

UNIVERSITY

FACULTY

STUDY PROGRAMME:

MASTER THESIS

Catchy Title

Submitted by:

John DOE

Matriculation Number: 123456

Street 123

789654 City

Initial Examiner:

Dr. Supervisor

Secondary Examiner:

Prof. Dr. Examiner

December 19, 2023

GitHub

First printed interior sheet: Recto: same text and formatting as front cover. o Verso:
ISBN, theme code, NUR code(s) and legal depot number. Position: bottom of the page.

Second printed interior sheet, recto: Members of the examination board - Title: 'Members of the Examination Board' / 'Leden van de examencommissie' (same language as front cover), - Members, per category: - Chair / Voorzitter - Other members entitled to vote / Andere stemgerechtigde leden - Supervisor(s) / Promotor(en) (they are members!) - (it is no longer needed to indicate the secretary) - Every member is listed as Title Name, Affiliation (incl. country if not BE)

UNIVERSITY

Abstract

Faculty

Master Thesis

Catchy Title

by John DOE

Write your abstract here.

Keywords: L^AT_EX, Master Thesis, Programming

Acknowledgements

Ich bedanke mich bei Dr. Supervisor für die Betreuung dieser Arbeit.

Contents

| | |
|---|------------|
| Abstract | III |
| Acknowledgements | IV |
| List of Figures | VII |
| List of Tables | X |
| List of Abbreviations | XI |
| 1 Introduction | 1 |
| 2 Neural Network-Based Segmentation of Biomedical Images | 2 |
| 2.1 Common Types of Biomedical Images | 2 |
| 2.1.1 3D Modalities | 2 |
| 2.1.2 2D Modalities | 6 |
| 2.2 Image Segmentation: From Images to Segmentation Maps | 7 |
| 2.2.1 Traditional Image Processing Methods | 8 |
| 2.2.2 Machine Learning | 10 |
| 2.3 Deep Learning-Based Segmentation Methods | 12 |
| 2.3.1 Neural Network Training, Validation and Testing | 13 |
| 2.3.2 Encoders and Decoders | 14 |
| 2.3.3 Convolutional Neural Networks | 15 |
| 2.4 CNN Architectures for Medical Image Segmentation | 19 |
| 2.4.1 Fully Convolutional Network (FCN) | 19 |
| 2.4.2 U-Net and Its Variants | 20 |
| 2.4.3 Mask R-CNN | 23 |
| 2.4.4 Other Notable Segmentation CNNs | 24 |
| 2.5 Fully Connected Transformers for Medical Image Segmentation | 25 |
| 2.5.1 Attention | 25 |
| 2.5.2 Positional Encoding | 26 |
| 2.5.3 Adapting Transformers to Image Segmentation | 27 |
| 3 Data Efficiency in Neural Network-Based Image Segmentation | 29 |
| 3.1 Transfer Learning | 32 |
| 3.1.1 Simple Transfer Learning | 32 |
| 3.1.2 Domain Adaptation | 33 |
| 3.1.3 Semi-Supervised and Self-Supervised Learning | 34 |
| 3.2 Synthetic Data | 35 |
| 3.3 Regularization | 35 |
| 3.3.1 Conclusion | 36 |

| | |
|--|-----------|
| 4 Data Efficiency via Model-Driven Preprocessing | 37 |
| 4.1 Motivation: Using The Polar Transform for Preprocessing | 37 |
| 4.2 Model-driven Preprocessing | 39 |
| 4.2.1 Training Model-Driven Preprocessing Networks | 41 |
| 4.3 Preprocessing With the Polar Transform Using Centerpoint Prediction | 42 |
| 4.3.1 Previous work using the polar transform for medical image processing | 44 |
| 4.3.2 Methodology | 44 |
| 4.3.3 Centerpoint prediction | 45 |
| 4.3.4 Experiments | 47 |
| 4.3.5 Datasets description | 48 |
| 4.3.6 Implementation details | 49 |
| 4.3.7 Results | 50 |
| 4.4 Discussion | 50 |
| 4.5 Conclusion | 56 |
| Bibliography | 57 |

List of Figures

| | | |
|------|--|----|
| 2.1 | A cardiac CTA in its full range (left) and windowed (right). [5] | 4 |
| 2.2 | An example T1-weighted MRI (left) and a T2-weighted MRI (right) showing a diffuse glioma. [6] | 5 |
| 2.3 | A demonstration of region growing for delineating the internal and external areas of the pericardium on a CT slice, set to the adipose tissue intensity range. The left image is the original input, and the right image depicts the segmented outcome with the heart exterior in red and the interior in blue. The green dots represent the manually chosen seed points initiating the region-growing technique. [4] | 9 |
| 2.4 | A demonstration of employing active contours to complete the absent segments of the pericardium line, displayed in white. The contour, illustrated in blue, starts as a complete circle surrounding the image. With every iteration, the contour adapts more closely to the image's shape. [4] | 10 |
| 2.5 | A schematic representation of the registration procedure. Initially, input and target images are chosen. Throughout the registration phase, the input image (shown in green) undergoes deformation to align with the fixed target image (shown in red). [4] | 11 |
| 2.6 | A schematic of a supervised linear classifier in a machine learning workflow. The upper section illustrates the training phase. Here, features are color-coded according to their known class from training data, depicted in red and blue. The parameters of the decision boundary, which demarcates the zones of the two classes (highlighted in light red and grey), are determined during training. The lower part of the diagram depicts the inference stage. In this phase, features are extracted from new images, and the trained model is employed to classify each pixel in the image. [4] | 12 |
| 2.7 | A small neural network with three layers, wherein each layer is connected to every neuron in the next layer. [16] | 13 |
| 2.8 | A visual depiction of a convolution procedure step-by-step. In each step, the kernel slides over the image (shown in blue). The overlapping elements between the kernel and the image are multiplied and then summed to produce a value in the resultant image (shown in green). [18] | 15 |
| 2.9 | An example of an input image (left) convolved with the Prewitt operator (right). Note that vertical edges are accentuated in the convolution result. | 16 |
| 2.10 | A view of one step of a single convolution operation inside a convolutional neural network (CNN) layer. The layer performs multiple convolutions, each with a different kernel that has an equal number of channels as the input image. In each step, the whole kernel slides over the width and height of the image, and the overlapping channels are multiplied together and summed to produce a single output value. The output of the convolution is one channel of a n -channel image where n is the number of different kernels in the layer. | 18 |

| | | |
|------|---|----|
| 2.11 | A typical architecture of a CNN encoder. The encoder consists of consecutive convolutional and pooling layers that gradually increase the feature map depth and decrease its width and height. The result is a map of features that tells the decoder what features are on the image but does not provide much spatial information about the location of those features. [19] | 18 |
| 2.12 | A diagram of how FCN forms predictions based on the output of different encoder layers. Encoder layers are shown on the left and the grid represents the coarseness of the feature map. The maps are combined at three different levels to produce three predictions. Each prediction is compared with the ground truth during training, but for inference only the 8x upsampled prediction is used. [20] | 20 |
| 2.13 | A diagram of the U-Net model. The output of each layer of the encoder is concatenated to the input of its corresponding layer in the decoder. [21] | 21 |
| 2.14 | A diagram of the nnU-Net procedure of creating a training configuration. [22] | 22 |
| 2.15 | A comparison between U-Net (left) and U-Net++ (right). Each node in the graph represents a convolutional layer. The dashed arrows represent skip connections, while full arrows are downsampling or upsampling operations. [24] | 23 |
| 2.16 | A diagram of the Mask R-CNN architecture. Two parallel decoder branches are used to achieve segmentation and object detection simultaneously. [17] | 24 |
| 2.17 | A visualization of a single encoder layer in a transformer network. This shows the encoding process for one element of the input sequence. | 26 |
| 2.18 | The ViT architecture for image classification. [29] | 28 |
| 3.1 | A simulated example of heart disease risk prediction using simple linear regression. The plots on the right show how the approximated function depends on the sample size. | 30 |
| 3.2 | Fitting a polygonal function of various degrees on three different sample sizes. | 31 |
| 3.3 | An overview of the simple transfer learning procedure. First, a model is pretrained on some related segmentation task. Then, part of the trained model's weights are copied to a new model that is then fine-tuned on the target segmentation task. | 33 |
| 3.4 | A diagram of the domain adaptation approach in Ganin and Lempitsky [40]. The gradients of the domain classification head which are applied to the encoder are reversed during backpropagation. | 34 |
| 3.5 | A self-supervised learning approach using shuffling image tiles as a pretext task presented by Carr, Berthet, Blondel, <i>et al.</i> [43]. The trained feature encoder learns to extract relevant features and its parameters are transferred to a model trained to perform the downstream task. | 34 |
| 4.1 | An example of how a 2-class classification neural network bends an input space (left) such that it is linearly separable with a hyperplane. | 38 |
| 4.2 | The decision boundary (shown as the background color) according to a fully connected neural network with one hidden layer where the hidden layer has one, two, or three neurons. | 38 |
| 4.3 | An example circular dataset transformed using the polar transform. The background color shows the decision boundary of a fully connected network with one hidden layer with a single neuron. | 39 |

| | | |
|------|--|----|
| 4.4 | The polar transformation of the simulated data using different polar origins (C) | 39 |
| 4.5 | An example of a polar transformation. | 43 |
| 4.6 | An example image and label from the lesion dataset and their corresponding polar transformation. | 43 |
| 4.7 | A diagram of the approach of predicting polar origins from a Cartesian network. The first network performs an initial segmentation, which is then used to extract a polar origin for the polar transformation. The method does not rely on any specific neural network architecture. The Polar and Cartesian network can be any neural network which takes an input image and produces a binary segmentation mask as output. The red point shows the extracted polar origin. The Polar network is trained on polar image transformations. The polar transformation is not part of the network itself, but happens as a preprocessing step for the Polar network. | 46 |
| 4.8 | A diagram of the approach of using a centerpoint prediction network. The first network can be any neural network which predicts a heatmap from an input image, which is then used to extract polar origin, shown as a red point. The Polar network can be any semantic segmentation neural network which produces a binary mask output from an input image. The Polar network is trained on polar image transformations. The polar transformation is not part of the network itself, but happens as a preprocessing step for the Polar network. | 46 |
| 4.9 | Examples of heatmaps generated from the ground truth data. The heatmap is a gaussian centered on the center of mass of the ground-truth label, shown as a blue point on the input images. | 47 |
| 4.10 | Example input images and ground-truth labels for each dataset used in our experiments. | 48 |
| 4.11 | The training and validation Dice coefficient (DSC) of the polar and Cartesian U-Net models during training. | 50 |
| 4.12 | The relationship between mean squared errors of the centers used for the polar transformation and segmentation performance of the polar network on the lesion dataset. The mean squared errors are calculated compared to the ground-truth centers. | 54 |
| 4.13 | The best Dice coefficient by epoch 50 for models trained on subsets of the lesion training dataset. | 54 |
| 4.14 | A random sampling of inverse polar transformed predictions from the polar network with the polar origins predicted from the centerpoint predictor for various datasets. The prediction is shown in green and overlaid on top of the original input image. EAT predictions (d) are cropped and zoomed to better show the predictions. | 55 |

List of Tables

| | | |
|-----|---|----|
| 2.1 | Approximate Hounsfield values of various tissues [1], [2]. | 4 |
| 4.1 | Results of our proposed approaches for different tasks for three different neural network architectures. The cartesian network is the network trained on Cartesian images. "GT centers" refers to obtaining a polar origin from the ground-truth labels and segmentation using the polar network. "Cartesian centers" refers to predicting the polar origins from the Cartesian network and then performing segmentation using the polar network. "Model centers" refers to using the center-point predictor to obtain polar origins. (Continued on the next page.) | 51 |
| 4.1 | Results of our proposed approaches (continued). | 52 |
| 4.2 | A comparison between our method (approach with best results) and the state of the art on the same datasets. | 53 |
| 4.3 | Ablation study of our approach for the polyp dataset. | 55 |

List of Abbreviations

CNN convolutional neural network

1 Introduction

2 Neural Network-Based Segmentation of Biomedical Images

This chapter introduces biomedical images as well as the process of segmenting them. We will first provide an overview of the world of biomedical images and its diverse modalities and imaging techniques. We will touch on the technical details of how these images are stored and used in segmentation algorithms. Then, we will more formally describe the process of image segmentation from traditional segmentation methods to the most current deep learning methods.

Biomedical imaging covers a wide range of techniques used in both biology and medicine. This includes complex modalities like CT scans and simpler ones like photographs of a patient's skin or plants. Because of this variety, a method designed for one modality might not work for another without changes. So, understanding each modality's specifics is crucial when considering segmentation methods. The next section will provide an overview of different biomedical imaging modalities, focusing on those most relevant to this thesis.

2.1 Common Types of Biomedical Images

Broadly, we may classify biomedical images into 3D and 2D images. 3D images include modalities such as CT and MRI. Examples of 2D image modalities, among others, include X-ray imaging, most microscopic images, and dermatoscopic images. However, it's worth noting that biomedical images don't always fit neatly into distinct categories. While our classification provides a general overview, there are emerging techniques that merge different modalities, blurring the classification boundaries presented here.

2.1.1 3D Modalities

It is hard to overstate the importance of 3D modalities such as CT and MRI in contemporary medical practice, scientific advancements, and patient outcomes. These techniques offer medical professionals a non-invasive window into the patient's body, allowing much safer and more reliable diagnosis, treatment, and surgical planning.

The foundational principle behind these modalities revolves around emitting a signal capable of penetrating tissues. This signal interacts variably with different tissue types. After transmission through the body, the residual signal is detected and used to create an image. By repeating this process across various planes of the body and leveraging computational algorithms, a 3D volume of the patient's anatomy is constructed. From there, these modalities are often stored as voxel-based 3D files, where a voxel is the smallest 3D unit equivalent to a pixel in 2D images. They are usually stored alongside patient data such as age, sex, imaging parameters, and other relevant information.

In the subsequent sections, we will delve deeper into the specifics of the two predominant 3D biomedical imaging modalities: CT and MRI.

Computed Tomography (CT)

Computed tomography is an X-ray-based imaging technique where different planes of the subject are captured and then reconstructed into a 3D image using a process called tomography. CT, as opposed to MRI, uses ionizing radiation that can be harmful in large doses. As such, the benefit of each scan should be weighed against potential hazards. Notably, there is a relationship between radiation dose and image quality. The radiation dose has to be carefully balanced to achieve the required image quality without subjecting the patient to undue radiation.

This balance between radiation dose and image quality has implications for segmentation methods. Low-dose and high-dose CT images visually differ — enough to cause issues in the segmentation model's performance across these two domains. High-quality images result in better neural network models, but their availability is much more limited than low-dose scans.

Another way to enhance the quality of CT images is the use of contrast agents. Typically administered intravenously, these agents are designed to have a high intensity on the resulting CT image. A common example of this is CT angiography where contrast is used to make it easier to delineate blood vessels from surrounding tissue.

CT images are generally stored as $W \times H \times D$ matrices, with each element indicating the grayscale intensity of a voxel (*volumetric pixel*) within the 3D volume. In contrast with regular pixels, voxels can have different lengths along each dimension and usually aren't cubes. The voxel's depth (along the z-axis) is known as slice thickness, and it is chosen based on the task at hand. Thinner slices increase the spatial detail in the image making it easier to see minute details. However, thinner slices generally also increase the noise in the image. Slice thickness is selected to maintain a tradeoff between spatial resolution and level of noise. The width and height of the voxel are governed by the field of view and the scanner itself.

Voxel values in CT images are typically stored as 12-bit signed integers. These values represent the attenuation of the X-ray as it passes through various tissues. Denser structures such as bones attenuate the radiation more strongly than fat tissue and thus result in higher intensity values.

The values of the voxels are usually stored as 12-bit signed integer values. The value corresponds to the attenuation of the X-ray as it passes through the tissue. Denser structures such as bones attenuate the radiation more strongly than fat tissue and thus result in higher intensity values.

To maintain consistency across scans and machines, the attenuation values in a CT are linearly transformed. Specifically, air is assigned a value of -1000, while water is attributed a value of 0, under standard temperature and pressure conditions. Mathematically, this transformation is expressed as:

$$HU(\mu) = 1000 \cdot \frac{\mu - \mu_{\text{water}}}{\mu_{\text{water}} - \mu_{\text{air}}}, \quad (2.1)$$

where μ is the attenuation value of a given voxel while μ_{water} and μ_{air} are attenuation values of water and air, respectively. This transformation is called the **Hounsfield unit** (HU) scale. Calibration according to this scale ensures that the attenuation of any tissue is always relative to water's attenuation. As a result, values remain standardized across different images, CT machines, imaging parameters, and institutions. The approximate Hounsfield unit values corresponding to various tissues are shown in [Table 2.1](#).

The Hounsfield unit values play a pivotal role in the analysis of CT scans. Beyond facilitating the identification of specific tissues or structures, variations in Hounsfield

| Tissue | Hounsfield value |
|--------|------------------|
| Fat | -30 to -70 HU |
| Muscle | 20 to 40 HU |
| Bone | 1000 HU |
| Blood | 13 to 50 HU |

Table 2.1: Approximate Hounsfield values of various tissues [1], [2].

values of the same tissues across patients or scans can hold significant diagnostic implications. For example, clotted blood exhibits a higher HU value than its unclotted counterpart, establishing the Hounsfield value as a potential marker for intracranial hemorrhage [2]. In another context, the mean HU value of epicardial fat can be indicative of myocardial infarction [3].

The range of gray levels in a CT scan is far wider than what the human eye can discern. Therefore, it's standard practice for professionals to employ a technique called "windowing" when reviewing scans. In essence, windowing compresses the value range by applying two specific thresholds (t_1, t_2), as defined by:

$$y(x) = \begin{cases} t_1, & \text{if } x \leq t_1 \\ x, & \text{if } t_1 < x < t_2 \\ t_2, & \text{if } x \geq t_2 \end{cases} \quad (2.2)$$

where x is a given HU value. This allows CT visualization software to visually stretch the remaining range of gray values, thus enhancing the differences in attenuation, as seen in [Figure 2.1](#). This technique isn't limited to human experts; it finds applications in image segmentation as well. It is common to window CT image inputs into segmentation neural networks. For example, if segmenting fat, it is often beneficial to discard all voxels outside of the fat tissue range [4], reducing the number of voxels that the network needs to process.

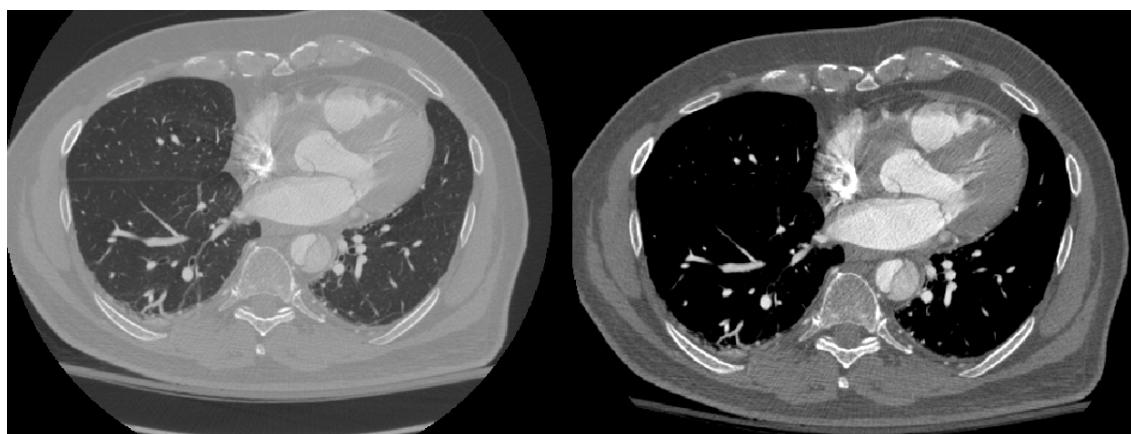


Figure 2.1: A cardiac CTA in its full range (left) and windowed (right). [5]

Magnetic Resonance Imaging (MRI)

Much like CT, MRI produces voxel-based images that visualize the insides of an object. The process of obtaining an MRI image can be separated into several steps:

1. The machine applies a strong magnetic field to the subject. This causes protons in the body to align in the same direction — parallel to the z-axis.

2. Next, the machine emits a brief radio frequency pulse, perturbing the aligned protons and causing them to deviate from their uniform alignment.
3. After the pulse is turned off, the protons gradually return back to alignment. As they realign, the movement of their positive charge induces an electric current in a coil inside the MRI machine. In contrast to CT, which gauges the attenuation of X-rays, MRI measures this induced current.

The time to reach the equilibrium state depends on the specific tissue type, so the measured current allows the differentiation of tissues in the image. The equilibrium state is achieved through two independent processes: T1 and T2. T1 measures the time it takes the protons to reach equilibrium longitudinal alignment, while T2 measures the time to regain its equilibrium transverse alignment. Different tissues return to their equilibrium states at varying rates for T1 and T2, enabling two different types of MRI images: T1- or T2-weighted. For example, water exhibits a long T1 time, while fat has a shorter one. This means that in T1-weighted images, fat appears brighter than water, while the opposite is true in T2-weighted images. This can be seen in [Figure 3.1](#).

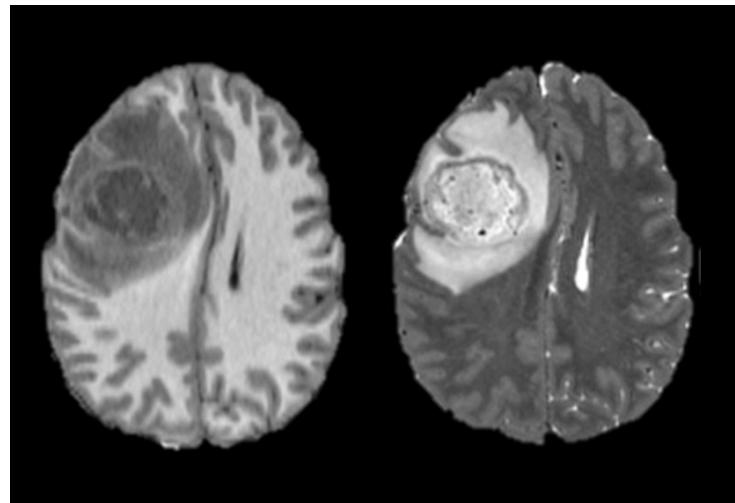


Figure 2.2: An example T1-weighted MRI (left) and a T2-weighted MRI (right) showing a diffuse glioma. [6]

The relative differences in T1 and T2 images make each more suitable for imaging specific tissues. For instance, T1-weighted images allow for identifying fatty tissue or detecting liver lesions, while T2-weighted images are more effective for identifying white matter lesions and edemas.

The voxels in an MRI, much like a CT, are generally rectangular solids and can be different in width, height, and depth. The image resolution is determined by several factors, including the scanner's field of view, the MRI machine itself, and specific imaging parameters. Generally, increasing resolution leads to higher levels of noise and a longer acquisition time. Consequently, an optimal balance between resolution, noise, and scan time must be found, tailored to the tissue being imaged and the specific diagnostic objective.

