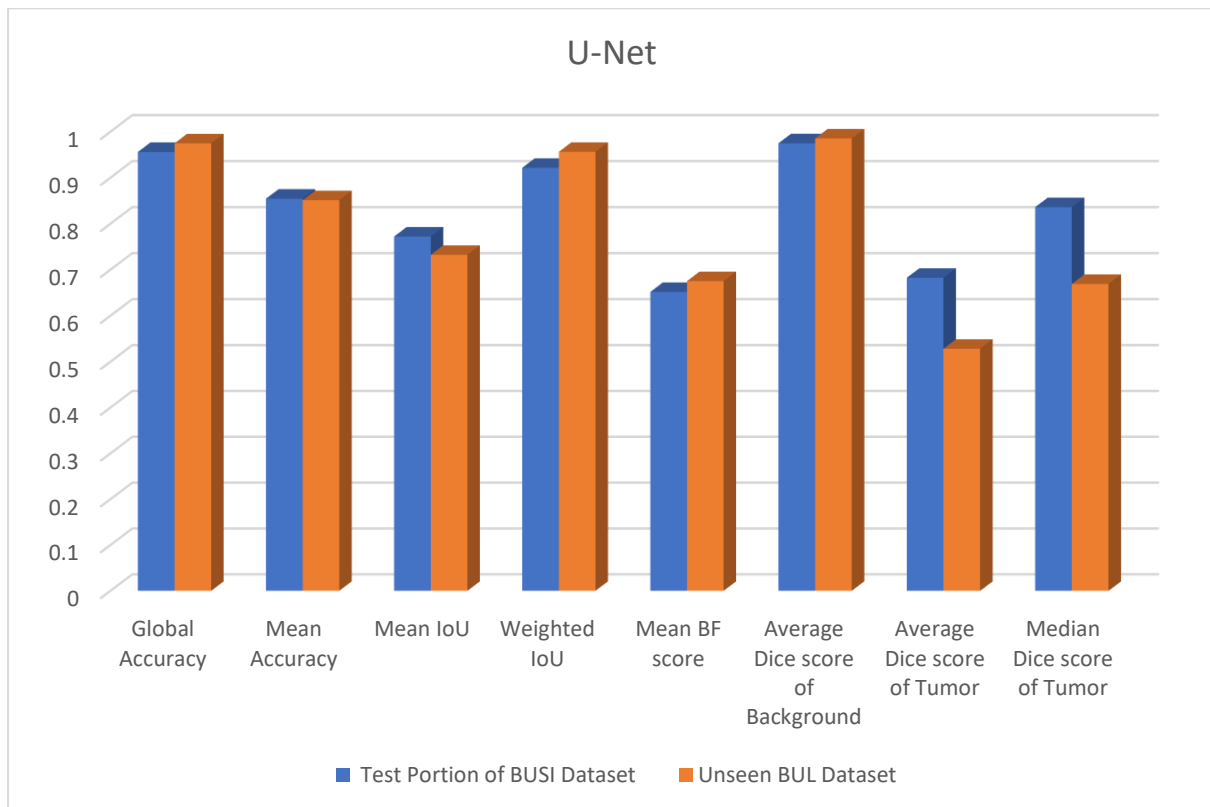


Figure 5.5 Evaluation results for U-Net.

## 5.6 Comparative summary

This section provides a comparative summary of the five networks under consideration based on the five criteria: Median Dice score of Tumor, Average Dice score of Tumor, Mean Accuracy, Mean IoU, and Mean BF score. The results are shown in Fig. 5.6 for the unseen BUL dataset and Fig. 5.7 for the test portion of the BUSI dataset.