MRI has several advantages over CT. For one, MRI usually offers superior contrast, especially when imaging soft tissues. This has established MRI as the gold standard for a large number of diagnostic procedures. Much like with CT scans, MRI contrast can be further enhanced using contrast agents. Furthermore, as MRI uses magnetic fields instead of ionizing radiation, it poses no radiation-induced risks to patients.

However, MRI has its own set of limitations. MRI scans typically take longer than CT scans, which can be uncomfortable for patients, particularly those who are claustrophobic or have mental and cognitive disabilities. MRI is contraindicated for individuals with certain non-removable magnetic objects, like coronary pacemakers and specific implants. A notable drawback is the lack of standardization in MRI values across different scans and machines. Contrary to CT scans, an MRI voxel's intensity isn't consistent across different scans and can only be reliably interpreted in relation to adjacent tissue in the same scan.

2.1.2 2D Modalities

3D modalities, while offering a detailed view of a subject, often entail time-consuming and infeasible procedures in common clinical use. 2D modalities such as X-ray and diagnostic ultrasound are often quicker and more readily available. This ease and speed of 2D modalities contribute to the large amount of publicly available datasets in 2D modalities compared to 3D ones. For example, analyzing X-ray images is one of the most active fields in computerized medical image analysis research [7], [8].

X-ray imaging

X-ray imaging, also known as radiography, can be thought of as the 2D equivalent of a CT as it also relies on ionizing radiation. To obtain an image, X-rays are emitted on one side of the subject and captured on the opposite side. The intensity of the pixels corresponds to the attenuation of the emitted radiation. Denser materials appear with a higher intensity on the resulting image due to their high attenuation.

When capturing an X-ray image, an expert manually positions the generator. The relative positioning of this generator and the subject directly influences the image's magnification and field of view. If the subject is closer to the detector than the generator, it will appear magnified on the resulting image. Such magnification makes the scale on an X-ray image non-standard and precludes the possibility of making objective length and area measurements.

The expert also determines various parameters that ultimately change the quality and quantity of the X-ray beams. The quality measures the ability of an X-ray to penetrate tissue and is proportional to the X-ray energy level. Quantity, on the other hand, measures the number of photons constituting the beam. A high beam quality is advantageous for imaging denser tissues, like bones, or traversing larger volumes. However, high-quality beams reduce the contrast in soft tissues. Owing to the variability introduced by these parameters, X-ray image intensities lack the uniform standardization found in CT scans. Consequently, there's no standardized unit for gauging X-ray intensity across images.

Dermatological images (clinical and dermatoscopic)

The application of deep learning models to dermatology, especially skin lesion analysis, has gained significant momentum. This is driven in part by organizations such as the International Skin Imaging Collaboration (ISIC) which curates large dermatological datasets [9]. These datasets consist of two primary types of dermatological images: dermatoscopic and clinical. Clinical images are regular photographs of skin lesions, while dermatoscopic images are captured with a device called a digital epiluminescence dermatoscope. A dermatoscope consists of a camera attached to a magnifying lens with a

built-in light source. It allows capturing detailed and magnified images of a skin lesion while filtering out skin reflections.

The primary application of deep learning in dermatological images is for classification, e.g. predicting whether a lesion is benign or not or detecting the type of skin disease. However, by precisely outlining the boundary of a skin lesion, segmentation techniques can yield more consistent and objective descriptors of the lesion, aiding classification algorithms [9].

Microscopy

In biomedicine, one of the predominant uses of computer vision is in segmenting, analyzing, and quantifying microscopic images [10]. This covers a broad array of tasks in digital pathology, from identifying cancerous cells and segmenting nuclei to quantifying white blood cell counts.

Publicly available microscopic images of cells are abundant due to how frequently they are captured and that they don't contain any personally identifiable information. However, the large size of these images can present a challenge. Often spanning multiple megapixels, they exceed the capacity of current deep-learning models and machines. Consequently, to make them more manageable, these images are typically divided into smaller patches, downscaled, or processed using a coarse-to-fine approach [11].

As can be seen, both 2D and 3D biomedical images are exceedingly diverse in their appearance, use, technical details, and format. Yet, many segmentation techniques prove versatile enough to be applied across different modalities. In the next section, we will provide a general outline of how image segmentation works, highlighting commonly used methods.

2.2 Image Segmentation: From Images to Segmentation Maps

Image segmentation is the process of categorizing each pixel (or voxel) of an image into one of several predefined classes. Consider a 3-class segmentation scenario for CT images, with classes being liver, liver tumor, and background. Each of these classes is assigned a distinct numeric identifier, termed the **class label**. For illustration, the labels could be '0', '1', and '2' for the background, liver, and tumor respectively. The segmentation process outputs an image identical in dimensions to the original but with each pixel's value corresponding to the class label of that position in the original image. As an example, pixels corresponding to the liver will all have the value '1'. This resulting image is called a **segmentation map** since it, like a map, delineates important regions on the original image.

In medical image segmentation, it's common to extract only a single tissue type against the backdrop, a technique known as **binary segmentation**. Additionally, multi-class segmentation challenges can be decomposed into multiple binary segmentation tasks, with each class having its individual class-vs-background segmentation map. Therefore, one can frame any segmentation problem as a set of binary segmentation problems.

Mathematically, given a set of K classes, and an input d -dimensional image of N channels $I(A)$, $I \in \mathbb{R}^N$, $A \in \mathbb{N}^d$ where A is the location of each voxel, the segmentation map $M : \mathbb{N}^d \rightarrow \mathbb{R}^K$ maps each pixel location to a vector of class probabilities:

$$M(A) = (\Pr(C_1 | I(A)), \Pr(C_2 | I(A)), \dots, \Pr(C_K | I(A))), \quad (2.3)$$

where $\Pr(C_i | I(A))$ is the probability that the voxel $I(A)$ contains an object of class C_i . Expressed this way, the segmentation map is a K -channel image of the same size as $I(A)$. Each channel of the image corresponds to a probability map of finding an object of a given class at a given voxel location.

In binary segmentation, where $K = 2$, the classes can be denoted as C_1 for the background and C_2 for the target object. Given this, for all voxel locations A , the relation $\Pr(C_2 | I(A)) = 1 - \Pr(C_1 | I(A))$ holds true. Therefore $M(A)$ simplifies to a scalar value $M(A) = \Pr(C_2 | I(A))$. This representation is very common in medical image segmentation problems such as segmenting organs, cell nuclei, and skin lesions, among others.

Often, the next step in the process is to binarize the segmentation map $M(A)$. Voxels corresponding to the target object are set to ‘1’, while the background is marked as ‘0’. In such cases, $M(A)$ is frequently termed a **segmentation mask**. In computer vision, “mask” refers to a binary image $M_{01}(A) \in \{0, 1\}$ that hides (masks) regions in another image, resulting in a masked image $I_m(A) = I(A)M_{01}(A)$. Within the context of deep learning-based image segmentation, the terms “segmentation map” and “segmentation mask” are often used interchangeably.

Having laid the foundation for the mathematical framework of image segmentation, we can now transition to exploring specific segmentation techniques. The next section presents an overview of common segmentation methods, from traditional ones based on heuristics to complex model-based approaches prevalent today.

2.2.1 Traditional Image Processing Methods

Despite the popularity of deep learning, traditional techniques, rooted in fundamental image characteristics, still play a vital role in modern medical image segmentation. Traditional techniques are often used as methods for pre-processing and data augmentation, as well as refining deep learning model outputs. As we will show later in the dissertation, traditional methods can increase the robustness and data efficiency of deep learning-based segmentation models.

Rather than relying on data-intensive training phases, these techniques often operate deterministically, using explicit algorithms that manipulate image characteristics. They directly use image properties such as intensity and texture, combined with heuristic strategies that draw from empirical observations and domain knowledge. Heuristics are best-practice rules derived from previous samples and experiences. For instance, the longest component in the bone intensity range of an X-ray image is usually the femur. These methods provide a clear, interpretable pathway to segmentation.

However, traditional methods are usually developed with a very specific application in mind and are hard to translate to other tasks and domains without significant changes. They are also sensitive to parameter selection and properties of the image such as intensity level. Helpfully, since they don’t use a learning component their limitations can be known ahead of time. Their validity can also be confirmed using fewer samples than is the case for learning-based methods.

Image Thresholding

Thresholding is a fundamental technique that isolates regions in an image based on a specified range of intensities. In essence, it removes or retains portions of the image where the intensity either falls outside or within a given threshold.

To illustrate the use of thresholding, consider an example of segmenting adipose tissue on a CT scan. As mentioned earlier in this chapter, voxel intensities on CT scans

quantify a tissue's X-ray radiation attenuation as measured in Hounsfield units. Typically, fatty tissue lies between -250 HU and -30 HU. Thresholding the entire scan to this range segments fatty tissue from all other tissues on the scan with no need for additional complex models. For preprocessing, Hounsfield unit thresholding is an efficient way to discard irrelevant voxels, allowing the rest of the segmentation process to focus on fewer voxels.

Thresholding can be also used to greatly simplify a model's task using domain knowledge. Take, for instance, the segmentation of epicardial fat, which refers to the fat in proximity to the heart wall. Epicardial fat is sparsely distributed in a complex shape inside the pericardium. Yet, the pericardium itself has a smooth, elliptical shape and is comparatively easy to segment. By first segmenting the pericardium and subsequently thresholding the pericardium region to the fatty tissue range, we can find epicardial fat without having to segment its complex shape directly [4].

Region Growing Techniques

Region growing is a voxel-based image segmentation method [12]. Beginning from a designated seed voxel, the method assesses neighboring pixels in successive steps. If these pixels meet a specified criterion, often related to pixel values or textures, they are integrated into the region. This process continues until all image pixels are assessed. Despite being a straightforward method, region growing is very sensitive to the selected seed point and may struggle with nuanced transitions between regions. This process is shown in [Figure 2.3](#).

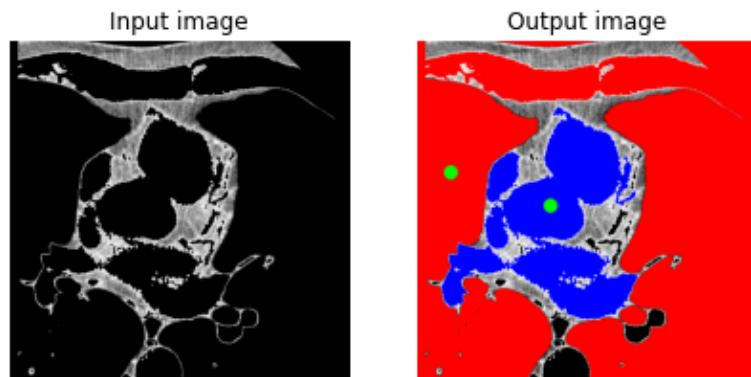


Figure 2.3: A demonstration of region growing for delineating the internal and external areas of the pericardium on a CT slice, set to the adipose tissue intensity range. The left image is the original input, and the right image depicts the segmented outcome with the heart exterior in red and the interior in blue. The green dots represent the manually chosen seed points initiating the region-growing technique. [4]

Active Contours or Snakes

Active contours, often termed "snakes", are segmentation methods that employ dynamic curves to outline image parts [13]. The process involves tightening a preliminary curve around an object iteratively until it conforms to the object's shape. The adaptation is guided by an energy function that evaluates the curve's smoothness and proximity to edges. A practical depiction of active contours is displayed in [Figure 2.4](#).

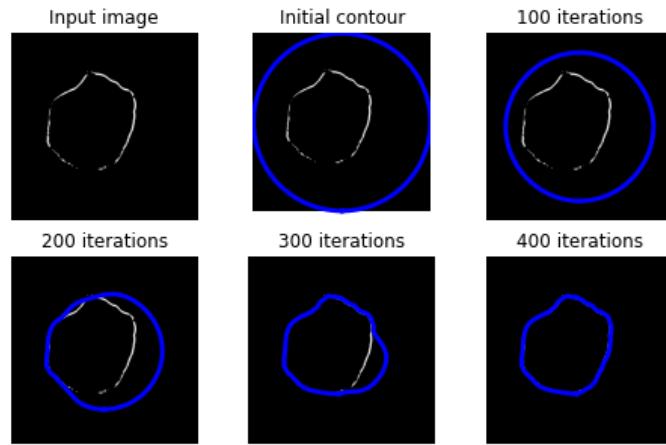


Figure 2.4: A demonstration of employing active contours to complete the absent segments of the pericardium line, displayed in white. The contour, illustrated in blue, starts as a complete circle surrounding the image. With every iteration, the contour adapts more closely to the image's shape. [4]

Atlas-Based Segmentation

Atlas-based methods, differing from contour-based ones, leverage the spatial relationships among identified structures in an image [14]. First, a template image is selected and an expert creates an atlas (a segmentation map) by manually segmenting and labeling structures in the template image. Due to anatomical variations, multiple representative images are often merged to produce a template image. The atlas can then be employed to segment new images using a registration algorithm. Image registration is an optimization problem where one image, called the moving image, is deformed to best align with a target image according to some scoring function. The scoring function usually uses the distance between heuristic-based landmarks in the target and moving images to determine how well the two images align. In atlas-based segmentation, a new moving image is deformed to be aligned with the template image that was used to construct the atlas. The atlas can then be used as a segmentation map for the moving image, as it is now aligned with the atlas. The atlas-based registration process is illustrated in Figure 2.5.

A disadvantage of this approach is that the process can lead to a complete failure to segment the image if the target and moving images are too dissimilar. Therefore, atlas-based segmentation has fallen out of favor due to the emergence of deep learning-based methods. However, recently there has been significant progress in image registration and atlas-based segmentation using deep learning-based techniques [15]. These approaches offer good potential for merging traditional and newer approaches will be discussed later in the chapter.

2.2.2 Machine Learning

While deep learning generally falls within the machine learning umbrella term, in this dissertation the term “machine learning” will be used to refer to techniques that are not based on deep neural networks. This encompasses methods such as support vector machines, random forests, and other statistical learning techniques that rely on manually engineered image features.

In this context, image segmentation can be viewed as a voxel-wise classification problem expressed as:

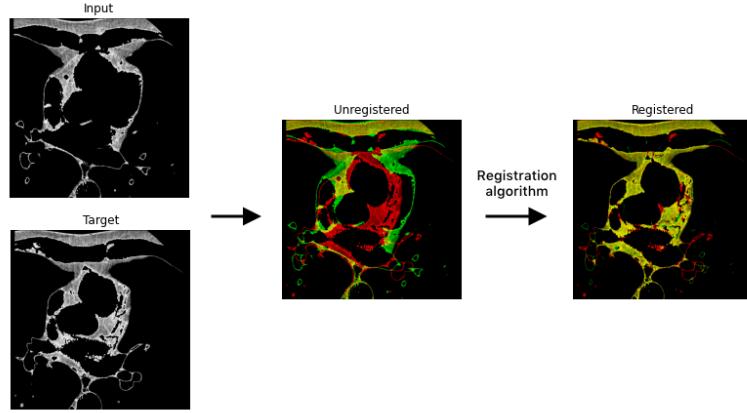


Figure 2.5: A schematic representation of the registration procedure. Initially, input and target images are chosen. Throughout the registration phase, the input image (shown in green) undergoes deformation to align with the fixed target image (shown in red). [4]

$$H(x; \theta) = (\Pr(C_0 | x), \Pr(C_1 | x), \dots, \Pr(C_{K-1} | x)), \quad (2.4)$$

where H is a classification function parameterized by θ , which maps an input vector x of voxel-level features to the probability that the voxel contains the class C_i . The features are manually constructed based on domain knowledge and represent each pixel and its surrounding region. Commonly, these include the pixel's intensity, mean intensity of the area, image moments, and other information deemed relevant for the classification. The classifier is trained to minimize a predefined loss function by feeding each pixel's features to the classifier and comparing the output to the ground truth output. This is presented visually in [Figure 2.6](#).

Another common approach is to divide the image into smaller patches and then employ a machine learning classifier to categorize each patch into one of the predefined K classes. The classifications are then fused together to form a segmentation map. Both patch- and voxel-based machine learning approaches can be more data efficient than deep learning-based approaches [4], but they often lack the ability to model complex features and dependencies.

It is clear that these traditional techniques have inspired and paved the way for more complex deep learning methods — while still playing a vital role as parts of segmentation pipelines today. In the next section, we will provide an overview of how deep learning builds on traditional machine learning to achieve more powerful models that can segment very complex regions.

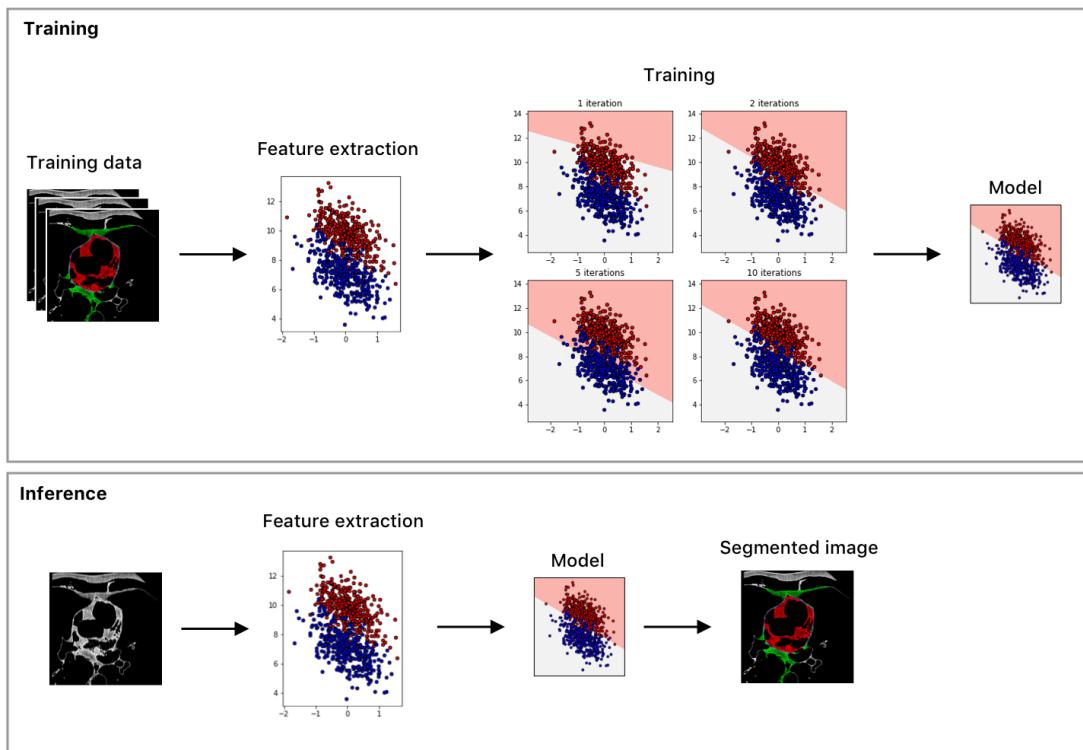


Figure 2.6: A schematic of a supervised linear classifier in a machine learning workflow. The upper section illustrates the training phase. Here, features are color-coded according to their known class from training data, depicted in red and blue. The parameters of the decision boundary, which demarcates the zones of the two classes (highlighted in light red and grey), are determined during training. The lower part of the diagram depicts the inference stage. In this phase, features are extracted from new images, and the trained model is employed to classify each pixel in the image. [4]

2.3 Deep Learning-Based Segmentation Methods

Machine learning encompasses a broad range of techniques for building statistical models, which includes deep learning — a subset of machine learning that focuses on using artificial neural networks. Artificial neural networks consist of simple nodes called **neurons**. Each neuron is a non-linear function of the sum of its inputs. This non-linear function is called the **activation function**. The neurons are all arranged in a graph where the outputs of one set of neurons are connected to the input of another set of neurons. Each connection has an associated **weight** and **bias** parameter. The value on the connection is multiplied by the weight and the bias is added before it's fed into the subsequent neuron. The exact configuration of the graph, including the number of neurons and how they are connected, is determined by hand and is referred to as the **neural network architecture**. Typically neurons are arranged into **layers**, where neurons of one layer connect only to the next layer. Current deep-learning networks usually have between ten and 100 layers. A typical neural network architecture can be seen in [Figure 2.7](#).

Mathematically, a neuron is a function that maps some input vector to a scalar value. First, the dot product between the input vector $x = [x_0 \ x_1 \ \dots \ x_{n-1}]$ and the weights column vector w is calculated and added with the bias vector b :

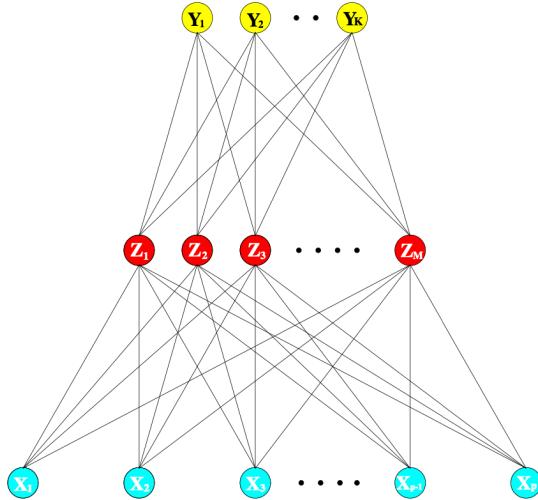


Figure 2.7: A small neural network with three layers, wherein each layer is connected to every neuron in the next layer. [16]

$$z(x; w, b) = [x_0 \ x_1 \ \cdots \ x_{n-1}] \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_{n-1} \end{bmatrix} + [b_0 \ b_1 \ \cdots \ b_{n-1}] \quad (2.5)$$

Then, all of the inputs are summed and the neuron output is produced using the activation function af :

$$f(x; w, b) = af \left(\sum_{i=0}^{n-1} z(x; w, b)_i \right). \quad (2.6)$$

The parameters b and w associated with each neuron can be consolidated into a large matrix, denoted as θ . This matrix is called the **parameters of the network**. Current deep learning networks usually have multiple millions of parameters.

One example of an activation function is the rectified linear unit (ReLU) function:

$$\text{ReLU}(x) = \begin{cases} x, & \text{if } x > 0, \\ 0, & \text{otherwise.} \end{cases} \quad (2.7)$$

ReLU can be seen as a simple thresholding function that sets all negative values to zero. Despite its simplicity, ReLU is one of the most prevalent activation functions in neural networks.

2.3.1 Neural Network Training, Validation and Testing

Initially, the parameters of a neural network are often initialized with random values. These values are then fine-tuned during a procedure called **training**, which iteratively adjusts the parameters to minimize a **loss function**. The loss function measures how well the network is performing its task by comparing its output to known correct values. For example, in image segmentation, the loss function would measure the similarity between the network's segmentation map and a hand-labeled counterpart. Through a process called **backpropagation**, each parameter is updated in the direction that will decrease the loss. This is repeated multiple times for each image in a **training dataset**.