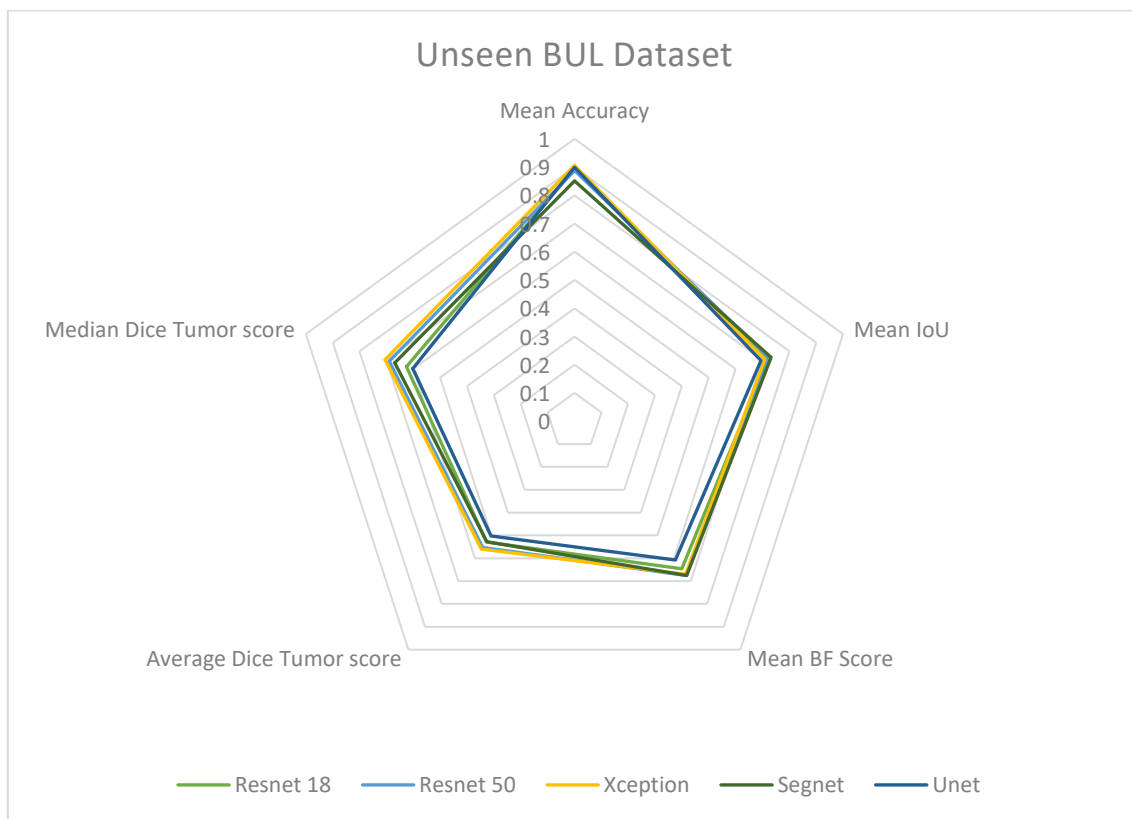


Figure 5.6 Network comparison results: Unseen BUL dataset.

Across the evaluated models applied to the BUL dataset—ResNet 18, ResNet 50, Xception, SegNet, and U-Net—distinct performances were observed across various segmentation metrics. U-Net demonstrated superior performance in terms of the mean IoU and mean BF score equal to 73.2% and 67.5% (tied with ResNet-50), respectively. Xception exhibited the highest mean accuracy with a score of 90.6%, showcasing its proficiency in overall pixel classification. It also achieved the highest average Dice score of tumor (56.0%) and the highest median Dice score of tumor (70.5%). These observations indicate that U-Net and Xception-based DeepLab V3+ are the two main competitors, whose selection depends on the targeted segmentation quality metrics under consideration (e.g., mean IoU vs. mean accuracy).