However, there's a caveat. The training dataset, much like any statistical data, is a sampling of the "real world" data from some unknown distribution. Depending on the size and quality of the dataset, it is possible that the training dataset does not represent the true distribution of the data. Moreover, the network can learn to produce correct solutions for each training image individually instead of learning general patterns in the data. This is called **overfitting** — the model excels on the training data but fails on unseen data. To detect and mitigate this, two additional datasets are introduced: the **validation** and **test** datasets.

The **validation dataset** serves a critical role during model development. After training, the model is assessed using this dataset to inform decisions about its architecture, preprocessing, and other facets. During the development cycle, the network is constantly modified and re-evaluated on the validation dataset. Since adjustments to the model are based on its performance on the validation dataset, there's a risk of inadvertently tailoring the model to the specific distribution of the validation data. To counteract this potential bias, a **test dataset** (often termed a hold-out dataset) is used. This dataset is reserved exclusively for a final evaluation, offering the least biased estimate of the model's real-world performance.

The training, validation, and testing datasets are created before the model development process. Usually, they are sampled randomly from a larger dataset with ratios such as 80%, 10% and 10% for the training, validation, and testing datasets, respectively.

In the next chapter, we will further discuss overfitting in neural networks. For now, we will move on to describing neural network architectures in more detail.

2.3.2 Encoders and Decoders

Machine and deep learning-based segmentation methods can be conceptualized as a two-stage process. The first stage, termed the **encoder** stage, is a function $En : \mathbb{R}^{D \times C} \rightarrow \mathbb{R}^{n \times m \times d}$. In other words, a D -dimensional image with C channels is mapped to a feature vector. This vector has an arbitrary size and depth determined by the network architecture. The encoder effectively compresses the image by distilling salient information into a smaller representation, the feature vector. In deep learning parlance sometimes this is referred to as the **backbone** of the network.

The output of the encoder then goes to the second stage called the **decoder** or the **head**. The decoder is a function $De : \mathbb{R}^{n \times m \times d} \rightarrow \mathbb{R}^{D \times K}$ that maps a given feature vector into a segmentation map.

Given an image $I(A)$, the segmentation process can then be written as:

$$M(A) = (En(\theta_{En}) \circ De(\theta_{De}))(I(A)), \quad (2.8)$$

where En and De are functions parameterized by a θ_{En} and θ_{De} , respectively. In conventional machine learning, En is not a trainable function. Instead, θ_{En} consists of hand-selected parameters to extract pre-selected features from the image. The value of θ_{De} , on the other hand, is determined by minimizing a loss function. In deep learning, both θ_{En} and θ_{De} are determined by minimizing a loss function.

Thus, the difference between traditional machine learning and deep learning is in the encoder stage. Deep learning utilizes neural networks to define features, whereas traditional machine learning uses handcrafted functions.

The separation of segmentation models into encoder and decoder stages is a crucial aspect of current research in deep learning. This separation allows for independent improvement of both stages and an easy combination of different encoders and decoders.

For instance, a classification and segmentation network could use the same encoder architecture. While the classifier would utilize a simpler decoder to map input features to a class probability vector, the segmentation model would employ a more complex decoder to produce a segmentation map. However, both models would use the exact same encoder. In fact, there are models such as Mask R-CNN [17] that use two parallel decoders, one for segmentation and another for object detection, both connected to the same encoder. Throughout the rest of this thesis, when describing neural network architectures, we will describe them in terms of their encoder and decoder and how the two are connected.

2.3.3 Convolutional Neural Networks

Convolutional neural networks have profoundly impacted computer vision, establishing deep learning as the prevailing approach for complex computer vision tasks. The foundational element of a CNN is its convolutional layer, which uses a convolution-like operation in place of the conventional neurons mentioned earlier in this chapter.

Convolution

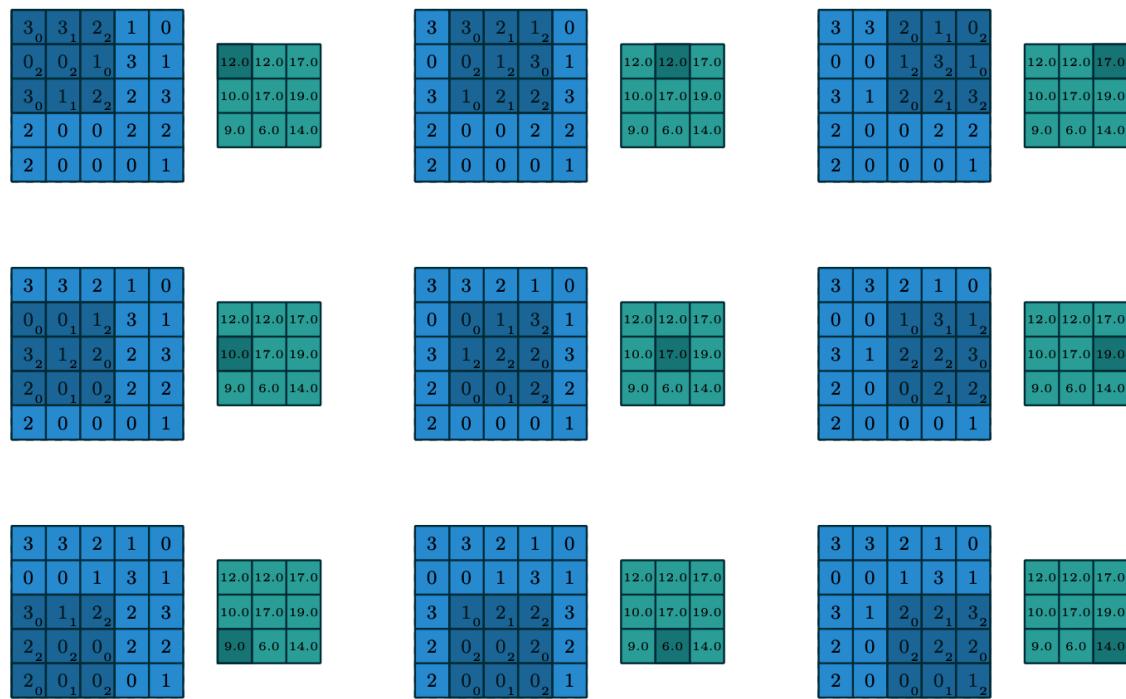


Figure 2.8: A visual depiction of a convolution procedure step-by-step. In each step, the kernel slides over the image (shown in blue). The overlapping elements between the kernel and the image are multiplied and then summed to produce a value in the resultant image (shown in green). [18]

Generally, convolution is a mathematical operation between two functions. In the context of this thesis, however, we will focus on discrete 2D convolution between two square images, as that is most relevant for image segmentation. Convolution is denoted as $I(A) \star k(B)$ where $I(A)$ is an image $I(A) \in \mathbb{R}^{W \times H}$ and $k(B)$ is a matrix $k(B) \in \mathbb{R}^{w \times h}$ indexed by locations $B \in \mathbb{N}^2$ for called the **convolutional kernel**. At position $A = (a_x, a_y)$, the convolution operation is defined as:

$$(I \star k)(a_x, a_y) = \sum_{j=-\infty}^{\infty} \sum_{i=-\infty}^{\infty} I[i, j]k[a_x - i, a_y - j]. \quad (2.9)$$

Explained differently, the resulting image is produced by sliding the kernel over the input image pixel by pixel. At each pixel location, values where the kernel and the image overlap are multiplied, and all of the products are summed together to form the corresponding pixel's value in the output image. This is shown visually in [Figure 2.8](#).

While mathematically a simple operation, convolution is exceedingly powerful and can produce almost endless transformations of an image. It's most commonly used for filtering — a convolution can elegantly find patterns in the image and increase their intensity in the image. One such example is the convolution with a kernel called the Prewitt operator:

$$I_y(A) = I(A) \star \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} \quad (2.10)$$

When convolved with this kernel, the resulting image has high-intensity pixels in regions where vertical edges are present, and low intensity everywhere else. This can be seen in [Figure 2.9](#).

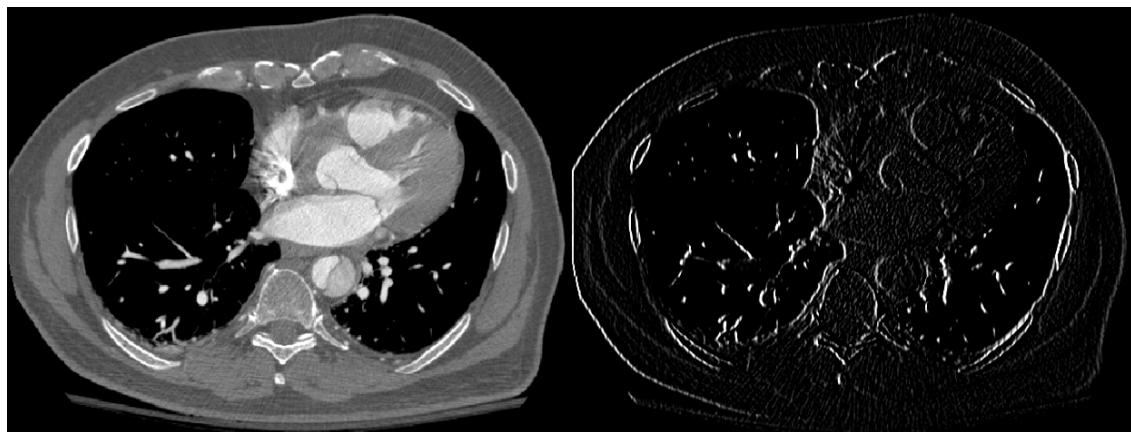


Figure 2.9: An example of an input image (left) convolved with the Prewitt operator (right). Note that vertical edges are accentuated in the convolution result.

Vertical edges necessarily have to have a large jump in values going from left to right or right to left. Otherwise, there would be no perceptible edge. This kernel takes advantage of that fact to accentuate parts of the image where there is such a jump. It does this by replacing each pixel with the difference between the pixels on its left and its right.

This process happens as follows. For each pixel of the input image, the kernel is placed such that it is centered on that pixel. This means that the values of the pixel as well as its neighbors above and below are all multiplied by zero. The neighbors on the left are multiplied by -1, and the ones on the right are multiplied by 1. Summed together, the result represents the sum of the values on the right of the pixel, minus the sum of the values on the left.

To illustrate this, let us consider 1×3 region of the image where no vertical edges are present:

$$\begin{bmatrix} 128 & 130 & 136 \end{bmatrix} * \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} = \sum_{i,j} \begin{bmatrix} -1 \times 0 + 0 \times 128 + 1 \times 130 \\ -1 \times 128 + 0 \times 130 + 1 \times 136 \\ -1 \times 130 + 0 \times 136 + 1 \times 0 \end{bmatrix} = 8$$

This section of the image does not contain a vertical edge, so the convolution result is a relatively low value. In a standard image with values in $[0, 255]$, 8 would appear almost completely black.

However, consider some section of the image where a vertical edge is indeed present:

$$\begin{bmatrix} 63 & 66 & 132 \end{bmatrix} * \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} = \sum_{i,j} \begin{bmatrix} -1 \times 0 + 0 \times 64 + 1 \times 66 \\ -1 \times 64 + 0 \times 66 + 1 \times 132 \\ -1 \times 66 + 0 \times 132 + 1 \times 0 \end{bmatrix} = 68$$

The value is now much larger due to the difference between the left and right sides of the image. This example demonstrates how a relatively simple kernel can capture complex features of an image.

Beyond edge detection, there are many commonly used convolution kernels to perform tasks such as blurring, sharpening, or denoising images. A CNN can leverage the power of the convolution by stringing together sequences of intricate kernels to match complex patterns in the image.

Convolutional Layers

A CNN operates by passing an image through a sequence of convolutions. The output from one convolution serves as the input for the next. Interspersed between these convolutions are non-linear transformations of the outputs. This layered approach enables the network to successively find more and more intricate features by combining simpler ones. For instance, combining edges into corners, corners into shapes, and ultimately detecting objects from shapes. The addition of non-linearity drastically increases the complexity of features the network is able to detect and combine.

More technically, in a CNN, convolutional layers are connected to one another in a similar fashion to neurons in a regular neural network. Each layer performs a set number of convolutions, denoted by n . The results are stored in an n -channeled **feature map**. Each channel corresponds to a convolution with a distinct kernel. Hence, a convolutional layer contains n unique kernels. The values of the kernels are the parameters of the layer, optimized during training to minimize a loss function like weights and biases of standard neurons. Finally, a convolutional layer also applies a non-linear activation function to the resulting feature map.

It's worth noting that the operation within a convolutional layer differs from a standard 2D discrete convolution. Here, both the kernel and the input are three-dimensional and have the same depth, but the kernel is usually much smaller in width and height. Since the kernel has the same width as the input, during the operation the kernel only slides along the width and height dimensions. The output of each sliding window step is still a single scalar since all overlapping elements in the sliding window volume are multiplied and summed together. This means that the output of the whole operation is a two-dimensional image. This is visualized in [Figure 2.10](#).

Pooling Layers

Most CNN-based encoders follow a pattern of gradually reducing the width and height of the feature maps while increasing their depth. Depth is increased by giving each

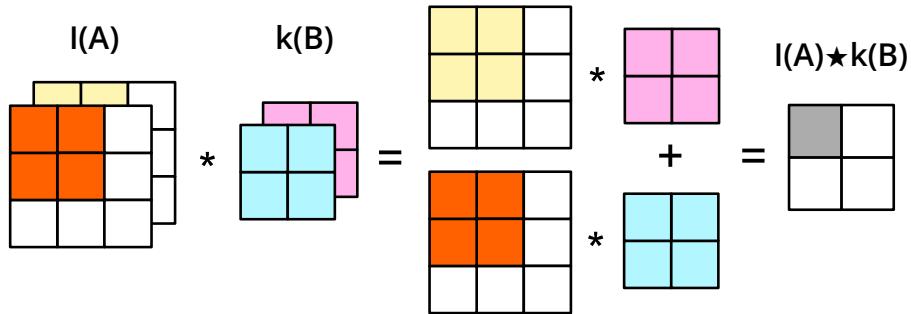


Figure 2.10: A view of one step of a single convolution operation inside a CNN layer. The layer performs multiple convolutions, each with a different kernel that has an equal number of channels as the input image. In each step, the whole kernel slides over the width and height of the image, and the overlapping channels are multiplied together and summed to produce a single output value. The output of the convolution is one channel of a n -channel image where n is the number of different kernels in the layer.

successive convolutional layer more kernels. Since each channel represents the result of a convolution with a specific kernel, having more kernels directly increases the output depth. To reduce the width and height, **pooling layers** are employed. These layers downsample the image, often by averaging pixel values within a defined neighborhood. This downsampling serves dual purposes in a CNN. Firstly, it compresses the image by gradually removing spatial information and retaining relevant semantic information. Secondly, the network is gradually able to build up large complex features by combining small general features. Such an architecture is shown in [Figure 2.11](#).

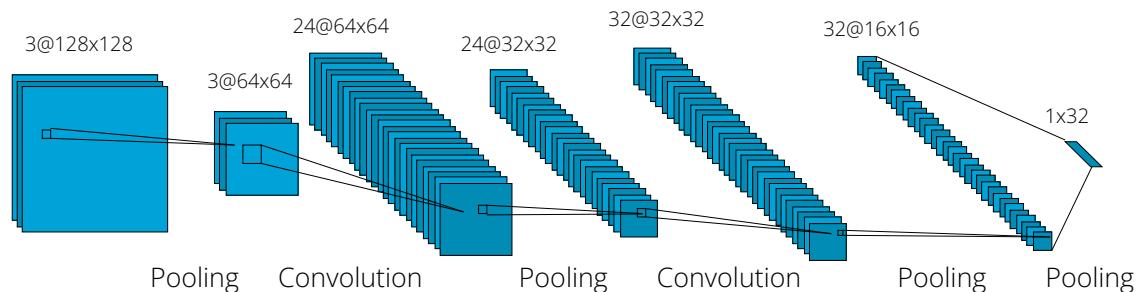


Figure 2.11: A typical architecture of a CNN encoder. The encoder consists of consecutive convolutional and pooling layers that gradually increase the feature map depth and decrease its width and height. The result is a map of features that tells the decoder what features are on the image but does not provide much spatial information about the location of those features. [19]

To illustrate this, consider a contrived example of a CNN encoder trained to detect grapes on a vine. The first few layers might focus on simple operations such as edge detection, requiring a large amount of spatial information. The next layer could use the edge information to detect circles. This layer doesn't need as much spatial detail since the edges are already identified. A third layer might then combine the locations of the circles to detect a cluster of grapes. By this stage, the spatial information is even less important, with the overall arrangement of circles being the primary focus.

A convolutional decoder essentially reverses the encoder process. Instead of adding channels, a convolutional decoder successively removes channels up to and upsamples

the image. By the end, the output matches the input image's width and height, and the number of output channels equals the number of classes.

One way in which upsampling is implemented in a convolutional decoder is with a **transposed convolutional layer**. These layers are very common in segmentation neural networks and allow for dynamic, learned upsampling of the image. Regular convolutions cannot increase the image dimension width and height. To overcome this, a transposed convolution produces an output matrix for each step of the sliding window process, instead of just a scalar value. For instance, a transposed convolution with a 2×2 kernel produces 4 new values for each pixel in the input image. This is done by performing scalar multiplication of the whole kernel and the pixel it is sliding over during each sliding window step. All of the sliding window results are joined as a grid to form a final, larger, image.

To better understand this, consider an example of upscaling a 2×2 2D input image m with a 2×2 kernel k . The transposed convolution can be calculated as follows:

$$m \star^T k = \begin{bmatrix} m_{11}k_{11} & m_{11}k_{12} & 0 \\ m_{11}k_{21} & m_{11}k_{22} & 0 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & m_{12}k_{11} & m_{12}k_{12} \\ 0 & m_{12}k_{21} & m_{12}k_{22} \\ 0 & 0 & 0 \end{bmatrix} + \\ \begin{bmatrix} 0 & 0 & 0 \\ m_{21}k_{11} & m_{21}k_{12} & 0 \\ m_{21}k_{21} & m_{21}k_{22} & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & m_{22}k_{11} & m_{22}k_{12} \\ 0 & m_{22}k_{21} & m_{22}k_{22} \end{bmatrix}$$

This method can be similarly applied to images and kernels of any size. While transposed convolution is often referred to as *deconvolution*, in this thesis, we will avoid using the term to prevent confusion with the meaning of deconvolution outside of the context of deep learning.

Segmentation models that use only convolutional and pooling layers are called **fully convolutional models** and are one of the most widely used types of models for medical image segmentation. In the following section, we will describe these and other commonly used CNN architectures in medical image segmentation.

2.4 CNN Architectures for Medical Image Segmentation

Deep learning's application in medical image segmentation has predominantly seen the use of convolutional neural networks. This field has witnessed a rapid evolution, giving rise to a diverse range of architectures. While some of these architectures were initially designed for general image classification and segmentation tasks, others have been specifically tailored to overcome challenges specific to medical imaging. In this section, we will highlight the most influential CNN architectures that have shaped medical image segmentation, presented in their chronological development.

2.4.1 Fully Convolutional Network (FCN)

The fully convolutional network (FCN), one of the more straightforward contemporary CNN architectures for image segmentation, was introduced in [20]. As mentioned earlier, convolutional encoders gradually downsample the image while increasing the number of channels in the feature maps. There is then a need to reverse the process and decode the features into an image of equal width and height as the original input image. FCN achieves this using convolutional layers and upscaling.

In FCN, the output of the encoder is interpreted as a coarse heatmap of detected features. An output segmentation map needs to be constructed based on these features and have a number of channels equal to the number of classes. FCN does this by using a convolutional layer with 1×1 kernels — one kernel for each class. The resulting convolutions combine the features depth-wise to form a class prediction for each location of the image. However, the resulting segmentation map has a very low spatial resolution. To be useful, it needs to be upsampled using transposed convolution to the original image resolution.

As noted earlier, the encoder's output feature map is only a rudimentary representation of spatial details. To inject additional spatial information into the final segmentation, FCN combines predictions from multiple encoder layers to form the final prediction. This is done by upsampling 1×1 convolution results from different encoder layers to a uniform size and summing them together. The composite prediction is then upsampled to the original image resolution. This process is shown in [Figure 2.12](#).

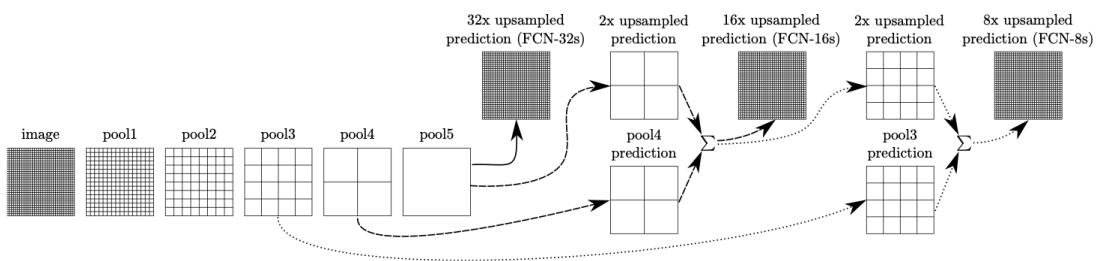


Figure 2.12: A diagram of how FCN forms predictions based on the output of different encoder layers. Encoder layers are shown on the left and the grid represents the coarseness of the feature map. The maps are combined at three different levels to produce three predictions. Each prediction is compared with the ground truth during training, but for inference only the 8x upsampled prediction is used. [20]

2.4.2 U-Net and Its Variants

U-Net [21] is one of the most influential and commonly used architectures in medical image segmentation. Aside from being a standard baseline model, a well-tuned U-Net network with data preprocessing has reached state-of-the-art performance in various segmentation challenges across different domains and modalities [22].

At its core, U-Net follows a straightforward architectural pattern, comprised of a fully convolutional encoder and decoder with only convolutional and pooling layers. The encoder and decoder have a symmetrical structure: the encoder gradually decreases the width and height of the feature maps while increasing the number of channels, while the decoder does the opposite using upsampling or transposed convolutions. While FCN uses only one 1×1 convolutional layer and upsampling to form a prediction, U-Net uses a sequence of several convolutional and upsampling layers to gradually build a segmentation map. This progressive structure allows U-Net to produce more detailed feature maps.

Earlier we saw that spatial information is important for the precise localization of the features. Like FCN, U-Net also uses the outputs of different encoder layers to maintain spatial information in the decoder but does so in a different way. In U-Net, the output of every encoder layer is added to the input of its corresponding decoder layer on the other side of the network. This gives each decoder layer concurrent access to information

about which features are on the image and where they are located. U-Net's architecture can be seen in [Figure 2.13](#).

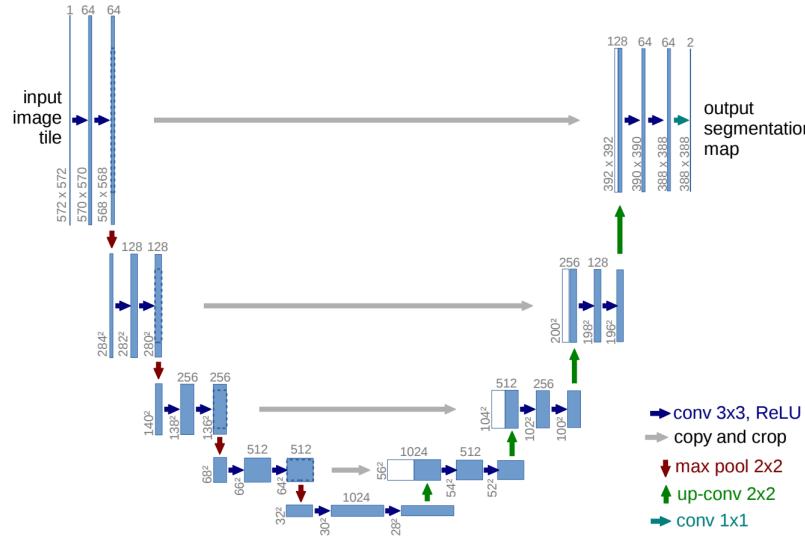


Figure 2.13: A diagram of the U-Net model. The output of each layer of the encoder is concatenated to the input of its corresponding layer in the decoder. [21]

Because of its symmetrical structure and skip connections, U-Net is often visualized in the shape of the letter U, giving it its name. U-Net is only one of a class of networks that are sometimes called U-shaped networks. Aside from being a state-of-the-art architecture for medical image segmentation, U-Net is also widely used as a component of influential deep learning models in other domains such as image generation [23] and registration [15].

nnU-Net

U-Net is a powerful baseline model that can be carefully tuned using various heuristics and combined with beneficial data preprocessing to achieve state-of-the-art results, even when compared to much more complex models. However, tuning the parameters of the model and preprocessing is a laborious process that requires expertise and empirical research. nnU-Net [22] aims to automate this tuning process by producing an optimal data pre- and postprocessing pipeline and a U-Net-based architecture for a given dataset. This is achieved through a set of heuristics, rules, and optimization techniques.

The nnU-Net framework begins by calculating key parameters of the data distribution, such as intensity distribution and median image size. These parameters feed into a set of predefined rules that determine the network's architecture and training strategy, including the number of layers, training parameters, and preprocessing approach.

Using the determined parameters, nnU-Net trains three models: one using 2D slices of the data, another using the entire 3D volume, and the third using a cropped 3D region of interest. During inference, outputs from these models are combined through preprocessing to form an ensemble prediction, enhancing the overall accuracy and robustness of the segmentation. This process is shown in [Figure 2.14](#).

nnU-Net has demonstrated exceptional performance, establishing itself as a state-of-the-art solution in a variety of CT and MRI segmentation tasks. nnU-Net is a hugely

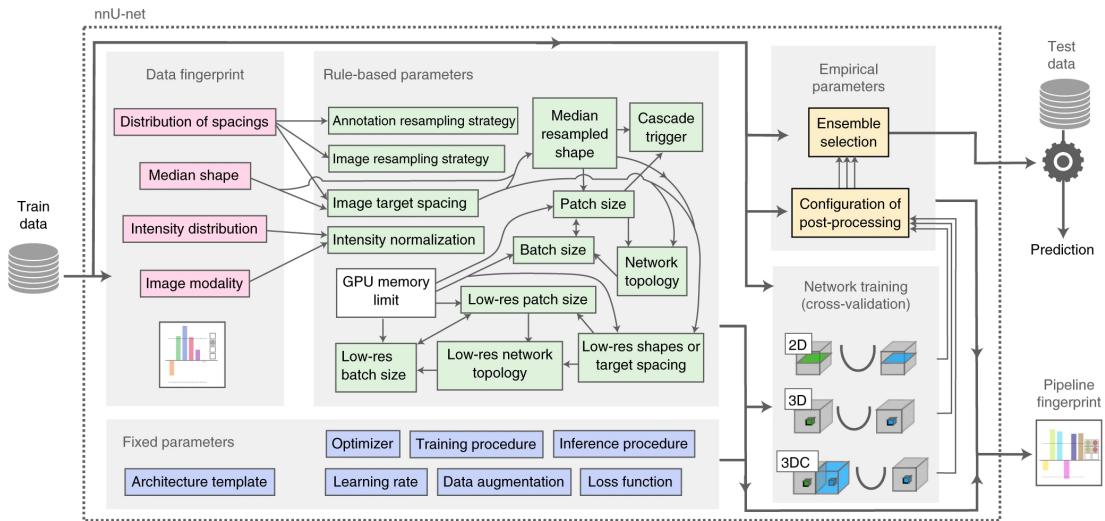


Figure 2.14: A diagram of the nnU-Net procedure of creating a training configuration. [22]

important tool for both academia and industry, as it allows for creating well-tuned baseline models with very little manual intervention or technical expertise.

Furthermore, nnU-Net exemplifies the critical role of preprocessing and parameter tuning in improving model performance, particularly when working with small datasets. This is the cornerstone of this thesis: improving data efficiency through preprocessing. The success of nnU-Net serves as a validation of this approach.

U-Net++

In standard U-Net, skip connections link encoder layer outputs to decoder layer inputs through basic concatenation operations. This approach leaves room for more intricate integration of the decoder and encoder outputs. Additionally, U-Net only merges features at corresponding levels of the encoder and decoder. An enhancement to this could be to allow each decoder layer access to features from multiple encoder layers. U-Net++ [24] takes advantage of these opportunities with a flexible architectural design.

In U-Net++, skip connections are reimaged as a network of convolutional layers. This setup connects each encoder layer to its respective decoder layer through multiple convolutional layers. This enhances the information fed into the skip convolutional layers. Additionally, the outputs of these skip connections are further upsampled and fed into the subsequent decoder layers. This design results in a significantly more complex network architecture. Here, each decoder layer not only receives features from its direct encoder counterpart but also from all deeper encoder layers. This multifaceted connection system is visually represented in Figure 2.15.

In U-Net++, skip connections are reimaged as a network of convolutional layers instead of simple concatenation. Moreover, outputs from each encoder layer (except the first) are upsampled and then supplied as additional inputs to the skip connections of the preceding layer. The outputs of the skip connections themselves are also upsampled and given to the next corresponding decoder layer. Thus, each decoder layer not only receives features from its direct encoder counterpart but also from all deeper encoder layers. This architecture is presented in Figure 2.15.

There are several advantages to U-Net++ over U-Net. Each decoder layer has access to multi-scale features of the input image and thus can produce more accurate outputs.

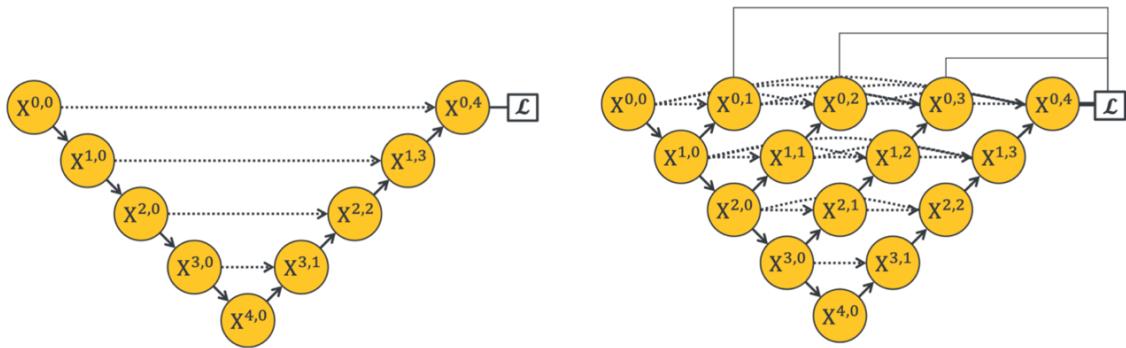


Figure 2.15: A comparison between U-Net (left) and U-Net++ (right). Each node in the graph represents a convolutional layer. The dashed arrows represent skip connections, while full arrows are downsampling or upsampling operations. [24]

Conceptually, U-Net++ can be viewed as an ensemble of U-Nets with varying depths, which helps in faster convergence compared to a standard U-Net of similar depth. Further, instead of finding the optimal U-Net depth, due to the flexibility of neural networks, U-Net++ learns the optimal depth during training. This makes U-Net++ more flexible, requiring fewer design decisions. However, U-Net++ has a higher parameter count and does not necessarily enhance data efficiency. In some cases, its performance benefits can be matched or surpassed by a well-preprocessed and tuned U-Net [22].

2.4.3 Mask R-CNN

Earlier in the chapter we mentioned that CNN-based segmentation networks generally consist of encoder and decoder stages. Segmentation architectures such as U-Net employ a decoder that progressively upsamples the feature maps to construct a segmentation map. This differs from how object detection decoders work.

Object detection aims to find the location of an object in an image. Instead of classifying each pixel as in segmentation, object detection typically outputs the coordinates of a bounding box surrounding the object. In technical terms, the output is not an image but rather a vector of four numbers for each object. As such, the decoders in object detection networks are usually much shallower and employ fully connected layers instead of convolutional layers.

Despite these differences, both object detection and segmentation networks share similar encoder designs. Mask R-CNN [17] takes advantage of this to combine object detection and segmentation into a single network. In Mask R-CNN, a standard CNN encoder generates a feature map from an input image. This map is then given to two parallel decoders: one for segmentation and the other for object detection. The segmentation decoder, akin to those in standard segmentation networks, uses convolutional layers and upsampling to generate a segmentation map for each object. The detection decoder, in contrast, employs downsampling and fully connected layers to output the bounding box and class label of each object. The Mask R-CNN architecture can be seen in Figure ??.

Consequently, for every detected object, Mask R-CNN provides a bounding box, a class label, and a segmentation mask. The network is trained using a composite loss function, optimizing for segmentation, detection, and classification simultaneously.

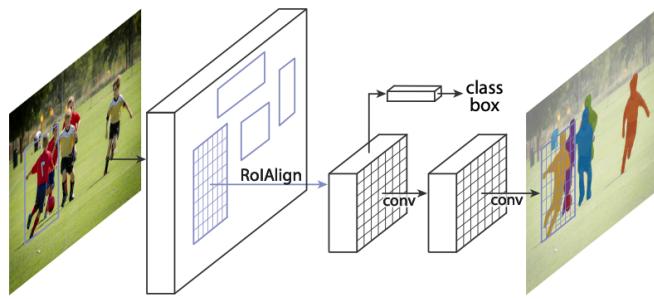


Figure 2.16: A diagram of the Mask R-CNN architecture. Two parallel decoder branches are used to achieve segmentation and object detection simultaneously. [17]

The integration of these tasks in Mask R-CNN takes advantage of synergies between them. The information learned for classifying an object is relevant to its detection. Similarly, by using a bounding box the segmentation decoder can focus on the relevant portion of the image and be invariant to the scale of the object on the image. This joint learning approach enhances performance across all tasks compared to training separate networks for each. However, Mask R-CNN’s complexity can pose challenges, particularly when trained on very small datasets, where simpler segmentation-focused architectures like U-Net might be more effective.

2.4.4 Other Notable Segmentation CNNs

A central challenge in CNN-based segmentation is balancing spatial information with feature complexity. As the network increases the feature map depth and decreases width and height, spatial information is erased but more complex features can be detected on the image. Effective segmentation requires both elements. While U-Net addresses this with skip connections, other architectures use different methods to integrate multi-scale features.

One approach involves pyramid pooling, as utilized in PSPNet [25]. Initially, PSPNet uses a standard encode to create a feature map. This map is then rescaled to various sizes, and each rescaled map undergoes processing through a series of convolutional layers. The resulting outputs from these layers are resized to a uniform spatial dimension and concatenated. This combined feature map is then fed into a decoder to produce the final segmentation map.

Pyramid pooling is also employed by DeepLabV3 [26] albeit in a different way. Like in PSPNet, the encoder first creates a feature map. This feature map then goes through a series of parallel layers that each use an atrous convolution. In atrous convolution, the kernel is dilated by introducing gaps between its values. During each sliding window step, the gaps are ignored and the result is convolved regularly for non-empty kernel values. This technique enables covering a larger area of the image without increasing kernel parameters. Each parallel layer uses different gap sizes in the kernel, creating multi-scale features. These features are upscaled, merged, and further decoded into a segmentation map.

MA-Net [27] takes a different route by incorporating an attention mechanism for combining multi-scale features. Rather than directly concatenating in skip connections like U-Net, MA-Net introduces a multi-scale attention block. This block aims to capture the interdependencies of different features to focus on more salient ones before they

are concatenated to the decoder layer input. Additionally, the final decoder layer is processed through a position-wise attention block to highlight significant features in specific image locations. MA-Net has shown performance improvements in liver tumor segmentation on CT images compared to U-Net and U-Net++.

2.5 Fully Connected Transformers for Medical Image Segmentation

So far most of this chapter dealt with CNNs. However, another class of networks — transformers [28] — has recently gained traction and even surpassed the segmentation quality of CNNs for many tasks. The key to the transformer’s impressive performance is the ability to support very deep neural networks. Empirical observations suggest that with enough data, network performance tends to increase with the depth of the network. Hence, transformers have become state-of-the-art in tasks where a large amount of data is available. Notably, current transformers used in natural language processing contain more than 100 layers and trillions of parameters, orders of magnitude more than CNN-based models.

Transformers first became broadly used in the field of natural language processing and achieved groundbreaking results in text generation, machine translation, and understanding natural language. Broadly, transformers employ a combination of regular (non-convolutional) fully connected layers and attention mechanisms, which are the key to their success. Attention allows a layer to learn the interdependencies between parts of its input or across different input values. In linguistic contexts, this leads to a more nuanced understanding of syntactic relationships, enabling the network to learn concepts such as the interplay between subjects and objects or the likelihood of certain adjectives preceding specific nouns.

2.5.1 Attention

For each element x_i of an input sequence of length n , an attention mechanism aims to produce a vector $v_i = (v_{i1}, v_{i2}, \dots, v_{in})$ where v_{ij} measures how dependent x_i is on x_j . In a natural language scenario, if we consider a sentence as our input sequence, this mechanism can be envisioned as going word by word, and for each word producing relevance scores for every other word in the sentence. For instance, in the sentence “Susan has a brown dog that she loves.”, the word “she” should have a high relevance score for “Susan”, and a low score for “brown” since it is not pertinent to encoding the meaning of “she” in this sentence.

Mathematically, attention can be expressed as a series of matrix multiplications. First, each input sequence element is transformed into a query vector. This transformation represents a learned encoding by the network, converting an input element into a lower-dimension vector. The query serves to identify the current input element in the sequence. Let X be an input matrix where each row is a new element of the sequence. Let W_Q be a matrix of weights learned by the network that encodes each element of the input sequences into a vector. The query matrix Q can then be calculated as:

$$Q = W_Q X \quad (2.11)$$

Where each row Q_i represents the query for X_i . Remember that the goal of attention is to score the relevance between an input element X_i and some other element X_j of the same input sequence. Having encoded each element X_i in the form of Q_i , we need

to similarly encode each X_j in a different matrix K where K_j is called the **key** for some input sequence element X_j . We can do so as:

$$K = W_K X, \quad (2.12)$$

where W_K is another weight matrix that the model optimizes during training. Given K and Q , we can compute the relevance scores for every combination of input elements x_i and x_j as:

$$A = QK^T, \quad (2.13)$$

Here, A_{ij} represents the relevance score between x_j and x_i . It is important to note that this is only one of many ways the matrix A can be calculated. In practice, A is scaled to be in $[0, 1]$ or transformed in some other way.

Finally, to get the output Z of the self-attention layer, we weigh each element of the original input sequence X by the attention scores A :

$$Z = AV, \quad (2.14)$$

Z consists of X elements weighted by how relevant they are to the encoding of the whole sequence. Visually, this process for one element of the input sequence can be seen in [Figure 2.17](#).

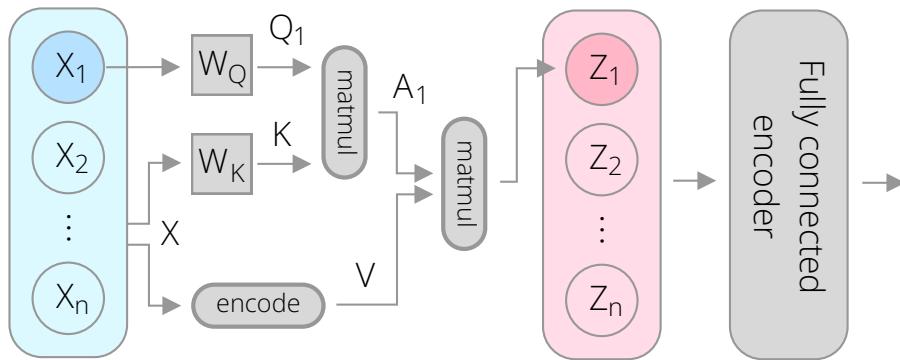


Figure 2.17: A visualization of a single encoder layer in a transformer network. This shows the encoding process for one element of the input sequence.

Encoder and decoder layers in a transformer include an attention layer followed by a fully connected layer. In the encoder, the attention layers are generally *self-attention*, where both the key and the query come from the input sequence. In contrast, the decoder layers employ both self-attention and **encoder-decoder attention**. In the latter, the keys and values come from the encoded features of the last layer of the encoder. Finally, each attention layer performs Z calculation multiple times using different W and K weights, and the weighted values are concatenated together. This design enables the network to simultaneously learn and apply various types of attention within each individual layer.

2.5.2 Positional Encoding

The position of each word in a sentence carries important information. The sentence “The yellow notebook is on top of the red notebook.” has a different meaning than “The red notebook is on top of the yellow notebook.”, even though both sentences contain the same words. To capture this in transformers, the position of each input element is encoded and included as part of the input.

Each input sequence element index is mapped to a matrix via a function $f : \mathbb{Z}^+ \rightarrow \mathbb{R}^d$ where d is the dimensionality of the layer's input. Mapped that way, the positional encoding can simply be concatenated to the layer's input with no additional changes in the network. While it's possible for this encoding function to be represented as neural network layers and learned during training, a common approach is to use a predefined, fixed function. For instance, in citeattnAllYouNeed, given an element index pos , the i -th element of the positional encoding is calculated as:

$$PE(pos, i) = \begin{cases} \sin(pos/1000^{2i/d}), & i \text{ is even} \\ \cos(pos/1000^{2i/d}), & i \text{ is odd.} \end{cases} \quad (2.15)$$

The choice of positional encoding functions is guided by their numerical characteristics, ensuring they are efficiently processed by the neural network.

2.5.3 Adapting Transformers to Image Segmentation

Convolutional Neural Networks (CNNs) have traditionally dominated image processing due to their efficiency in handling large inputs. Consider a 128×128 pixel image which, when flattened into a row vector, results in 16,384 elements. In a fully connected network, where each input element connects to every output element, this translates into $16,384^2$ parameters for just one layer. Even if the image were to be heavily down-scaled and ignore the loss of information this would create, managing such a network on current GPUs is infeasible.

Transformers that use fully connected layers face a challenge when applied to image data due to this very issue. The key solution, as utilized in the Vision Transformer (ViT) [29], is treating an image as a sequence of small patches, allowing the image to be processed similarly to text.

ViT, primarily designed for image classification, begins by partitioning each image into uniform square patches (e.g. 32×32 pixels). The patches, ordered starting from the top-left of the image, form an input sequence X . Each patch is then flattened to a row vector and given to a fully connected layer to encode it as a d -dimensional vector. Then, a d -dimensional positional encoding is added to the encoding as a function of the patch index. The resultant matrix serves as the input for the encoder, which operates identically to a standard transformer network. This architecture can be seen visually in Figure 2.18.

While ViT originally used a classification decoder, it is possible to combine the ViT encoder with a segmentation decoder to obtain a segmentation map. There are several ways this has been achieved. For instance, the Segmentation Transformer (SETR) [30] introduces two distinct CNN-based segmentation decoders connected to a ViT encoder. The first variant follows a conventional convolutional decoder design, where the encoder's output goes through a series of upsampling and convolutional layers, producing a segmentation map in the original image resolution. The second decoder variant is similar to the way FCN works — outputs from various encoder layers are resized to a uniform resolution, merged, and then decoded into a segmentation map.

Later transformer-based segmentation networks modified the ViT encoder as well. For instance, the authors of the Swin Transformer [31] argue that neural networks for computer vision need to be built with translation and scale invariance in mind. Swin Transformer starts with small patches, similar to ViT, but in each successive encoder layer neighboring patches are merged. Rather than computing self-attention within individual patches, Swin uses non-overlapping windows of multiple patches to perform self-attention. This ensures that Swin's computational demand increases linearly with

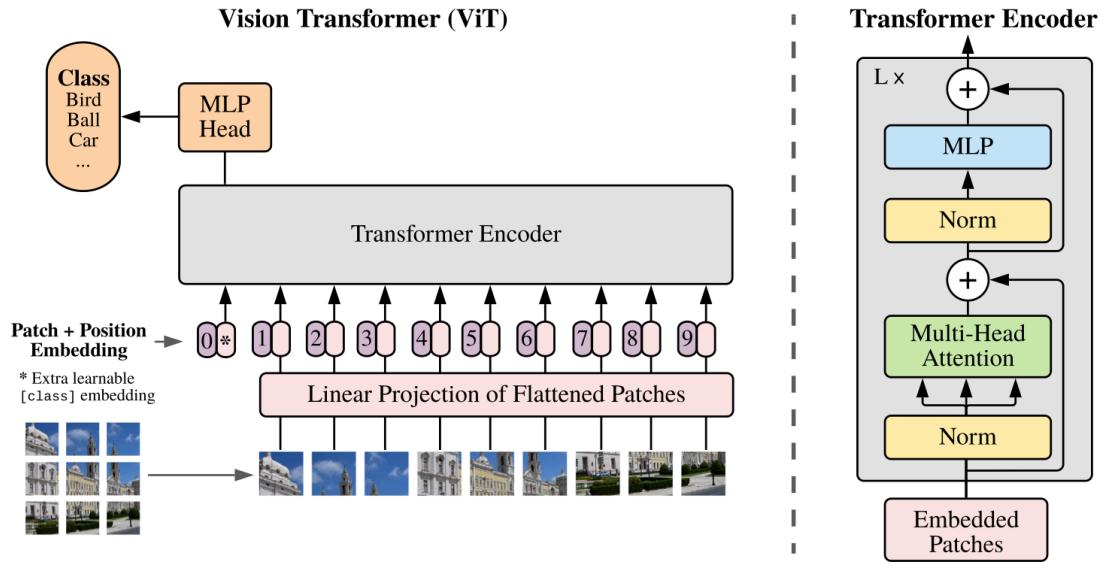


Figure 2.18: The ViT architecture for image classification. [29]

image size, a significant improvement over ViT's exponential growth. A key aspect of Swin is that each successive layer shifts the window by half of the window size. This allows self-attention between two different windows and more global information about the image. Finally, Swin's encoder output is fed into a DeepLab-like decoder to produce the segmentation map.

Transformers have demonstrated exceptional performance in image segmentation, particularly in large datasets and benchmarks with upwards of 20,000 natural images [30]–[32]. However, their advantage is less pronounced in smaller datasets. This aligns with observations from [28] — convolutional layers have properties that are naturally a good fit for extracting information out of images. Transformers need to compensate for their lack of convolutions by using more layers. However, as we will see in the next chapter, there is a relationship between the number of parameters of a network and the number of samples required to train that network.