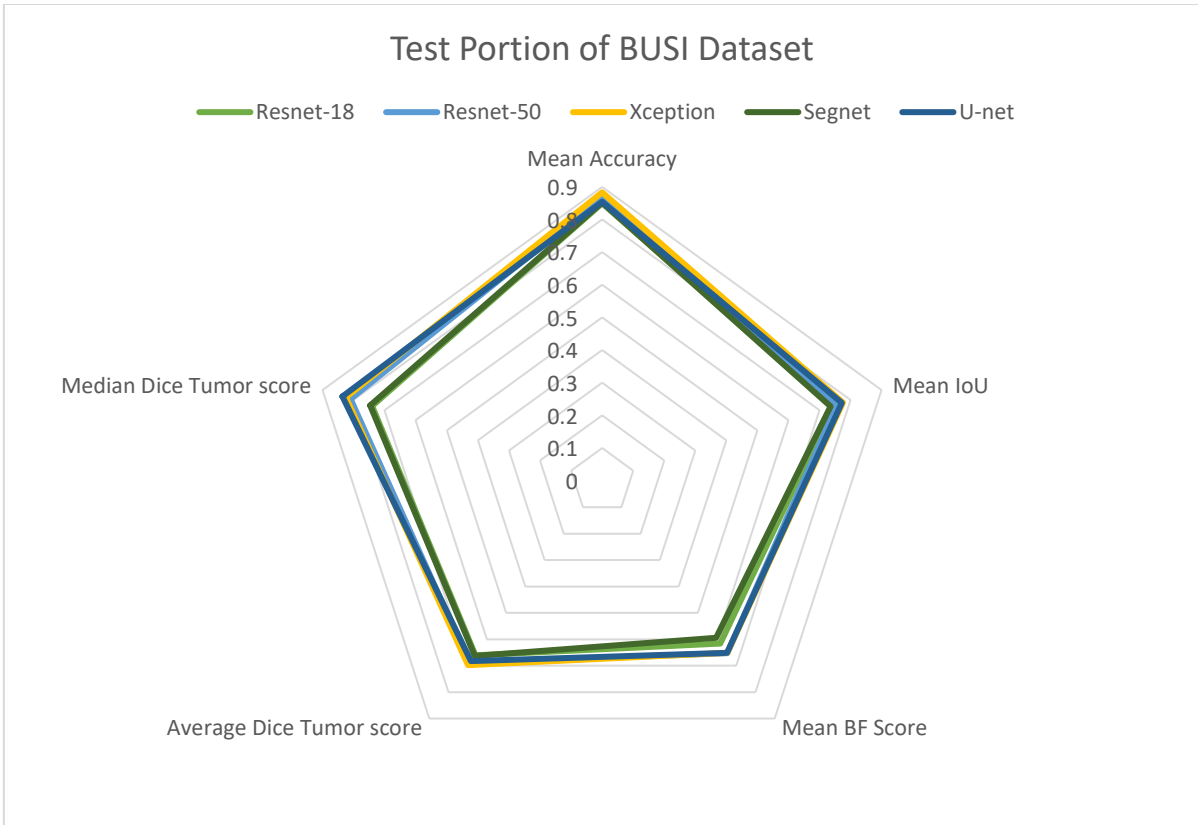


Figure 5.7 Network comparison results: Test portion of BUSI dataset.

As for the test portion of the BUSI dataset, Xception attained the highest mean accuracy score again (88.3%), indicating its superiority in overall pixel classification. Additionally, Xception outperformed the other in terms of the mean IoU metric, achieving the highest score of 77.6%, signifying its superior ability to accurately delineate object boundaries. Notably, Xception also achieved the highest average score of tumor, equal to 69.8%. In terms of the mean BF score, ResNet-50 had the highest score of 65.2% (tied with Xception and U-Net). Finally, U-Net surpassed the others in terms of the median Dice score of tumor, achieving the highest value of 83.6%. In summary, it appears that Xception-based DeepLab V3+ has the best overall performance on the BUSI dataset. Unlike the weighted IoU, global accuracy, average Dice score of background (exceeding 90%, 94%, and 96% in all cases), the segmentation quality metrics used in this section highlight the limitations of our trained models. The mean BF score and the average Dice score of tumor are particularly problematic, falling within the 50-to-70% range among all cases. This observation suggests that it may be worthwhile to explore different loss functions during the training process with a better focus on the tumor pixels and boundaries.