3 Data Efficiency in Neural Network-Based Image Segmentation

Above all, the accuracy of neural networks is governed by three factors: problem complexity, model complexity, and sample size. Each of these factors contributes to accuracy in different ways and they have complex interdependent relationships. Grasping the concept of data efficiency necessitates an understanding of these three aspects. In this chapter, we will elucidate how the three factors impact model error by examining neural networks as function approximators.

Consider a machine learning scenario where we estimate the 10-year risk of heart disease based on age. Here, some underlying function maps age to heart disease risk. However, there is no way to know what this function is. Our only recourse is to collect a set of data points and try to approximate the underlying function. Our goal is to iteratively search for a function that minimizes some criterion that measures the discrepancy between the function and our collected data points.

One can construct an infinite number of functions that fit a set of data points equally well under a selected criterion. Therefore, to derive a single solution, we have to choose a finite class of functions to evaluate. One of the simplest classes of functions we could use is lines — leading to a simple linear regression. A line is defined by two parameters, its slope θ_1 and intersect θ_0 :

$$f(x) = \theta_1 x + \theta_0. \quad (3.1)$$

In this equation, $f(x)$ represents our regression model. The objective is to determine the optimal value of $\theta = \{\theta_1, \theta_2\}$ that most accurately represents the collected data pairs (x_i, y_i) . This can be achieved by minimizing the mean squared distance between the line and each data point:

$$\min_{\theta} \sum_{i=1}^N (y_i - f(x_i))^2. \quad (3.2)$$

We will conduct a simulated experiment. We will construct an underlying function as a cosine function, with each point sampled with additive Gaussian noise:

$$y(x) = -(\cos(0.01\pi x) - 1) \cdot 0.4 + Z, Z \sim \mathcal{N}(0, 0.1). \quad (3.3)$$

We will generate several (x_i, y_i) pairs and apply the previously discussed linear model and criterion to fit a line through these points. In real-world conditions, the underlying function and the nature of the measurement noise remain unknown. However, for this simulation, knowing these details is useful to demonstrate the divergence between the fitted function and the underlying one. Figure 3.1 illustrates the fitted line using varying numbers of data points.

The simulation with only three points clearly demonstrates a mismatch between the approximated and underlying functions. This discrepancy is known as **overfitting**. Overfitting occurs when the model is overly influenced by the specific sample of data,

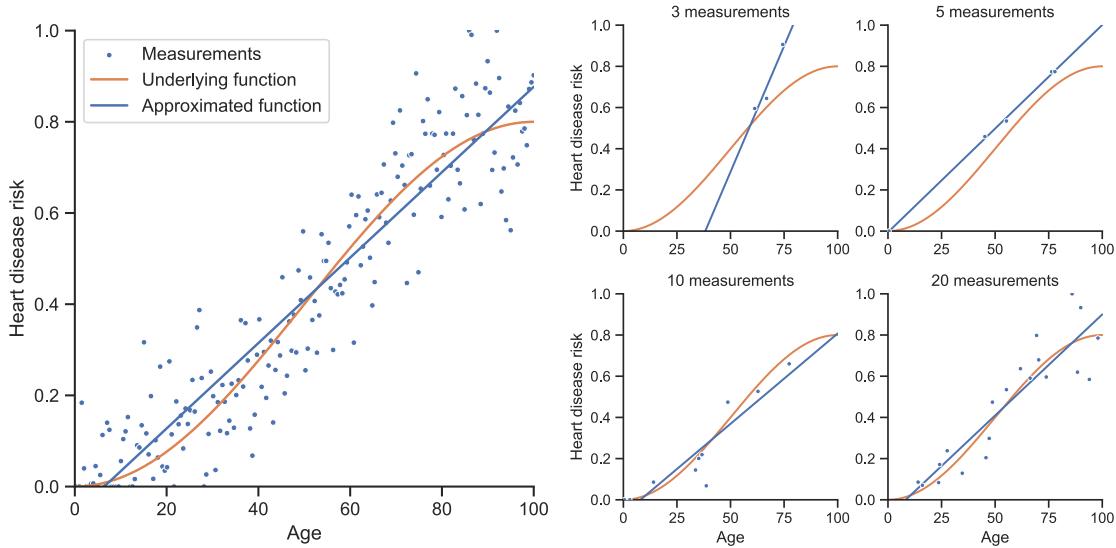


Figure 3.1: A simulated example of heart disease risk prediction using simple linear regression. The plots on the right show how the approximated function depends on the sample size.

particularly the noise, and fails to generalize to new, unseen examples. As the number of samples increases, the approximated function aligns more closely with the underlying one. However, even with 20 measurements, there is a notable difference between the approximated and underlying functions. The underlying function is a non-linear cosine function. No matter how many data points we use, a line cannot perfectly approximate this non-linear function. This is an example of **underfitting**, where the model lacks sufficient parameters to capture the complexity of the underlying function.

To address underfitting, a more complex model is required, meaning we need a broader class of functions for optimization. Our current regression can be viewed as fitting a polynomial function to the dataset. Initially, we have employed polynomials of degree 1 — lines. We can generalize our model to an n -degree polynomial:

$$f(x) = \theta_n x^n + \theta_{n-1} x^{n-1} + \cdots + \theta_1 x + \theta_0. \quad (3.4)$$

To increase our function class, we will choose a larger polynomial degree and again minimize the mean squared distance to obtain θ_n through θ_0 . By increasing the degree we are increasing the complexity of the model and the model will be able to represent more complex underlying functions. We can observe what happens when we choose degrees 2, 3, or 5 in [Figure 3.2](#).

In the case with 10 measurements, we observe that the approximated polynomial with $n = 2$ follows the underlying function somewhat well. However, as we increase the degree the function quickly begins to overfit. This experiment demonstrates that the likelihood of overfitting increases with model complexity. Yet, we need complex models to approximate complex functions.

This inherent conflict between overfitting and underfitting is known as the **bias-complexity tradeoff** (alternatively, the bias—variance tradeoff). Let us investigate this problem more formally. We can decompose the error of a deep learning model into two components [33]. Consider a trained model f_θ with optimal parameter values θ from a class of functions \mathcal{F} . Let there be a function $L(f)$ that measures the error between any model $f \in \mathcal{F}$ and the underlying function we are approximating. The total error of the

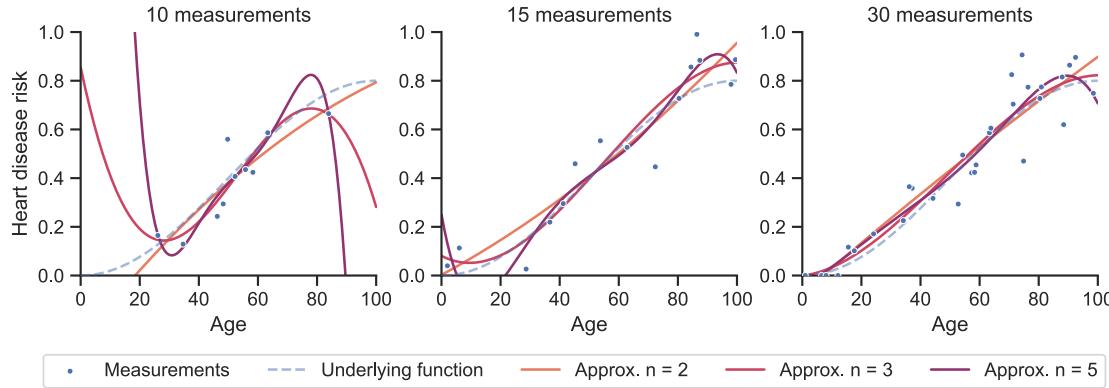


Figure 3.2: Fitting a polygonal function of various degrees on three different sample sizes.

trained model can be expressed as:

$$L(f_\theta) = \epsilon_{app} + \epsilon_{est} \text{ where: } \epsilon_{app} = \min_{f \in \mathcal{F}} L(f), \quad \epsilon_{est} = L(f_\theta) - \epsilon_{app}, \quad (3.5)$$

where ϵ_{app} is the **approximation error** and ϵ_{est} is the **estimation error**. The approximation error is the residual error when using the best possible approximator within our chosen class of functions \mathcal{F} . In other words, the error that remains given the best possible model parameters. For instance, in the earlier function approximation example with a 1st-degree polynomial, even with infinite data, there would always be some discrepancy in approximation compared to the underlying function. The approximation error depends entirely on the class of functions we use and can only be reduced by adopting a different class, such as polynomials of a higher degree. The remaining estimation error depends on the quality of our parameters θ . It is possible to have a class of functions capable of perfectly approximating the underlying function but still fail to precisely identify the optimal function within that class.

In our previous experiment, we observed that the estimation error is influenced by both the sample size and the complexity of the model. To formally analyze this, we can explore the concept of **sample complexity**. Sample complexity, denoted as $n(\epsilon, \delta)$ for chosen $\epsilon, \delta \in (0, 1)$ measures the number of samples needed for a model f chosen from a finite class of functions \mathcal{F} to be **probably approximately correct**. This means that, across different samplings of data points, there is a probability of $1 - \delta$ that the model is correct within some margin of error ϵ , i.e. $L(f) \leq \epsilon$. Assuming there exists some unknown $f' \in \mathcal{F}$ for which $L(f') = 0$, we can say that the approximation error is zero. Then, it can be shown [33] that:

$$n(\epsilon, \delta) \geq \frac{\log(|\mathcal{F}| \delta)}{\epsilon}. \quad (3.6)$$

We can see that an increased function class size (for instance, by choosing a higher-degree polynomial) necessitates a larger sample size n to probably approximately correctly approximate the underlying function. Thus, we conclude that the estimation error can be reduced either by adding more data samples or by decreasing the size of the function class we are optimizing over.

In medical image segmentation, we have an unfortunate stalemate between the three factors governing estimation and approximation errors. The complexity of medical image segmentation demands a large number of parameters, which in turn increases the

estimation error. Ideally, this error could be mitigated by adding more data samples, but in the realm of medical imaging, acquiring ample data is often impractical.

To reduce the estimation error, we are left with reducing the size of the function class. We can do so in two ways. The first involves reducing the parameters of the model itself, shrinking the class of functions, and allowing the model to more easily find the optimal function. This usually requires transforming the data such that it can be easily segmented with a less complex model, or using domain-specific knowledge to reframe the problem. The second method involves initializing the model with parameters that are already near-optimal. In this case, the model starts close to the optimal function, needing only to search within a local space around these parameters. This is typically achieved by first optimizing the model on similar data for a related task that involves learning features pertinent to segmentation, followed by further optimization for the specific segmentation task.

In the subsequent sections, we will delve into specific techniques from these two approaches, providing examples of their application in neural network-based medical image segmentation.

3.1 Transfer Learning

Consider the task of segmenting liver tumors in CT scans. Suppose we have a limited dataset of labeled liver tumor regions, but a more extensive collection of labeled liver regions. In such a scenario, the first step would be to train a neural network for liver segmentation. The rationale here is that most of the parameters learned for liver segmentation will be relevant and close to those needed for liver tumor segmentation. Therefore, when we are developing the liver tumor segmentation network, we copy the values of the parameters learned in the liver segmentation network. We can then say that the network was **pre-trained** or **initialized** on a liver segmentation dataset and **fine-tuned** on the liver tumor segmentation dataset.

In this thesis, we refer to this basic technique as “simple transfer learning”. There are many other more complex ways of achieving pre-training and fine-tuning. In this section, we will present an overview of some of these methods.

3.1.1 Simple Transfer Learning

Simple transfer learning is a prevalent technique in medical image segmentation. Typically, segmentation neural networks in this field are initialized with weights trained on the ImageNet dataset [34]. ImageNet comprises a collection of natural images, which significantly differ in distribution from medical images. As a result, initializing with ImageNet-trained weights doesn’t usually enhance accuracy but it does make the network converge faster. This is because such initialization ‘skips’ the phase of learning fundamental image features that are common across various types of images, including medical ones. Figure 3.3 illustrates this process.

Simple transfer learning is used extensively in medical image segmentation. Segmentation neural networks are commonly initialized using weights trained on the ImageNet dataset [34]. ImageNet is a dataset of natural images and its distribution is drastically different from medical images. Therefore, this initialization does not typically lead to increases in accuracy but can make the network converge faster as it “skips” having to learn to detect base-level image features common to all images. This process is shown in Figure 3.3.

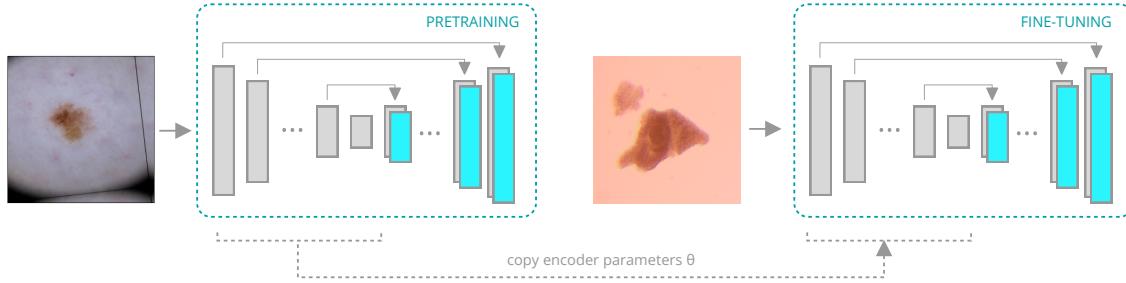


Figure 3.3: An overview of the simple transfer learning procedure. First, a model is pretrained on some related segmentation task. Then, part of the trained model’s weights are copied to a new model that is then fine-tuned on the target segmentation task.

Recent datasets, like Wasserthal, Breit, Meyer, *et al.* [35], comprising over 1000 subjects, are increasingly utilized for transfer learning, improving the accuracy of downstream tasks [36], [37]. While these approaches can result in impressive performance improvements, their application is constrained to fields where ample, similar data is accessible. In more specialized modalities or segmentation problems, the knowledge learned during pre-training may not be adequate for improving downstream accuracy. A potential solution to bridge this gap involves adapting the pre-trained model to the specific target domain, a topic we’ll explore in the following section.

3.1.2 Domain Adaptation

Neural networks often struggle to generalize across different datasets, even within similar domains [38]. This is referred to as the “domain shift” problem. Overcoming this problem could lead to an increase in data efficiency since a model trained on some similar large dataset can perform well on a smaller target dataset. Various domain adaptation techniques have been developed to tackle this challenge.

One class of methods aims to explicitly penalize domain shift by including an estimate of domain shift in the loss function. For instance, Liu, Lin, Wu, *et al.* [39] first train a network on a source domain with abundant data. This network is then adapted to perform segmentation on a target domain with limited data. During adaptation, the loss function includes a term for the maximum mean discrepancy between the encoded feature vectors of the source and target domains. Minimizing this term forces the encoder to maintain consistent feature representations rather than shifting the distribution of the encoded feature maps.

Another approach involves reducing the network’s ability to distinguish between source and target domains, thereby removing information that leads to domain shift. Ganin and Lempitsky [40] add a classification head to the network to identify the input’s domain. Typically, networks adjust their parameters to minimize loss during each iteration. However, in [40], the gradients from the domain classifier are inverted so the network updates its parameters to make it harder to differentiate between the two datasets. This method is illustrated in [Figure 3.4](#).

Domain adaptation has been effectively applied to medical image segmentation, as demonstrated in the work of Bermudez-Chacon, Marquez-Neila, Salzmann, *et al.* [41]. They utilized a U-Net-like architecture with dual encoder branches, each processing images from a different domain, paired with a shared decoder. Their loss function incorporates a maximum mean discrepancy term, comparing the decoder’s final feature maps

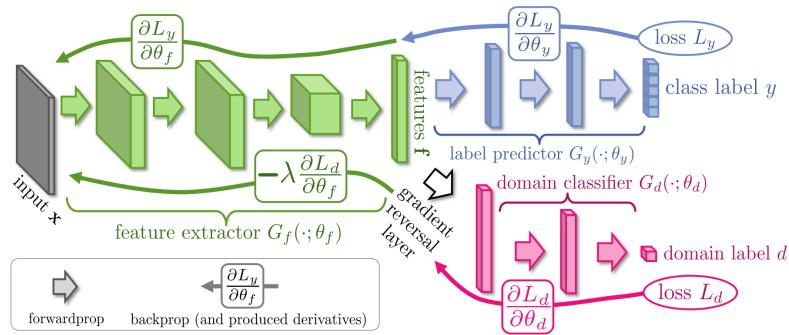


Figure 3.4: A diagram of the domain adaptation approach in Ganin and Lempitsky [40]. The gradients of the domain classification head which are applied to the encoder are reversed during backpropagation.

from inputs of both domains. This method was employed for segmenting mitochondria and synapses in microscopic images.

Although domain adaptation techniques enhance transfer learning performance, they generally require labeled data from the source domain. There is growing interest in methods that can utilize unlabeled data to identify general features in images from the target domain regardless of the specific task. We'll explore these methods in the following section.

3.1.3 Semi-Supervised and Self-Supervised Learning

Semi-supervised learning combines both labeled and unlabeled data to train neural networks. One such approach is **self-supervised learning**, where a feature extractor can be trained in a completely unsupervised manner using a large number of images. Then, that feature extractor can be employed via transfer learning to train a network for some specific task such as segmentation. One common strategy for this involves a **pretext task**, where the network is trained on a task for which solutions can be automatically generated from unlabeled data, using automatically generated labels to train the network. A pretext task could be, for instance, solving jigsaw puzzles [42]. An image is divided into 3×3 tiles, the tiles are randomly shuffled and fed into the network. The network is tasked with reordering these tiles to reconstruct the original image. Solving this task forces the network to learn to identify salient features such as recognizing a single object across multiple tiles. Such an approach is visualized in [Figure 3.5](#).

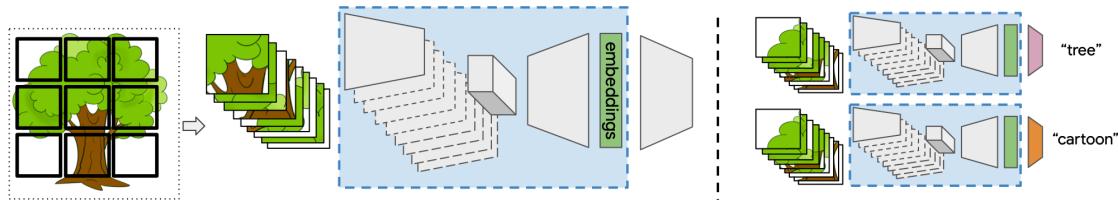


Figure 3.5: A self-supervised learning approach using shuffling image tiles as a pretext task presented by Carr, Berthet, Blondel, *et al.* [43]. The trained feature encoder learns to extract relevant features and its parameters are transferred to a model trained to perform the downstream task.

Recently, a more common approach to self-supervised learning is **contrastive learning**. This approach focuses on training the encoder to minimize the distance between feature vectors of similar (or positive) examples while maximizing the distance between

dissimilar (or negative) examples. Positive examples are generated in an unsupervised manner, often by applying random augmentations to an image to create two variants. The goal is for the feature vectors of these variants to be closely aligned. This teaches the network to recognize similar objects in a way that is invariant to the augmentations. Notable implementations of contrastive learning include SimCLR [44] and MoCo [45].

3.2 Synthetic Data

The data efficiency methods we've discussed so far revolve around pre-training networks with unlabeled or similar data. An alternative strategy involves generating synthetic data that mimics the target distribution. There are various techniques for this, each offering different levels of complexity and generation quality. A straightforward and commonly used method for expanding dataset size is **data augmentation**.

Data augmentation refers to the application of random transformations to input images, thereby artificially increasing their diversity. These transformations can be simple ones such as adjusting brightness or contrast, flipping the image horizontally, or more intricate, involving complex random deformations of the image. Generally, the choice of transformations is governed by domain knowledge — the network is forced to disregard transformations that we know it should be invariant to. Data augmentation is a standard practice in segmentation models, often employed alongside other data efficiency techniques. However, its implementation typically relies on simple deformation or transformation methods, which limits its capacity to create truly diverse samples.

Recently, the field of deep learning-based image generation, especially generative adversarial networks and diffusion models, has sparked a large interest in generating synthetic medical images. **Generative adversarial networks** [46] consist of two neural networks, a generator and a discriminator, that are trained simultaneously. The generator creates synthetic images intended to be indistinguishable from real images, while the discriminator learns to differentiate between the two. The process continues until the generator produces images so convincing that the discriminator cannot easily distinguish them from real images. This technique has been particularly effective in generating realistic medical images for data augmentation [47].

Diffusion models [48], a more recent development in the field of generative models, offer an alternative approach to image generation. These models, inspired by the physical process of diffusion, work by gradually adding noise to an image and then learning to reverse this process. The result is the generation of new images by reversing the noise addition process from a randomly sampled noise distribution. Diffusion models have been shown to generate high-quality images that can be particularly useful in medical imaging contexts [49].

3.3 Regularization

Aside from getting new data, the easiest way to reduce the estimation error is by limiting the size of the function class our model optimizes. In a traditional machine learning approach, this would be done by hand-selecting a class of functions with a lower number of parameters and then training the model. In deep learning, however, this approach is usually replaced with **regularization**, where the parameters themselves are manipulated to simplify the model.¹ Instead of selecting the number of parameters, we can penalize

¹Often, any method of reducing the estimation error of the model is broadly termed as regularization.

the network for the complexity and size of the parameters. In other words, the loss function between a predicted model output y_{pred} and a ground-truth label y has the form:

$$L(\theta; y_{pred}, y) = L_{seg}(\theta; y_{pred}, y) + \alpha\Omega(\theta) \quad (3.7)$$

where L_{seg} is the segmentation loss function while $\Omega(\theta)$ is a parameter norm penalty, i.e. a measure of the complexity of the model [50]. As we increase α , the strength of the regularization increases and the loss is more impacted by the complexity of the model. As the network is being trained, it will seek to both reduce $\Omega(\theta)$ as well as L_{seg} , leading to a more constrained model.

A commonly used parameter norm penalty is the L^2 regularization, also known as **weight decay** in deep learning contexts. Assuming θ consists of a set of weights w , $|w| = n$ and biases b , the L^2 norm is defined as $\Omega(w) = \frac{1}{2}\sqrt{(w_1^2 + w_2^2 + \dots + w_n^2)}$. In other words, the network is exponentially penalized for moving weights away from zero. This regularisation, in effect, forces the network to use fewer parameters.

Thus, a deep learning network can learn to use fewer parameters without necessarily changing the architecture of the network itself. However, what remains unsolved is how to transform the data such that the network can solve the problem with fewer parameters. We will discuss that in the next chapter, where we present another approach to achieving data efficiency.

3.3.1 Conclusion

As we have seen, there are several approaches to increasing the data efficiency of neural network-based segmentation of medical images. This thesis focuses on using regularisation by transforming the input data in such a way as to allow the network to use fewer parameters to achieve its segmentation task. We do so by combining predictions of neural networks with more traditional approaches from image processing and data augmentation. In the next chapter, we will introduce this idea and some specific applications.

4 Data Efficiency via Model-Driven Preprocessing

In this chapter, we will propose another way of achieving data efficiency in neural network-based medical image segmentation. This approach relies on transforming the input data such that it is segmentable with fewer parameters. We achieve this by leveraging domain knowledge and traditional image-processing techniques to select image transformations. However, instead of manually selecting transformation parameters, we will model the parameters as functions of the image, and train a neural network to estimate optimal transformation parameters. As described in the previous chapter, regularization allows a network to use as few parameters as required by the problem's complexity. By transforming the image we reduce the segmentation problem complexity and therefore achieve segmentation with fewer parameters, i.e. we improve data efficiency.

We will first present some motivation to show how it is possible to transform data to be easier to model. Next, we will present a general description of our approach. Finally, we will show an empirical study of our approach to various medical image segmentation tasks.

4.1 Motivation: Using The Polar Transform for Preprocessing

We will first describe yet another way to view neural networks — as algorithms to find a linearly separable transformation of an input space. To simplify, let us focus on a 2-class classification problem, with a dataset $(X_i, Y_i), X_i \in \mathbb{R}^2, Y_i \in \{0, 1\}$. All of our input vectors X_i come from an input space \mathcal{X} . The neural network is tasked with finding a decision boundary that splits the input space into two sides, one labeled 0 and the other labeled 2. For complex problems, this decision boundary is not linear in the input space. For instance, we can imagine a dataset where points belonging to class 0 are clustered around $(0, 0)$, while points belonging to class 1 are distributed along a circle with radius 2 with additive Gaussian noise, as is shown in [Figure 4.1\(left\)](#). In the input space, the optimal decision boundary is circular, and therefore the input space cannot be linearly separated in its current form. Instead, a neural network learns a transformation $\mathcal{X} \rightarrow \mathcal{X}'$ such that a hyperplane can be drawn in \mathcal{X}' to optimally separate points belonging to classes 0 and 1. This is illustrated in [Figure 4.1\(right\)](#).

To do this, the network needs sufficient parameters — more parameters allow for more complex transformations. Without enough parameters, no hyperplane in the transformed space can separate the two points well, leading to underfitting. To approach correctly approximating a circular decision boundary, the network needs at least three neurons in its hidden layer. We exemplify this in [Figure 4.2](#).

Therefore, we will focus on finding transformations of the data such that the problem becomes easier. To illustrate why this is effective, we will present yet another view of neural networks. Let there be some neural network that we are using to determine

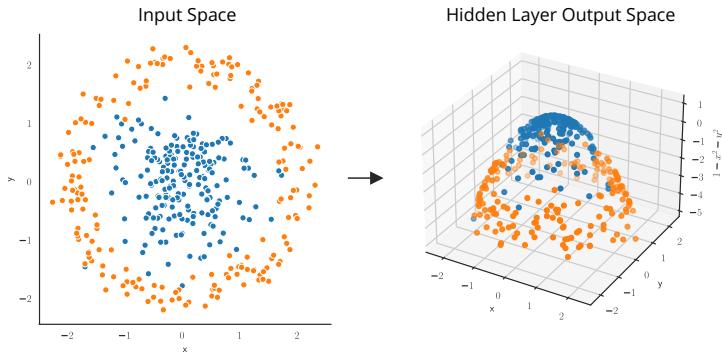


Figure 4.1: An example of how a 2-class classification neural network bends an input space (left) such that it is linearly separable with a hyperplane.

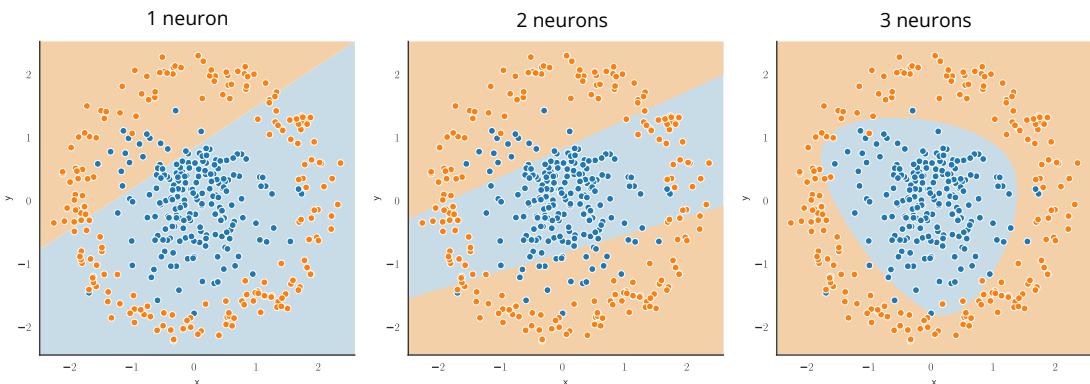


Figure 4.2: The decision boundary (shown as the background color) according to a fully connected neural network with one hidden layer where the hidden layer has one, two, or three neurons.

whether a given image region $I(x, y)$ belongs to class “0” or “1”. All possible image regions $I(x, y)$ come from a d -dimensional space called the input space where d is the number of pixels in the region, multiplied by the number of channels in each pixel. Through optimization, the neural network iteratively finds a transformation of the input space such that this space can be linearly separated by a d -dimensional hyperplane into two classes. This is illustrated in Figure 4.1.

The goal of data efficiency is to find a decision boundary given fewer data samples. As discussed in the previous chapter, this necessitates reducing the function class space we are searching over, or in other words, reducing the number of neural network parameters. In this and the following chapters, we will focus on finding ways to pre-transform the data into some intermediary form \mathcal{X}_{pre} that is more separable than the original input space \mathcal{X} . This allows networks with fewer parameters, as \mathcal{X}_{pre} is closer to X' , and thus the network has to search a smaller space to find X' . We will call this transformation \mathcal{T} .

A commonly used example of \mathcal{T} is the polar transform. In the two-dimensional case, the polar transform $\mathcal{T}(X_i; C)$, parameterized by a chosen origin $C = [c_0, c_1]$, transforms an input vector $X_i = [x_0, x_1]$ into a vector $\mathcal{T}(X_i; C) = [\theta, r]$ as:

$$\begin{aligned} r &= \sqrt{(x_0 - c_0)^2 + (x_1 - c_1)^2}, \\ \theta &= \text{atan2}(x_1 - c_1, x_0 - c_0). \end{aligned} \tag{4.1}$$

The polar transform turns a circular decision boundary into a linear one, allowing

linear separability of our example dataset even given only one neuron in the neural network, as can be seen in [Figure 4.3](#)

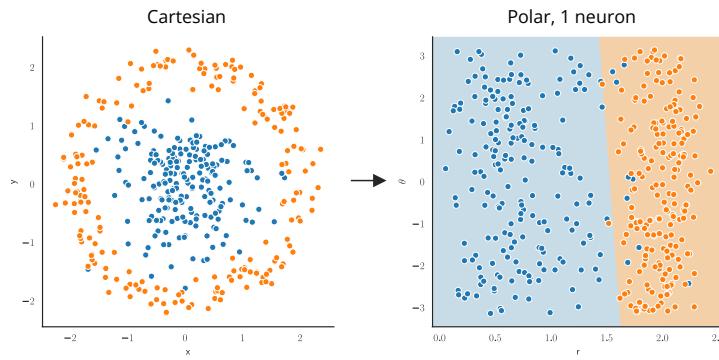


Figure 4.3: An example circular dataset transformed using the polar transform. The background color shows the decision boundary of a fully connected network with one hidden layer with a single neuron.

Here, we have shown how using the polar transform already transforms the input space into a linearly separable form. The network effectively only needs to search the space of lines to find one that optimally separates the two classes.

However, the effectiveness of the polar transform in easing separability rests in choosing the correct polar origin. As we shift the origin outside of the center of class 0, the decision boundary stops being linear in the polar space, as is shown in [Figure 4.4](#).

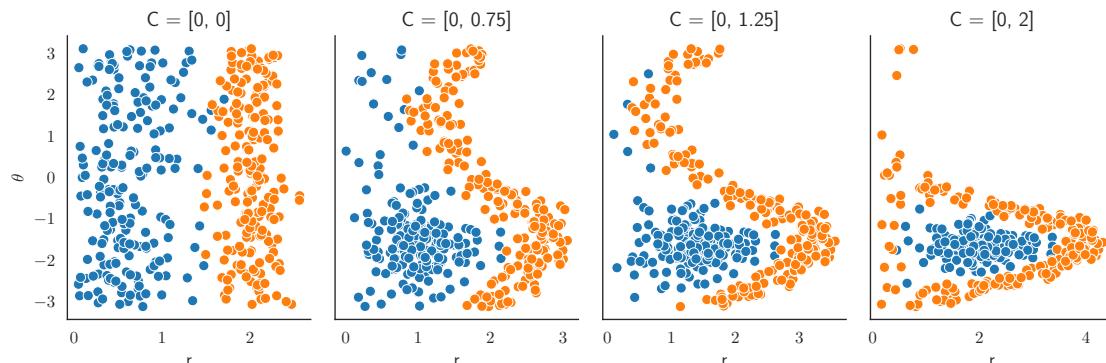


Figure 4.4: The polar transformation of the simulated data using different polar origins (C).

With this example in mind, let us present a general approach of using model-driven preprocessing to achieve data efficiency in medical image segmentation.

4.2 Model-driven Preprocessing

The example above with the polar transform can be extended to image segmentation. Similarly to the simulated data we used in the example, we can think of a single image as having pixels of one class arranged inside a circle, while the background lies outside the circle. The segmentation model would have to approximate the circular decision boundary to segment the image. When transformed with the polar transform, the circular decision boundary becomes a linear one, and the decision boundary can be predicted using a model with fewer parameters.

We propose a general formulation of this approach, wherein images are transformed to be more linearly separable using hand-constructed transformations using traditional image processing methods. Since the effectiveness of the transformation rests on choosing the correct parameters, we train neural networks to predict optimal transformation parameters. Thus, we use a model-driven approach to preprocess the image. The pre-processed image is then segmented using a regular neural network. However, because the image is easier to segment, the segmentation network can have fewer parameters than if we were segmenting the original, unprocessed images. In the remainder of this section, we will describe this approach in more detail.

To pre-process the image, we may choose a transformation $\mathcal{T}_\phi(I)$ of the image I , parameterized by transformation parameters ϕ , that will produce a transformed image I' segmentable by fewer parameters than I . The choice of transformation function \mathcal{T} depends on the problem at hand. When segmenting elliptically objects from their background, the polar transform is beneficial. When segmenting small objects, we might opt for an affine transformation that scales each object to a uniform scale. Just like designing the segmentation network architecture, we design \mathcal{T} using domain knowledge and verify its effectiveness empirically. We may then perform the segmentation on the transformed image as:

$$M(I) = (\mathcal{T}_\phi^{-1} \circ F_{seg} \circ \mathcal{T}_\phi)(I), \quad (4.2)$$

where $M(I)$ is the resulting segmentation mask for image I , F_{seg} is a segmentation neural network. The transformation can be any transformation function including affine transformations, deformation fields, thresholding as well as the polar transform, as in the above example. The transformation needs to be invertible so that the segmentation mask $M(I)$ can be obtained in the original image space.

Note that the correct value for ϕ matters greatly to the effectiveness of the preprocessing, as we have observed with the polar transform. Moreover, ϕ is image-dependent: The origin, scale, location, and distribution of the object vary from image to image. Given this, we may train a model to estimate an optimal ϕ given an image:

$$\phi(I) = F_{\mathcal{T}}(I; \theta_{\mathcal{T}}), \quad (4.3)$$

where $F_{\mathcal{T}}$ is a neural network with parameters $\theta_{\mathcal{T}}$ trained to estimate the correct parameters for \mathcal{T} given an image I . In our approach, we use various types of neural networks as $F_{\mathcal{T}}$, though it is also possible to use more traditional approaches.

More specifically, $F_{\mathcal{T}}$ is structured similarly to a regular segmentation network, producing a segmentation mask $M_\phi(I)$ in its second-to-last layer. The final layer of the network is a non-trainable function $\phi(I) = f_\phi(M_\phi(I))$, which produces segmentation parameters given a predicted segmentation map. $F_{\mathcal{T}}$ can then be expressed as:

$$F_{\mathcal{T}}(I; \theta_{\mathcal{T}}) = f_\phi(F_{\mathcal{T}seg}(I; \theta_{\mathcal{T}})), \quad (4.4)$$

where $F_{\mathcal{T}seg}$ is the segmentation sub-network producing $M_\phi(I)$. $F_{\mathcal{T}seg}$ can then be trained using segmentation-like loss functions, where M_ϕ can be compared to a ground truth segmentation mask $M(I)$.

The function f_ϕ is hand-constructed using knowledge of the transformation. For instance, we know that for the polar transform where ϕ is the origin of the polar transformation, the optimal origin lies in the center of the object. Therefore, we may use the centroid of the ground truth segmentation labels as f_ϕ . There are other examples of f_ϕ :

visual de-
scription

- For CT thresholding, the transformation parameters are the maximum and minimum Houndsfield values of some tissue. For f_ϕ we may use the minimum and maximum values within the predicted segmentation mask.
- For rescaling, the transformation parameters are the x and y scaling factors. We may use the predicted segmentation mask to determine a scaling factor such that the object is of some given fixed width and height.
- For translation, we can determine x and y translation parameters by calculating how far away the center of the object is from a desired location.

Note that $M_\phi(I)$ does not necessarily need to be a regular segmentation mask. Instead, we can construct pseudo-labels that better fit the task of predicting optimal transformation parameters. For instance, in cases where the transformation parameters depend on the centroid of the object, we can model $M_\phi(I)$ as a heatmap centered on the center point of the object. The network will then learn to predict object centroids instead of the object boundaries, and f_ϕ can calculate the centroid of the object by looking up the x and y coordinate of the maximal intensity pixel in $M_\phi(I)$.

The core idea of our approach is that \mathcal{T} will simplify the segmentation problem such that a regularized segmentation network F_{seg} will be able to perform the segmentation using fewer parameters. To achieve this, both $F_{\mathcal{T}}$ and F_{seg} need to be trained using specific training methods, as we will describe in the next section.

Using this approach, data efficiency is achieved by reducing the total function class size when compared to a standard segmentation neural network. We do so by using two smaller neural networks instead of one large one. Since the decision boundaries of the preprocessed images are simplified, the segmentation network needs fewer parameters to model them. Transformation parameters often have a much smaller dimensionality than the segmentation mask, and therefore the network predicting transformation parameters can naturally use much fewer parameters.

4.2.1 Training Model-Driven Preprocessing Networks

As mentioned earlier, using the approach of model-driven preprocessing allows and requires specific neural network training techniques. In our approach, two separately trained neural networks are used: the transformation parameter predictor $F_{\mathcal{T}}$ and the segmentation network F_{seg} .

Training the Transformation Parameter Predictor

$F_{\mathcal{T}}$ receives the original untransformed images as input and outputs the segmentation parameters $\phi(I)$. As mentioned earlier, the second-to-last layer of $F_{\mathcal{T}}$ outputs a segmentation-like mask $M_\phi(I)$. The data used to train $F_{\mathcal{T}}$ is in the form of (I, M_{gt}) . Given a predicted segmentation mask $M_\phi = F_{\mathcal{T}seg}(I)$ and the transformation parameter function $f_\phi(M(I))$, the transformation parameters loss $L_{\mathcal{T}}$ can be calculated as:

$$L_{\mathcal{T}}(I, M_\phi, M_{gt}) = \alpha L_\phi(f_\phi(M), f_\phi(M_{gt})) + \beta L_M(M_\phi, M_{gt}), \quad (4.5)$$

where L_ϕ measures the error in predicting the transformation parameters, while L_M is a segmentation-like loss that measures the error in predicting the segmentation mask used to calculate the transformation parameters. α and β control the strength of the two terms and are chosen empirically.

Training the Segmentation Network

To train the segmentation network, a dataset of ground truth transformed images. For each image and ground truth segmentation label pair (I, M_{gt}) , the transformation parameters are obtained as:

$$\phi = f_\phi(M_{gt}) + \gamma z_\phi, \gamma \in [0, 1] \quad (4.6)$$

where f_ϕ is the same transformation parameter function used in the parameter predictor network, and z is a random noise value of the same dimensions as ϕ , where γ determines the strength of the noise. During training, $\gamma > 0$ and the random noise acts as a type of augmentation. This makes the segmentation network more robust to suboptimal ϕ predictions. During inference, γ is set to zero.

Given ϕ , the input to the image is constructed as $I' = \mathcal{T}_\phi(I)$. The last layer of the network performs an inverse transformation

4.3 Preprocessing With the Polar Transform Using Centerpoint Prediction

Images are most commonly viewed in Cartesian coordinates, where the pixels are arranged along the x- and y-axes. The polar coordinate system has two axes: (1) the radial coordinate ρ , which is the distance of a point from the origin of the polar transformation; and (2) the angular coordinate ϕ , which is the angle between the point and the reference direction. In other words, the x-axis of the polar image represents the distance from an origin, while the y-axis represents the rotation around the origin. This makes polar coordinates invariant to rotation.

Our intuition is that polar transformations can be especially beneficial to segmenting images where an elliptical border must be found on the image. Consider a contrived example of predicting a circular decision boundary on a single-channel image with a linear model. A circular decision boundary must be modeled by a function of at least four dimensions. When transformed to polar coordinates, a perfect circle in Cartesian coordinates becomes a straight line, as shown in [Figure 4.5](#). This linear decision boundary can be modeled with a simpler linear function in two dimensions. The image in polar coordinates would require a less complex model to predict a border. It is possible that, even for more complex examples, the polar transform of an image of a roughly elliptical object reduces the required segmentation model complexity, as shown visually in [Figure 4.6](#). Furthermore, by transforming an image to polar coordinates using a polar origin that is the center of the object, we fix the location and standardize border distances in each training example. The model can then learn the distance of the border from the origin at each angle around the origin, without having to learn to localize the object.

To obtain a polar transformation, the angle and magnitude of each pixel (x, y) of the original image are calculated using [\(4.7\)](#):

$$\begin{aligned} magnitude(x, y) &= \sqrt{x^2 + y^2}, \\ angle(x, y) &= atan2(y, x)[\cdot 180 / \pi] \end{aligned} \quad (4.7)$$

Where $atan2$ is the 2-argument arctangent function.

Given a polar origin (c_x, c_y) of a Cartesian image $I(x, y)$ of resolution $H \times W$, we obtain each point (ρ, ϕ) of the polar transformation $I'(\rho, \phi)$ using [\(4.8\)](#).

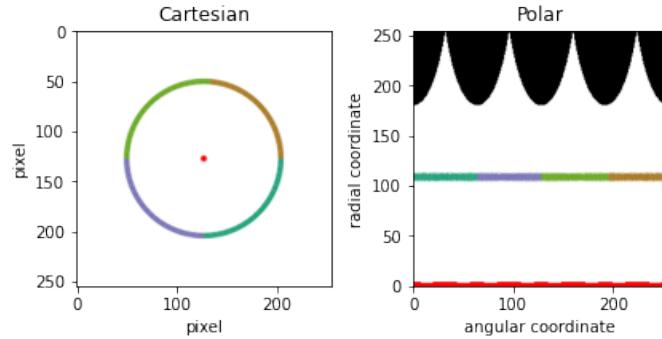


Figure 4.5: An example of a polar transformation.

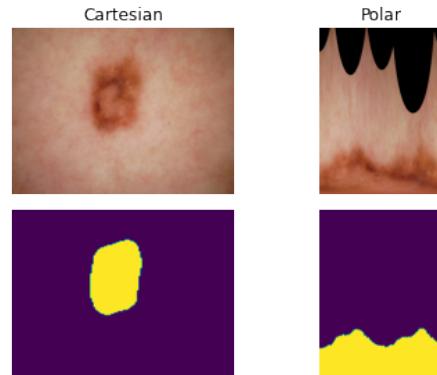


Figure 4.6: An example image and label from the lesion dataset and their corresponding polar transformation.

$$\begin{aligned} \rho &= \frac{H}{2\pi} \cdot \text{angle}(x - c_x, y - c_y) \\ \phi &= \frac{W}{\sqrt{(W/2)^2 + (H/2)^2}} \cdot \text{magnitude}(x - c_x, y - c_y) \end{aligned} \quad (4.8)$$

In this chapter, we propose a general way to improve neural network segmentation data efficiency and performance on medical imaging segmentation tasks where the goal is to segment roughly elliptically distributed objects.

We propose and explore ways to train neural networks for biomedical image segmentation on polar transformations of images. The polar transformation transforms an image from Cartesian coordinates into a new coordinate system where the two axes are the rotation around an origin and radius from that origin. When the regions to be segmented are elliptical in shape or distribution, this transformation results in a reduction of dimensionality, allowing convergence in fewer epochs and good performance even in models with a low number of parameters.

Experimentally, we observed that selecting a correct polar origin is one of the key parameters that determine segmentation performance. Therefore, we propose two different approaches of selecting an optimal polar origin: (1) estimation via a segmentation neural network trained on non-polar images and (2) estimation via a neural network trained to predict heatmaps. Our method is evaluated on the tasks of polyp segmentation, liver segmentation, skin lesion segmentation, and epicardial adipose tissue (EAT) segmentation. The proposed methods can be used as a pre-processing step for existing

formulate problem and introduce using the above framework

neural network architectures, so we evaluate the methods using common neural network architectures for medical image segmentation including U-Net [51], U-Net++ [24] with a ResNet [52] encoder, and DeepLabV3+ [53] with a ResNet [52] encoder.

Evaluation of our approach as a pre-processing step shows that it improves segmentation performance across different datasets and neural network architectures while making the networks more robust to small dataset sample sizes.

4.3.1 Previous work using the polar transform for medical image processing

Several image segmentation methods were proposed that utilize polar coordinates. Liu, Hong, Ke, *et al.* [54] proposed an approach they call Cartesian-polar dual-domain network (DDNet) to perform optic disc and cup segmentation in retinal fundus images. The neural network contains two encoding branches, one for a Cartesian input image and another for the polar transformation of the same input image. The predictions are fused into a single feature vector which is then decoded into a final segmentation. Salehinejad, Valaee, Dowdell, *et al.* [55] used the polar transformation as a way to augment training data by transforming each input image into multiple polar images at various polar origins, thus increasing the number of training data. Kim, Jung, Kim, *et al.* [56] proposed a convolutional neural network layer for images in polar coordinates to achieve rotational invariance. Their cylindrical convolution layer uses cylindrically sliding windows to perform a convolution. Kim, Sun, Kim, *et al.* [57] proposed a user-guided segmentation method where an expert selects the point used as the polar origin. The transformed image is then segmented using a convolutional neural network (CNN). Esteves, Allen-Blanchette, Zhou, *et al.* [58] proposed a polar transformer network for image classification. Note that "transformer network" here refers to spatial transformer networks [59] and not attention-based networks commonly called transformers. The network consist of a polar origin predictor and a neural network that predicts a heatmap. The centroid of the heatmap is then used as the origin for a polar transformation of the input image. The polar image is classified using a CNN. This approach is most similar to our proposed method, however, their approach focuses on image classification, not segmentation. Additionally, our approach differs in the ways the ground truth data is prepared, the used neural network architectures, as well as other details.