| <b>Network</b> | <b>Global Accuracy (BUL)</b> | <b>Observed Latency</b> | <b>Memory Requirement</b> | <b>Number of MACs (multiply-and-accumulate operations)</b> |
|----------------|------------------------------|-------------------------|---------------------------|--|
| ResNet-18      | 0.964                        | 56 ms per sample        | 62.5352 MB                | 7.80e+09   |
| ResNet-50      | 0.966                        | 65 ms per sample        | 151.5281 MB               | 12.99e+09  |
| Xception       | 0.968                        | 82 ms per sample        | 89.2042 MB                | 11.63e+09  |
| SegNet         | 0.956                        | 59 ms per sample        | 2.5425 MB                 | 6.24e+09   |
| U-Net          | 0.975                        | 58 ms per sample.       | 107.7437 MB               | 35.18e+09  |

Table 5.6 Latency, size, and number of MACs associated with each model.

Table 5.6 summarizes the cost of each network under consideration. ResNet-18 emerges as the fastest model with the smallest latency of 56 ms per sample. Following closely, U-Net and SegNet exhibit the latencies of 58 and 59 ms per sample, respectively.

Examining the number of MACs, SegNet demonstrates the lowest computational complexity with 6.24e+09 MACs, followed by ResNet-18 with 7.80e+09 MACs. Conversely, U-Net requires 35.18e+09 MACs, indicating the heaviest computational load. It should be noted that a higher number of MACs does not necessarily translate into a higher latency. One possible explanation for this discrepancy is that different models may have different degrees of parallelism exploited by the GPU hardware during the inference process.

In terms of storage, SegNet stands out with the smallest footprint of 2.5425 MB, making it highly memory efficient. Conversely, ResNet-50 exhibits the largest footprint of 151.5281 MB, which is significantly greater than the other models. ResNet-18 falls in the middle, with a memory requirement of 62.5352 MB.

In summary, SegNet emerges as a favorable choice for tasks prioritizing speed and memory efficiency; however, its segmentation performance was inferior to the other models for the

majority of metrics. Meanwhile, Xception-based DeepLab V3+ offers a more balanced option, providing a good compromise between segmentation quality and computational complexity. However, it should be mentioned that input image size for Xception is 299x299 as opposed to 224x224 used by the other models. Consequently, the latency, footprint, and number of MACs of Xception would be smaller when its input image size matches 224x224.

## Chapter 6. Conclusion and Future Work

### 6.1. Conclusion

This report explored a deep learning-based method for segmenting breast US images to distinguish between tumor and background regions. The Breast US Images (BUSI) and Breast US Lesions (BUL) datasets were utilized for performance evaluation purposes. Data augmentation was applied during training to increase the diversity of data samples.

DeepLab V3+, SegNet, and U-Net were chosen and optimized under similar training settings and dataset conditions. In the case of DeepLab V3+, three different backbones were used: ResNet-18, ResNet-50, and Xception. The segmentation performance was evaluated using two sources of test data: the unseen BUL dataset and the test portion of the BUSI dataset. Each model's output was assessed based on the following criteria: global accuracy, mean accuracy, mean IoU, weighted IoU, mean BF score, average Dice score of background, average Dice score of tumor, median Dice score of tumor, and the model cost.

The experiment results showed that DeepLab V3+ with Xception outperformed the other models in terms of the accuracy, average Dice score of tumor, and median Dice score of tumor on the unseen BUL dataset. However, the best global accuracy, mean IoU, weighted IoU, mean BF score, and average Dice score of background for the same dataset were due to U-Net. SegNet had the lowest memory requirements and the lowest computational complexity (in terms of the number of MACs). In contrast, Xception had the largest latency, while U-Net had the largest storage footprint and the largest MAC count, among all five models under consideration. Overall, our results indicate that Xception-based DeepLab V3+ and U-Net are the two leading competitors for handling segmentation tasks in breast US images.