While there are proposed methods that combine polar transformation with neural networks, most of them solve classification tasks. Some medical image segmentation methods use the polar transformation as a preprocessing step, however, the way they obtain the origin of the polar transformation is usually based on heuristics. To our knowledge, there is currently no work that explores using the polar transformation with a dynamic polar origin as a preprocessing step for semantic segmentation in a variety of medical image datasets.

4.3.2 Methodology

All of the proposed methods rely on training a neural network model to segment polar images. To train on polar images, the input images need to be transformed using a polar origin which is near the center of the segmented object. The correct origin is not known ahead of time, so a prerequisite for predictions on polar images is a method to determine the correct polar origin. We propose and evaluate two different methods for automatically obtaining the polar origin: (1) estimation via a segmentation trained on non-polar images and (2) training a center-point predictor that predicts heatmaps from input images. This section describes these methods, as well as methods to train the final segmentation model on the polar images.

In each of our approaches, the final segmentation is done using a neural network trained on polar transformations of the input images. In the rest of this paper, we refer to this network as the polar network. In all of the described approaches, the polar transformation is not part of the network architecture itself, but happens as a preprocessing step for the polar network. To transform each input image, the polar origin is determined as the center of mass of the ground truth label for that image. The center of mass of an image $I(x, y)$ is calculated by first calculating the spatial image moments matrix M , where the entry of the matrix at row i and column j is calculated using (4.9).

$$M_{ij} = \sum_{x,y} I(x, y) \cdot x^i \cdot y^j \quad (4.9)$$

The center of mass (c_x, c_y) of the image can then be calculated using (4.10).

$$\begin{aligned} c_x &= M_{10} / M_{00} \\ c_y &= M_{01} / M_{00} \end{aligned} \quad (4.10)$$

Finally, to increase the model's robustness to suboptimal center point predictions, we augment the calculated center for training images [58]. Each training image has a 30% chance of varying the center's x and y coordinates by a random value in the range $(-S \cdot 0.05, S \cdot 0.05)$, where S is the smallest resolution of the image, i.e. $S = \min(\text{width}, \text{height})$.

4.3.3 Centerpoint prediction

Once the polar network is trained, inference can be done by transforming an input image to polar coordinates. The polar network requires choosing a center that is close to the center of mass of the segmented object. Because a future input image is unlabeled, the correct center needs to be inferred from the image. We propose two ways to accomplish this, described in this section.

(A) Training the same neural network on cartesian and polar images

Our first approach is training the same neural network on cartesian and polar images. A summary of this approach is presented in Figure 4.7. With this approach, the inference is done by first feeding the original Cartesian input images into a neural network used for segmentation. We refer to this network as the Cartesian network. For an input image, the polar origin is calculated as the center of mass of the Cartesian network's prediction for that input image, using (4.10). This polar origin is used to transform the original input image to polar coordinates, and the transformed image is fed to the polar network. The output of the polar network is transformed back to Cartesian coordinates to obtain a final segmentation. We assume identical architectures for both the Cartesian and polar networks. This makes applying this framework to existing architectures very straightforward, as it does not require designing new neural network architectures or specific hyperparameter optimization, and allows for using transfer learning to initialize the networks.

(B) Training a centerpoint predictor

In the second approach for determining the optimal polar origin, we train a model specifically tasked with predicting the correct polar origin for each input image, which is then used to transform the input image. The approach is shown in Figure 4.8. We do this by training a neural network based on the stacked hourglass architecture [60] first used for human pose estimation. Instead of training a regressor network to predict key points

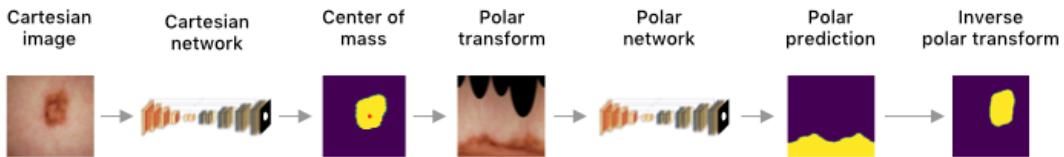


Figure 4.7: A diagram of the approach of predicting polar origins from a Cartesian network. The first network performs an initial segmentation, which is then used to extract a polar origin for the polar transformation. The method does not rely on any specific neural network architecture. The Polar and Cartesian network can be any neural network which takes an input image and produces a binary segmentation mask as output. The red point shows the extracted polar origin. The Polar network is trained on polar image transformations. The polar transformation is not part of the network itself, but happens as a preprocessing step for the Polar network.

in an image, the stacked hourglass architecture uses a series of stacked encoder-decoder networks, where the output of each stack is a heatmap centered on the key point to be predicted. The output of each stack is fed as input into the next stack, allowing successive refinement of the heatmap prediction. During training, the loss of each stack's output is averaged to produce the final loss, allowing deep supervision. The final prediction heatmap is the output of the last stack in the network. To predict the center point, we use 8 stacked hourglass blocks, which we empirically determined as the value providing the best results. The network receives images in Cartesian coordinates and predicts a heatmap of the image.

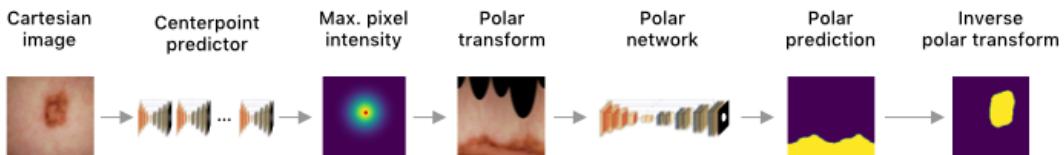


Figure 4.8: A diagram of the approach of using a centerpoint prediction network. The first network can be any neural network which predicts a heatmap from an input image, which is then used to extract polar origin, shown as a red point. The Polar network can be any semantic segmentation neural network which produces a binary mask output from an input image. The Polar network is trained on polar image transformations. The polar transformation is not part of the network itself, but happens as a preprocessing step for the Polar network.

The ground truth heatmaps were generated by calculating the center of mass of each ground truth label image using (4.10). We then create the heatmap as an image with a 2D gaussian with the mean on the center of mass on the image and a standard deviation of 8 pixels for all datasets except the liver, and 16 for the liver. Example heatmaps are shown in Figure 4.9. The optimal value for the standard deviation was determined empirically on the validation datasets. We found that the optimal value of the standard deviation is proportional to the size of the object.

Additionally, during training, we use augmentation to increase the number of training inputs. In particular, during training each input example the following random augmentations are applied:

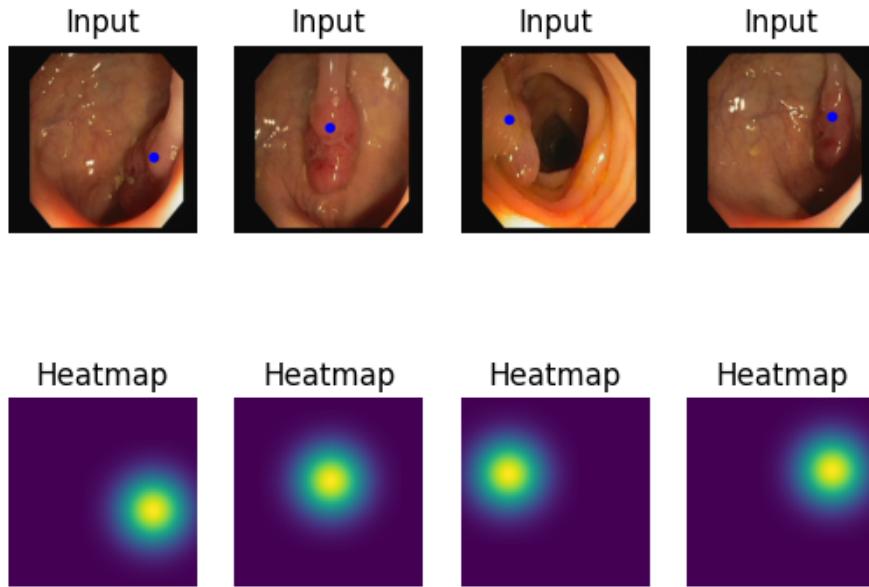


Figure 4.9: Examples of heatmaps generated from the ground truth data. The heatmap is a gaussian centered on the center of mass of the ground-truth label, shown as a blue point on the input images.

- A 50% chance of a horizontal flip.
- A 30% chance of a random combination of shifting up to 6.5% of the image dimensions, scaling up to 10% and rotating up to 45°.
- A 30% chance for a grid distortion, details of which are described in [61].

The center-point predictor outputs 8 separate heatmaps [60]. We calculate the predicted center as the coordinates of the pixel with the largest intensity in the heatmap predicted by the final layer of the model. This predicted center is then used to transform the input image to polar coordinates, and the transformed image is fed into the polar network to perform the segmentation. Finally, the segmentation label is transformed back to Cartesian coordinates.

4.3.4 Experiments

To validate the generality of our approach, we trained a variety of neural network architectures on multiple medical imaging datasets. In particular, we trained three different neural network architectures: U-Net [51], U-Net++ **zhouUNetNestedUNet2018a** with a ResNet encoder and DeepLabV3+ **chenEncoderDecoderAtrousSeparable2018a** with a ResNet encoder. Notably, each dataset we use presents a problem wherein almost all examples a single roughly elliptical object needs to be segmented. For each dataset and network architecture combination, we train a Cartesian and polar network, and we then perform four different experiments:

1. testing the Cartesian network using Cartesian images
2. testing the polar network using the ground-truth polar origin
3. testing the polar network using polar origins obtained from predictions of the Cartesian network, as outlined in 4.3.3

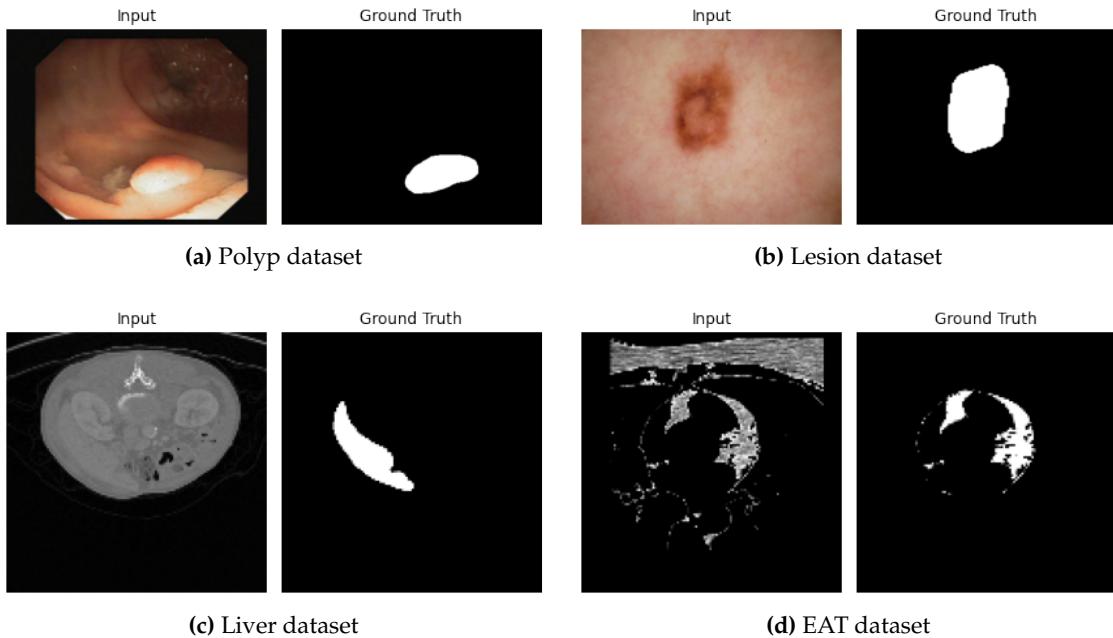


Figure 4.10: Example input images and ground-truth labels for each dataset used in our experiments.

4. testing the polar network using polar origins from the center-point predictor, as outlined in 4.3.3.

4.3.5 Datasets description

We used four different datasets to train the network. In this section, we give an overview of each used dataset and how it was preprocessed. Note that for training the center-point predictor, the input images were resized to a resolution of 256×256 , while the generated heatmaps were resized to 64×64 pixels. Otherwise, all preprocessing steps described here are applied to the center-point model datasets as well. Each dataset was normalized and zero-centered to better facilitate network convergence.

Polyp dataset

The CVC-ClinicDB dataset [62] contains 612 RGB colonoscopy images with the resolution 288×384 with labeled polyps from MICCAI 2015. We normalize each image to a range of $[-0.5, 0.5]$. We use the original image resolution to train all networks except the centerpoint network. As is used in [63], we use an 80%, 10% and 10% split for training, validation and testing datasets, respectively. An example of the dataset is shown in Figure 4.10(a).

Liver dataset

The second dataset we use is the LiTS dataset [64] from the Liver Tumor Segmentation Challenge from MICCAI 2017. The dataset contains 131 CT scans of patients with hepatocellular carcinoma, with the liver as well as tumor lesions labeled by experts. In our experiments, we disregard the lesion segmentation labels and treat the dataset as a binary liver segmentation problem. In addition, we removed all slices that did not contain a ground-truth liver segmentation label, resulting in a dataset of roughly 15,000 slices.

Each axial slice is thresholded to a Hounsfield scale range of $[0, 200]$ HU that contains the liver. Next, the slices are normalized to a $[0, 1]$ range and zero-centered by subtracting the global intensity of all training slices (0.1). We then proceed to train the networks on each axial slice separately. We use 101 scans for training, 15 scans for validation, and the remaining 15 scans for testing. Example liver segmentation images are shown in Figure 4.10(c).

Lesion dataset

The third dataset we use is the ISIC 2018 Lesion Boundary Segmentation dataset [65], [66] which contains 2,694 dermatoscopy images of skin lesions with expert labels of the lesions from various anatomic sites and several different institutions. We resize each image to a resolution of 384×512 and use a training, validation and test split of 80%, 10% and 10%, respectively. This is consistent with [63]. Additionally, we normalize each image to a range of $[-0.5, 0.5]$. An example of a lesion input image and its corresponding label is shown in Figure 4.10(b).

EAT dataset

Finally, we also train on a dataset of labeled EAT regions from 20 patients' cardiac CT scans from the Cardiac Fat Database [67]. The dataset has three classes labeled: the pericardium, EAT, and pericardial adipose tissue. We disregard all original labels except EAT and treat the dataset as a binary EAT segmentation dataset. The dataset is first split into training (10 patients), validation (5 patients), and test (5 patients) datasets. In the original dataset, each slice is thresholded to the adipose tissue range of $[-200, -30]$ HU and registered so that anatomical structures have the same locations. In addition to these original pre-processing steps, we normalize each slice to a $[0, 1]$ range and zero-center the dataset by subtracting a global mean intensity of the training set (0.1). We then train on each CT slice separately. An input image of the EAT dataset and its corresponding label is shown in Figure 4.10(d).

4.3.6 Implementation details

We use the OpenCV linear polar transformation implementation. Each model is implemented and trained using PyTorch 1.7.1 on an NVIDIA GeForce RTX 3080 GPU. For all networks, we use the Adam optimizer with a learning rate of 10^{-3} . A batch size of 8 was used for all networks except the center-point model, where a batch size of 6 was used for the lesion and liver datasets, and 8 for all remaining datasets. We trained all models up to a maximum of 200 epochs and used checkpoints after each epoch to store the model with the best validation loss. We modify the Dice coefficient to act as a loss function as shown in (4.11).

$$DSC_{loss} = 1 - \frac{2|X \cap Y| + \lambda}{|X| + |Y| + \lambda} \quad (4.11)$$

Where X and Y are the input and predicted images, respectively, and λ is a smoothing parameter set to 1 in our experiments. This loss function is used to train all models except the center-point model.

The centerpoint model outputs eight heatmaps [60]. We use a loss function that is the mean of the mean squared errors between each of the heatmaps and the ground truth heatmap. The code used for all experiments is available at github.com/marinbenc/medical-polar-training.

4.3.7 Results

We evaluate segmentation performance along with four key metrics: the Dice coefficient (DSC), the median intersection-over-union score (mIoU), precision, and accuracy. Precision and accuracy are both calculated pixel-wise. The results of training the different approaches presented in 4.3.4 are shown in Table 4.1 for polyp, lesion, liver and EAT segmentation. In all cases, training on polar coordinates improves the segmentation in all metrics when compared to training the same model on Cartesian coordinates. As is to be expected, testing the polar network on images transformed using the ground truth polar origins produces the best results. A close second is predicting the polar origin from the center-point predictor. Predicting polar origins from the Cartesian model leads to less accurate polar origins, and the results are worse, however, they are still better than using only the Cartesian model.

We also compare our methods to other state-of-the-art methods that use the same datasets, shown in Table 4.2. We achieve state-of-the-art results for the polyp and liver datasets. Additionally, we achieve state-of-the-art liver segmentation when compared to other per-slice methods, and nearly state-of-the-art results when compared to 3D-based methods. For EAT segmentation, our approach outperforms standard medical image segmentation networks but does not achieve state-of-the-art performance due to segmenting EAT directly and not first segmenting the pericardium.

A training graph for a polar and Cartesian U-Net-based network is shown in Figure 4.11. Additionally, we evaluate the accuracy of the different ways of obtaining the polar origin. This accuracy is compared with segmentation performance in Figure 4.12.

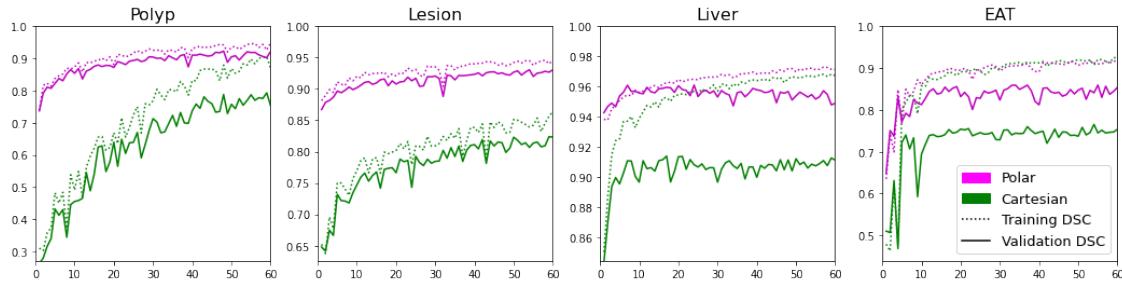


Figure 4.11: The training and validation Dice coefficient (DSC) of the polar and Cartesian U-Net models during training.

We also train several models with both polar and Cartesian coordinates on subsets of the training dataset. Namely, we trained models on 25%, 50%, 75%, and 100% of the lesion training dataset for 50 epochs. The results of this training are shown in Figure 4.13. The polar network is much more data efficient and achieves better results than the cartesian network even with only 25% of the data.

4.4 Discussion

We obtain state-of-the-art results for polyp and lesion segmentation by training common biomedical image segmentation models. In the liver dataset, we achieve state-of-the-art results when compared to other 2D methods, but 3D methods achieve the same or slightly better results [70]. The liver dataset is by far the largest dataset we evaluated. As such, improvements gained from encoding localization information and reducing dimensionality might not be as large as in smaller datasets, since the network has enough data to learn these complex structures. The EAT dataset is one where the task is not to find a single object, but instead, segment multiple smaller pockets of tissue around the

| Polyp Segmentation | | | | | |
|---------------------|---------------|--------|--------|--------|--------|
| Architecture | Method | DSC | mIoU | Prec. | Rec. |
| U-Net | Cartesian | 0.8315 | 0.7604 | 0.8513 | 0.8334 |
| | GT centers | 0.9484 | 0.9141 | 0.9563 | 0.9442 |
| | Cart. centers | 0.8973 | 0.8571 | 0.8996 | 0.8998 |
| | Model centers | 0.9374 | 0.8977 | 0.9488 | 0.9368 |
| Res-U-Net++ | Cartesian | 0.8356 | 0.7636 | 0.9004 | 0.8256 |
| | GT centers | 0.9557 | 0.9260 | 0.9583 | 0.9554 |
| | Cart. centers | 0.9063 | 0.8685 | 0.9243 | 0.9027 |
| | Model centers | 0.9332 | 0.8924 | 0.9477 | 0.9321 |
| DeepLabV3+ | Cartesian | 0.8706 | 0.8013 | 0.8857 | 0.8876 |
| | GT centers | 0.9593 | 0.9296 | 0.9576 | 0.9682 |
| | Cart. centers | 0.9212 | 0.8823 | 0.9179 | 0.9397 |
| | Model centers | 0.9338 | 0.8967 | 0.9436 | 0.9347 |
| Lesion Segmentation | | | | | |
| | Method | DSC | mIoU | Prec. | Rec. |
| U-Net | Cartesian | 0.8256 | 0.7393 | 0.8407 | 0.8712 |
| | GT centers | 0.9320 | 0.8824 | 0.9261 | 0.9541 |
| | Cart. centers | 0.8836 | 0.8317 | 0.8746 | 0.9492 |
| | Model centers | 0.9224 | 0.8699 | 0.9165 | 0.9494 |
| Res-U-Net++ | Cartesian | 0.8664 | 0.7925 | 0.8728 | 0.9122 |
| | GT centers | 0.9439 | 0.9014 | 0.9418 | 0.9584 |
| | Cart. centers | 0.9125 | 0.8653 | 0.9075 | 0.9540 |
| | Model centers | 0.9253 | 0.8743 | 0.9253 | 0.9464 |
| DeepLabV3+ | Cartesian | 0.8717 | 0.7984 | 0.8807 | 0.9068 |
| | GT centers | 0.9459 | 0.9059 | 0.9418 | 0.9632 |
| | Cart. centers | 0.9162 | 0.8686 | 0.9097 | 0.9536 |
| | Model centers | 0.9235 | 0.8721 | 0.9125 | 0.9570 |

Table 4.1: Results of our proposed approaches for different tasks for three different neural network architectures. The cartesian network is the network trained on Cartesian images. "GT centers" refers to obtaining a polar origin from the ground-truth labels and segmentation using the polar network. "Cartesian centers" refers to predicting the polar origins from the Cartesian network and then performing segmentation using the polar network. "Model centers" refers to using the center-point predictor to obtain polar origins. (Continued on the next page.)

| Liver Segmentation | | | | | |
|------------------------------------|---------------|--------|--------|--------|--------|
| Segm. net. | Method | DSC | mIoU | Prec. | Rec. |
| U-Net | Cartesian | 0.8976 | 0.8505 | 0.8997 | 0.9201 |
| | GT centers | 0.9553 | 0.9227 | 0.9595 | 0.9569 |
| | Cart. centers | 0.9302 | 0.8985 | 0.9279 | 0.9429 |
| | Model centers | 0.9125 | 0.8828 | 0.9108 | 0.9219 |
| Res-U-Net++ | Cartesian | 0.8908 | 0.8463 | 0.8936 | 0.9085 |
| | GT centers | 0.9548 | 0.9215 | 0.9492 | 0.9661 |
| | Cart. centers | 0.9219 | 0.8898 | 0.9119 | 0.9428 |
| | Model centers | 0.9109 | 0.8795 | 0.9009 | 0.9306 |
| DeepLabV3+ | Cartesian | 0.8868 | 0.8341 | 0.8995 | 0.8959 |
| | GT centers | 0.9518 | 0.9171 | 0.9547 | 0.9550 |
| | Cart. centers | 0.9253 | 0.8932 | 0.9244 | 0.9361 |
| | Model centers | 0.9092 | 0.8783 | 0.9075 | 0.9199 |
| Epicardial Fat Segmentation | | | | | |
| Segm. net. | Method | DSC | mIoU | Prec. | Rec. |
| U-Net | Cartesian | 0.7544 | 0.5812 | 0.7190 | 0.6949 |
| | GT centers | 0.8088 | 0.6607 | 0.7986 | 0.7675 |
| | Cart. centers | 0.7835 | 0.6227 | 0.7455 | 0.7208 |
| | Model centers | 0.7840 | 0.6252 | 0.7451 | 0.7302 |
| Res-U-Net++ | Cartesian | 0.3410 | 0.1743 | 0.2700 | 0.3294 |
| | GT centers | 0.8030 | 0.6827 | 0.7939 | 0.8043 |
| | Cart. centers | 0.5466 | 0.3980 | 0.5286 | 0.5066 |
| | Model centers | 0.7740 | 0.6140 | 0.7156 | 0.7453 |
| DeepLabV3+ | Cartesian | 0.6380 | 0.4246 | 0.5665 | 0.5940 |
| | GT centers | 0.6952 | 0.5123 | 0.6519 | 0.6779 |
| | Cart. centers | 0.6696 | 0.4716 | 0.5988 | 0.6454 |
| | Model centers | 0.6720 | 0.4779 | 0.6070 | 0.6488 |

Table 4.1: Results of our proposed approaches (continued).

| Dataset | Method | DSC | mIoU | Prec. | Rec. |
|---------|-------------------|---------------|---------------|---------------|---------------|
| Polyp | PraNet [68] | 0.8990 | 0.8490 | - | - |
| | FANet [69] | - | 0.8937 | 0.9401 | 0.9339 |
| | Our method | 0.9374 | 0.8977 | 0.9488 | 0.9368 |
| Lesion | DeepLabV3+ [53] | 0.8717 | 0.7984 | 0.8807 | 0.9068 |
| | DoubleU-Net [63] | 0.8962 | 0.8212 | 0.9459 | 0.8780 |
| | Our method | 0.9253 | 0.8743 | 0.9253 | 0.9464 |
| Liver | U-Net [51] | 0.8976 | 0.8505 | 0.8997 | 0.9201 |
| | KiU-Net 3D [70] | 0.9423 | 0.8946 | - | - |
| | Our method | 0.9302 | 0.8985 | 0.9279 | 0.9429 |
| EAT | U-Net. [51] | 0.7544 | 0.5812 | 0.7190 | 0.6949 |
| | Zhang et al. [71] | 0.9119 | 0.8425 | - | - |
| | Our method | 0.7840 | 0.6252 | 0.7451 | 0.7302 |

Table 4.2: A comparison between our method (approach with best results) and the state of the art on the same datasets.

heart. This task is more challenging for common models like U-Net and requires a more complex approach [71]. It is possible that combining these existing approaches, namely segmenting the pericardium first, with training on polar coordinates would lead to an improvement in the state of the art.