### 6.2. Future work

For future work, the training process could be extended to incorporate alternative loss functions and/or additional datasets targeting other tumor segmentation tasks, such as lung and brain tumors. Additionally, to address the computational complexity and memory requirements of deep learning neural networks, the pruning technique proposed in [28], for example, could be employed. This technique aims to create a simpler and more power-efficient model while maintaining good segmentation performance. Its layer-by-layer pruning algorithm prunes more

aggressively than previous methods by minimizing the error in the output feature maps rather than the filter weights.

The images used in this project had tumors that were relatively small in size compared to the background. Therefore, one could apply additional techniques to address the class imbalance issue. As mentioned previously, this could involve alternative loss functions (e.g., Tversky loss). This issue could also be addressed with randomized resampling techniques, either by oversampling the minority class (e.g., duplicating samples) or undersampling the majority class (e.g., removing samples), until a more balanced distribution is achieved.

Furthermore, the work could be extended by exploring other segmentation models beyond DeepLab V3+, SegNet, and U-Net. For instance, models like DeepLab V3+ with MobileNet-v2 [38] or Inception-ResNet-v2 [39] bondnets, Mask R-CNN [40], and Faster R-CNN [41] could be investigated for their suitability in breast tumor segmentation tasks. Additionally, exploring unsupervised learning approaches, such as autoencoders or Generative Adversarial Networks (GANs), could offer alternative ways for tumor segmentation, particularly in scenarios where the amount of labeled data is limited.

## Bibliography

- [1] D. Avola, L. Cinque, A. Fagioli, G. Foresti, and A. Mecca, “US Medical Imaging Techniques: A Survey”, *ACM Computing Surveys*, vol. 54, no. 3, pp. 1–38, 2022, doi: 0.1145/3447243.
- [2] R. W. Cootney, “US Imaging: Principles and Applications in Rodent Research”, *ILAR Journal*, vol. 42 (3), pp. 233-247, Jan 2001.
- [3] H. D. Cheng, X. Cai, X. Chen, L. Hu, and X. Lou, “Computer-aided detection and classification of microcalcifications in mammograms: a survey”, *Pattern Recognition*, vol. 36, no. 12, pp. 2967–91, 2003, doi:10.1016/S0031-3203(03)00192-4.
- [4] J. A. Noble and D. Boukerroui, "Ultrasound image segmentation: a survey", *IEEE Transactions on Medical Imaging*, vol. 25, no. 8, pp. 987-1010, Aug 2006, doi: 10.1109/TMI.2006.877092.
- [5] M. Franke and H. P. Kohl, “Second-generation real-time 3-D echocardiography: A revolutionary new technology”, *Medicamundi*, vol. 47, no. 2, pp. 3440, 2003.
- [6] M. Mulet-Parada, “Intensity independent feature extraction and tracking in echocardiographic sequences”, *Ph.D. dissertation*, Dept. Eng. Sci., Univ. Oxford, Oxford, U.K., 2000.
- [7] B. Liu, H. D. Cheng, J. Huang, J. Tian, X. Tang, and J. Liu, “Probability density difference-based active contour for Ultrasound image segmentation”, *Pattern Recognition*, vol. 43, no. 6, pp. 2028–42, 2010, doi:10.1016/j.patcog.2010.01.002.
- [8] K. D. Marcomini, H. Schiabel, and A. A. O. Carneiro, “Quantitative evaluation of automatic methods for lesions detection in breast US images”, *SPIE Medical Imaging*, International Society for Optics and Photonics, vol. 8670, pp. 27, Feb 2013, doi: 10.1117/12.2008056.
- [9] T. Binder, M. Süssner, D. Moertl, T. Strohmer, H. Baumgartner, G. Maurer, and G. Porenta, “Artificial neural networks and spatial temporal contour linking for automated endocardial contour detection on echocardiograms: A novel approach to determine left ventricular contractile function”, *Ultrasound in Medicine & Biology*, vol. 25, no. 7, pp. 1069–76, 1999, doi:10.1016/S0301-5629(99)00059-9.
- [10] Vincent Chan and Anahi Perlas, “Basics of Ultrasound Imaging”, *Atlas of Ultrasound-Guided Procedures in Interventional Pain Management*, Springer New York, pp. 13–19, 2010, doi:10.1007/978-1-4419-1681-5\_2.
- [11] W. S. McCulloch and Walter Pitts., “A logical calculus of the ideas immanent in nervous activity”, *Bulletin of Mathematical Biology*, vol. 52, no. 1, pp. 99–115, 1990, doi:10.1016/S0092-8240(05)80006-0.
- [12] G. E. Hinton, S. Osindero, and Yee-Whye Teh., “A fast learning algorithm for deep belief nets”, *Neural Computation*, vol. 18, no. 7, pp. 1527–54, 2006, doi:10.1162/neco.2006.18.7.1527.

- [13] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning”, *Nature (London)*, vol. 521, no. 7553, pp. 436–44, 2015, doi:10.1038/nature14539.
- [14] A. Carovac, F. Smajlovic, and D. Junuzovic., “Application of Ultrasound in Medicine”, *Acta Informatica Medica*, vol. 19, no. 3, pp. 168–71, 2011, doi:10.5455/aim.2011.19.168-171.
- [15] G. Stetten, et al., “C-Mode Real-time Tomographic Reflection for a Matrix Array Ultrasound Sonic Flashlight”, *Academic Radiology*, vol. 12, no. 5, pp. 535–43, 2005, doi:10.1016/j.acra.2004.06.011.
- [16] Y. Yang, X. Xia, D. Lo, and J. Grundy, “A Survey on Deep Learning for Software Engineering”, *ACM Computing Surveys*, vol. 54, no. 10S, pp. 1–73, 2022, doi:10.1145/3505243.
- [17] Shi Dong, Ping Wang, and Khushnood Abba, “A survey on deep learning and its application”, *Computer Science Review*, vol. 40, pp. 100379-, 2021, doi:10.1016/j.cosrev.2021.100379.
- [18] G. Csurka, D. Larlus, and F. Perronnin, "What is a good evaluation measure for semantic segmentation?", *British Machine Vision Conference*, pp. 32.1–32.11, 2013.
- [19] Diederik Kingma and Jimmy Ba., " Adam: A method for stochastic optimization", *International Conference on Learning Representations*, Dec 2014.
- [20] Al-Dhabyani, Walid, et al., “Dataset of Breast Ultrasound Images”, *Data in Brief*, vol. 28, pp. 104863–104863, 2020, doi:10.1016/j.dib.2019.104863.
- [21] Yap, Moi Hoon, et al., “Automated Breast Ultrasound Lesions Detection Using Convolutional Neural Networks”, *IEEE Journal of Biomedical and Health Informatics*, vol. 22, no. 4, pp. 1218–26, 2018, doi:10.1109/JBHI.2017.2731873.
- [22] Chen, L., Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, “Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation”, *arXiv.Org*, 2018, doi:10.48550/arxiv.1802.02611.
- [23] Liu, Wei, et al., “ParseNet: Looking Wider to See Better”, *arXiv.Org*, 2015, doi:10.48550/arxiv.1506.04579.
- [24] He, Kaiming, et al., “Deep Residual Learning for Image Recognition”, *2016 IEEE CONFERENCE ON COMPUTER VISION AND PATTERN RECOGNITION (CVPR)*, IEEE, pp. 770–78, 2016, doi:10.1109/CVPR.2016.90.
- [25] Olaf Ronneberger, et al., “U-Net: Convolutional Networks for Biomedical Image Segmentation”, *arXiv.Org*, 2015, doi:10.48550/arxiv.1505.04597.
- [26] J. Wang and X. Liu, “Medical image recognition and segmentation of pathological slices of gastric cancer based on DeepLab V3+ neural network”, *Computer Methods and Programs in Biomedicine*, vol. 207, pp. 106210–106210, 2021, doi:10.1016/j.cmpb.2021.106210.

- [27] Su, Feng, et al., “Development and Validation of a Deep Learning System for Ascites Cytopathology Interpretation.” *Gastric Cancer: Official Journal of the International Gastric Cancer Association and the Japanese Gastric Cancer Association*, vol. 23, no. 6, pp. 1041–50, 2020, doi:10.1007/s10120-020-01093-1.
- [28] Tien-Ju, Yang, et al., “Designing Energy-Efficient Convolutional Neural Networks Using Energy-Aware Pruning” *arXiv.Org*, 2017, doi:10.48550/arxiv.1611.05128.
- [29] Chollet, François, “Xception: Deep Learning with Depthwise Separable Convolutions”, *arXiv.Org*, 2016, doi:10.48550/arxiv.1610.02357.
- [30] Srinivasan, Kathiravan, et al., “Performance Comparison of Deep CNN Models for Detecting Driver’s Distraction”, *Computers, Materials & Continua*, vol. 68, no. 3, pp. 4109–24, 2021, doi:10.32604/cmc.2021.016736.
- [31] Badrinarayanan, Vijay, et al., “SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 12, pp. 2481–95, 2017, doi:10.1109/TPAMI.2016.2644615.
- [32] “Evaluate semantic segmentation data set against ground truth - MATLAB evaluateSemanticSegmentation”, <https://www.mathworks.com/help/vision/ref/evaluatesemanticsegmentation.html>.
- [33] D. T. Tzou, et al., “Ultrasound Use in Urinary Stones: Adapting Old Technology for A Modern-day Disease”, *Journal of Endourology*, vol. 31 (S1), Apr 2017.
- [34] S. Liu and W. Deng, “Very deep convolutional neural network based image classification using small training sample size”, 3<sup>rd</sup> IAPR Asian Conference on Pattern Recognition (ACPR), Kuala Lumpur, Malaysia, pp. 730-734, 2015, doi: 10.1109/ACPR.2015.7486599.
- [35] M. H. Beale, M. T. Hagan, and H. B. Demuth, “Deep Learning Toolbox”, <https://www.mathworks.com/help/deeplearning/index.html>.
- [36] “Computer Vision Toolbox”, <https://www.mathworks.com/help/vision/index.html>.
- [37] “Image Processing Toolbox”, <https://www.mathworks.com/help/images/index.html>.
- [38] Mark Sandler, et al., “MobileNetV2: Inverted Residuals and Linear Bottlenecks”, *2018 IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, IEEE, pp. 4510–20, 2018, doi:10.1109/CVPR.2018.00474.
- [39] Christian Szegedy, et al., “Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning”, *arXiv.Org*, 2017, doi:10.48550/arxiv.1602.07261.
- [40] Kaiming He, et al., “Mask R-CNN”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 2, pp. 386–97, 2020, doi:10.1109/TPAMI.2018.2844175.

- [41] Shaoqing Ren, et al., “Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137–49, 2017, doi:10.1109/TPAMI.2016.257703.
- [42] Forqan Ali Wahhab, “Lecture: Ultrasound Imaging”, *Al-Mustaqbal University College*, [https://uomus.edu.iq/img/lectures21/MUCLecture\\_2022\\_112516382.pdf](https://uomus.edu.iq/img/lectures21/MUCLecture_2022_112516382.pdf).
- [43] “What causes throbbing tinnitus?”, [read://https\\_belkashop.com.ua/?url=https%3A%2F%2Fbelkashop.com.ua%2Fvid-chogovinikaye-pulsuyuchij-shum-u-vusi%2F](read://https_belkashop.com.ua/?url=https%3A%2F%2Fbelkashop.com.ua%2Fvid-chogovinikaye-pulsuyuchij-shum-u-vusi%2F).