We also show that training on polar images leads to a significant improvement in segmentation performance when compared to training on Cartesian images for the same network architecture. Additionally, as shown in Figure 4.11, the polar network portions of our approach converge in much fewer epochs than the Cartesian networks. This is in part due to the location information being encoded in the image itself via the polar origin, and in part due to a possible data dimensionality reduction, allowing the network to more easily optimize the loss function. The polar networks are also more robust to low dataset sample size. This is especially important in biomedical image segmentation where the availability of large labeled datasets is often very limited.

Predicting the center point from the Cartesian model, while still an improvement over the plain Cartesian network, leads to worse results than those obtained by the center point predictor model. We conclude that segmentation is highly dependent on choosing the correct polar origin. This dependency is somewhat loosened by adding polar origin augmentation when training the polar network.

Figure 4.14 shows a random sampling of predictions from the polar network using the center point predictor for polar origins. Qualitatively, we conclude that the network achieves very good segmentation results, leading to a very high overlap with the target object. The network successfully segments both small and elliptical as well as large and unevenly shaped polyps. On the lesion images, the network predicts a smooth border when sometimes the actual border of the lesion is rough, as shown in the left-most example on Figure 4.14(b), however, the network still does a good job of delineating a lesion

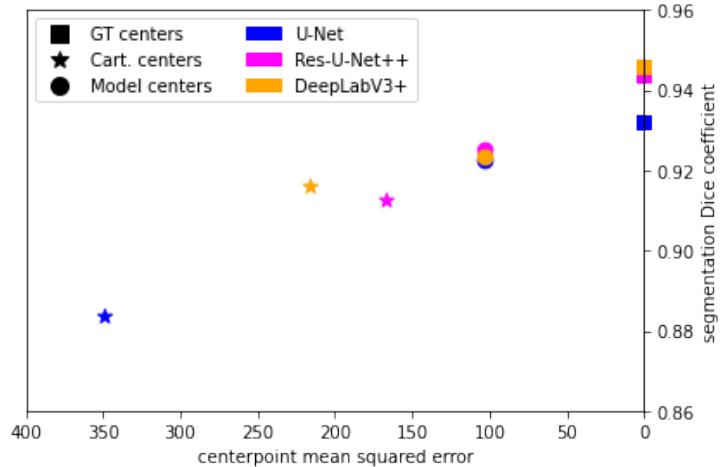


Figure 4.12: The relationship between mean squared errors of the centers used for the polar transformation and segmentation performance of the polar network on the lesion dataset. The mean squared errors are calculated compared to the ground-truth centers.

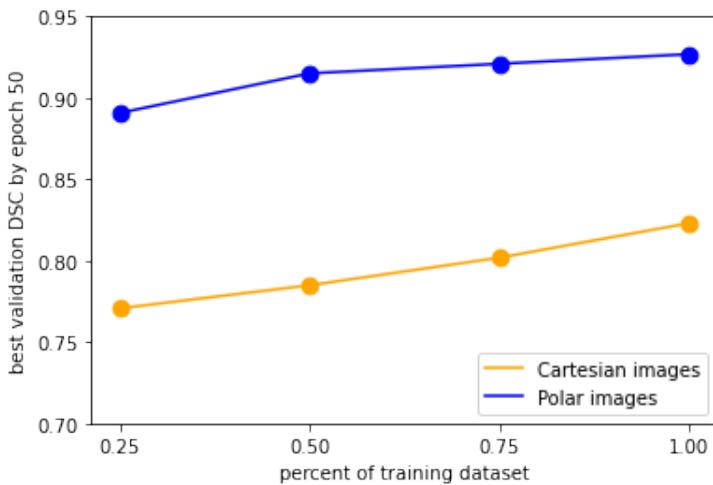


Figure 4.13: The best Dice coefficient by epoch 50 for models trained on subsets of the lesion training dataset.

border even when the color of the lesion is very similar to the surrounding skin. The network successfully predicts a liver border both when the liver is very large and very small on the image, showing good scale invariance, but sometimes under segments the liver when multiple connected components are needed. On the EAT dataset, the network successfully learns to segment EAT despite its highly discontinuous and sparse distribution. However, the network sometimes under segments EAT.

Finally, we also perform an ablation study shown in [Table 4.3](#). Training on the polar coordinates with the polar origins predicted from the cartesian network yields the largest performance improvement. Predicting the polar origin from the center point predictor as well as adding center point augmentation to the predictor play a roughly equally important role in the performance. Lastly, a small performance improvement is further achieved by using data augmentation when training the polar network.

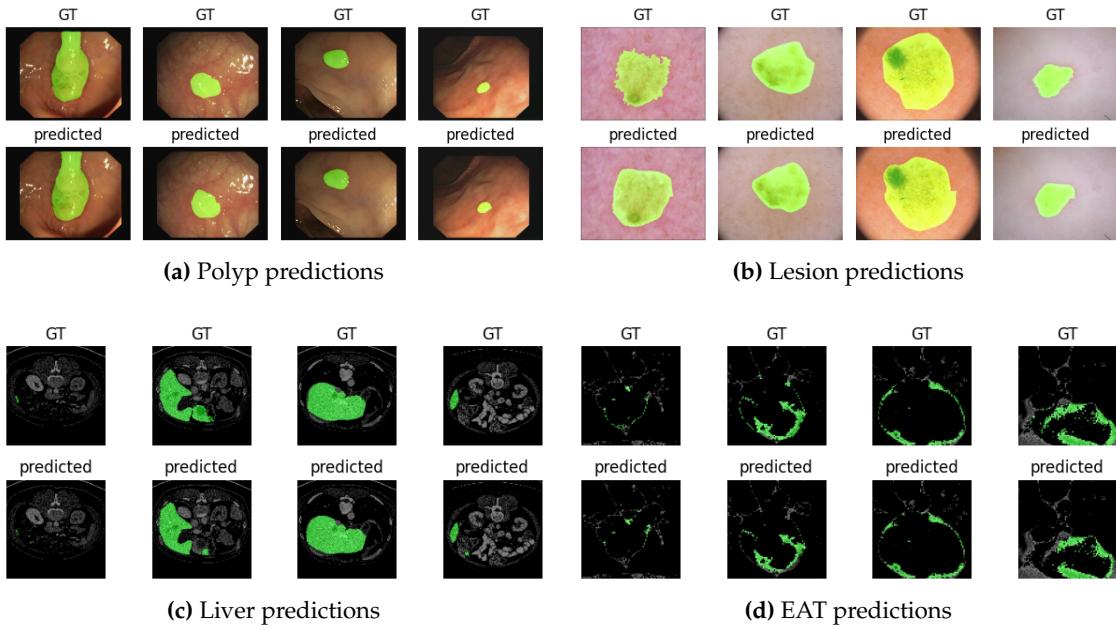


Figure 4.14: A random sampling of inverse polar transformed predictions from the polar network with the polar origins predicted from the center-point predictor for various datasets. The prediction is shown in green and overlaid on top of the original input image. EAT predictions (d) are cropped and zoomed to better show the predictions.

| Method | DSC | Difference |
|---------------------------------------|--------|------------|
| Cartesian | 0.8315 | - |
| Polar (Cartesian origins) | 0.8918 | +0.0603 |
| Polar (Centerpoint predictor) | 0.9094 | +0.0176 |
| + centerpoint augmentation | 0.9288 | +0.0194 |
| + polar network training augmentation | 0.9374 | +0.0086 |

Table 4.3: Ablation study of our approach for the polyp dataset.

A potential improvement of our method is to train a single neural network that combines the center-point predictor and the segmentation network and is trained end-to-end. In our approach, polar origins are always optimized towards the center of mass of the segmented object. Training an end-to-end network would allow the polar origins to be optimized for that specific segmentation task. Additionally, the center points could be obtained manually from experts, creating a user-guided segmentation approach similar to [57]. The center points could also be obtained by a more basic segmentation approach like thresholding or other traditional image processing method, leading to a possible reduction in the number of required neural network parameters to achieve good segmentation. Furthermore, in our experiments, we found that the segmentation is dependent on choosing the correct standard deviation of the generated heatmaps for training the center point predictor. An improvement to our method could be made by developing a method to automatically estimate the standard deviation from the training or validation data without needing to first train the center point predictor.

4.5 Conclusion

We explored training neural networks for biomedical image segmentation on polar transformations on images. We hypothesized that polar transformations would reduce the dimensionality of the input images, and allow the network to separately learn localization and fine segmentation of an object. We showed that training time improves when training on polar images for tasks where a single object which is roughly elliptical in shape or distribution needs to be segmented. Additionally, we show that training on polar images achieves state-of-the-art results on small datasets, and achieves near state-of-the-art results on larger datasets using generic low-parameter-count models like U-Net. We also noted that choosing the correct polar origin is essential for improving performance on polar images. Therefore, we proposed two different ways of obtaining the polar origin automatically from unlabeled input images. We trained a center-point predictor which predicts a heatmap to produce a polar origin, and showed that its performance is better than predicting the origin from a segmentation network trained on Cartesian images. We noted that sometimes our method under segments in examples where multiple objects need to be segmented.

rewrite
this using
formula-
tion

While our approach already produces state-of-the-art results in some cases, our results could be further improved. Our approach can be used as a pre-processing step for existing and future semantic segmentation methods that use neural networks to provide additional segmentation improvement. Therefore, it is possible that our approach could be used in a variety of different biomedical and non-medical segmentation applications.

Bibliography

- [1] R. Fosbinder and D. Orth, *Essentials of Radiologic Science*. Wolters Kluwer Health/Lippincott Williams & Wilkins, 2011, ISBN: 978-0-7817-7554-0.
- [2] S. Kamalian, M. H. Lev, and R. Gupta, “Computed tomography imaging and angiography – principles”, in *Handbook of Clinical Neurology*, vol. 135, Elsevier, 2016, pp. 3–20, ISBN: 978-0-444-53485-9. DOI: [10.1016/B978-0-444-53485-9.00001-5](https://doi.org/10.1016/B978-0-444-53485-9.00001-5). (visited on 09/29/2023).
- [3] A. A. Mahabadi, B. Balcer, I. Dykun, M. Forsting, T. Schlosser, G. Heusch, and T. Rassaf, “Cardiac computed tomography-derived epicardial fat volume and attenuation independently distinguish patients with and without myocardial infarction”, *PLOS ONE*, vol. 12, no. 8, M. W. Merx, Ed., e0183514, Aug. 2017, ISSN: 1932-6203. DOI: [10.1371/journal.pone.0183514](https://doi.org/10.1371/journal.pone.0183514). (visited on 09/29/2023).
- [4] M. Benčević, I. Galić, M. Habijan, and A. Pižurica, “Recent Progress in Epicardial and Pericardial Adipose Tissue Segmentation and Quantification Based on Deep Learning: A Systematic Review”, *Applied Sciences*, vol. 12, no. 10, p. 5217, May 2022, ISSN: 2076-3417. DOI: [10.3390/app12105217](https://doi.org/10.3390/app12105217). (visited on 09/29/2023).
- [5] L. Radl, Y. Jin, A. Pepe, J. Li, C. Gsaxner, F.-h. Zhao, and J. Egger, “AVT: Multicenter aortic vessel tree CTA dataset collection with ground truth segmentation masks”, *Data in Brief*, vol. 40, p. 107801, Feb. 2022, ISSN: 23523409. DOI: [10.1016/j.dib.2022.107801](https://doi.org/10.1016/j.dib.2022.107801). (visited on 10/05/2023).
- [6] E. Calabrese, J. E. Villanueva-Meyer, J. D. Rudie, A. M. Rauschecker, U. Baid, S. Bakas, S. Cha, J. T. Mongan, and C. P. Hess, “The University of California San Francisco Preoperative Diffuse Glioma MRI (UCSF-PDGM) Dataset”, *Radiology: Artificial Intelligence*, vol. 4, no. 6, e220058, Nov. 2022, ISSN: 2638-6100. DOI: [10.1148/ryai.220058](https://doi.org/10.1148/ryai.220058). arXiv: [2109.00356 \[cs, eess\]](https://arxiv.org/abs/2109.00356). (visited on 10/05/2023).
- [7] H. Q. Nguyen, K. Lam, L. T. Le, H. H. Pham, D. Q. Tran, D. B. Nguyen, D. D. Le, C. M. Pham, H. T. T. Tong, D. H. Dinh, C. D. Do, L. T. Doan, C. N. Nguyen, B. T. Nguyen, Q. V. Nguyen, A. D. Hoang, H. N. Phan, A. T. Nguyen, P. H. Ho, D. T. Ngo, N. T. Nguyen, N. T. Nguyen, M. Dao, and V. Vu, *VinDr-CXR: An open dataset of chest X-rays with radiologist's annotations*, 2020. arXiv: [2012.15029 \[eess.IV\]](https://arxiv.org/abs/2012.15029).
- [8] J. Irvin, P. Rajpurkar, M. Ko, Y. Yu, S. Ciurea-Ilcus, C. Chute, H. Marklund, B. Haghgoo, R. Ball, K. Shpanskaya, J. Seekins, D. A. Mong, S. S. Halabi, J. K. Sandberg, R. Jones, D. B. Larson, C. P. Langlotz, B. N. Patel, M. P. Lungren, and A. Y. Ng, *CheXpert: A Large Chest Radiograph Dataset with Uncertainty Labels and Expert Comparison*, Jan. 2019. arXiv: [1901.07031 \[cs, eess\]](https://arxiv.org/abs/1901.07031). (visited on 10/02/2023).
- [9] V. Rotemberg, N. Kurtansky, B. Betz-Stablein, L. Caffery, E. Chousakos, N. Codella, M. Combalia, S. Dusza, P. Guitera, D. Gutman, A. Halpern, B. Helba, H. Kittler, K. Kose, S. Langer, K. Lioprys, J. Malvehy, S. Musthaq, J. Nanda, O. Reiter, G. Shih, A. Stratigos, P. Tschandl, J. Weber, and H. P. Soyer, “A patient-centric dataset of images and metadata for identifying melanomas using clinical context”, *Scientific Data*, vol. 8, no. 1, p. 34, Jan. 2021, ISSN: 2052-4463. DOI: [10.1038/s41597-021-00815-z](https://doi.org/10.1038/s41597-021-00815-z). (visited on 10/02/2023).

- [10] A. U. M. Khan, A. Torelli, I. Wolf, and N. Gretz, "AutoCellSeg: Robust automatic colony forming unit (CFU)/cell analysis using adaptive image segmentation and easy-to-use post-editing techniques", *Scientific Reports*, vol. 8, no. 1, p. 7302, May 2018, ISSN: 2045-2322. DOI: [10.1038/s41598-018-24916-9](https://doi.org/10.1038/s41598-018-24916-9). (visited on 10/16/2023).
- [11] A. Jha, H. Yang, R. Deng, M. E. Kapp, A. B. Fogo, and Y. Huo, "Instance segmentation for whole slide imaging: End-to-end or detect-then-segment", *Journal of Medical Imaging*, vol. 8, no. 01, Jan. 2021, ISSN: 2329-4302. DOI: [10.1117/1.JMI.8.1.014001](https://doi.org/10.1117/1.JMI.8.1.014001). (visited on 10/02/2023).
- [12] R. Adams and L. Bischof, "Seeded region growing", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 16, no. 6, pp. 641–647, 1994. DOI: [10.1109/34.295913](https://doi.org/10.1109/34.295913).
- [13] M. Kass, A. Witkin, and D. Terzopoulos, "Snakes: Active contour models", *International Journal of Computer Vision*, vol. 1, no. 4, pp. 321–331, Jan. 1988. DOI: [10.1007/bf00133570](https://doi.org/10.1007/bf00133570).
- [14] T. Rohlfing, R. Brandt, R. Menzel, D. B. Russakoff, and C. R. Maurer, "Quo vadis, atlas-based segmentation?", in *Handbook of Biomedical Image Analysis*, Springer US, 2005, pp. 435–486. DOI: [10.1007/0-306-48608-3_11](https://doi.org/10.1007/0-306-48608-3_11).
- [15] M. Sinclair, A. Schuh, K. Hahn, K. Petersen, Y. Bai, J. Batten, M. Schaap, and B. Glocker, "Atlas-ISTN: Joint segmentation, registration and atlas construction with image-and-spatial transformer networks", *Medical Image Analysis*, vol. 78, p. 102383, May 2022, ISSN: 13618415. DOI: [10.1016/j.media.2022.102383](https://doi.org/10.1016/j.media.2022.102383). (visited on 10/03/2023).
- [16] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, ser. Springer Series in Statistics. New York, NY: Springer New York, 2009, ISBN: 978-0-387-84857-0 978-0-387-84858-7. DOI: [10.1007/978-0-387-84858-7](https://doi.org/10.1007/978-0-387-84858-7). (visited on 10/21/2023).
- [17] K. He, G. Gkioxari, P. Dollar, and R. Girshick, "Mask R-CNN", in *2017 IEEE International Conference on Computer Vision (ICCV)*, Venice: IEEE, Oct. 2017, pp. 2980–2988, ISBN: 978-1-5386-1032-9. DOI: [10.1109/ICCV.2017.322](https://doi.org/10.1109/ICCV.2017.322). (visited on 10/09/2023).
- [18] V. Dumoulin and F. Visin, *A guide to convolution arithmetic for deep learning*, Jan. 2018. arXiv: [1603.07285 \[cs, stat\]](https://arxiv.org/abs/1603.07285). (visited on 10/19/2023).
- [19] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition", *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov./1998, ISSN: 00189219. DOI: [10.1109/5.726791](https://doi.org/10.1109/5.726791). (visited on 10/05/2023).
- [20] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation", in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.
- [21] O. Ronneberger, P. Fischer, and T. Brox, *U-Net: Convolutional Networks for Biomedical Image Segmentation*, May 2015. arXiv: [1505.04597 \[cs\]](https://arxiv.org/abs/1505.04597). (visited on 10/06/2023).
- [22] F. Isensee, P. F. Jaeger, S. A. A. Kohl, J. Petersen, and K. H. Maier-Hein, "nnU-Net: A self-configuring method for deep learning-based biomedical image segmentation", *Nature Methods*, vol. 18, no. 2, pp. 203–211, Feb. 2021, ISSN: 1548-7091, 1548-7105. DOI: [10.1038/s41592-020-01008-z](https://doi.org/10.1038/s41592-020-01008-z). (visited on 10/06/2023).
- [23] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, *High-resolution image synthesis with latent diffusion models*, 2021. arXiv: [2112.10752 \[cs.CV\]](https://arxiv.org/abs/2112.10752).

- [24] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang, "UNet++: Redesigning skip connections to exploit multiscale features in image segmentation", *IEEE Transactions on Medical Imaging*, 2019.
- [25] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network", in *CVPR*, 2017.
- [26] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, *Rethinking atrous convolution for semantic image segmentation*, 2017. arXiv: [1706.05587 \[cs.CV\]](https://arxiv.org/abs/1706.05587).
- [27] T. Fan, G. Wang, Y. Li, and H. Wang, "MA-Net: A multi-scale attention network for liver and tumor segmentation", *IEEE access : practical innovations, open solutions*, vol. 8, pp. 179 656–179 665, 2020. DOI: [10.1109/ACCESS.2020.3025372](https://doi.org/10.1109/ACCESS.2020.3025372).
- [28] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need", in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30, Curran Associates, Inc., 2017.
- [29] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale", in *International Conference on Learning Representations*, 2021.
- [30] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr, and L. Zhang, "Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers", in *CVPR*, 2021.
- [31] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows", in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [32] Z. Chen, Y. Duan, W. Wang, J. He, T. Lu, J. Dai, and Y. Qiao, "Vision transformer adapter for dense predictions", *arXiv preprint arXiv:2205.08534*, 2022. arXiv: [2205.08534](https://arxiv.org/abs/2205.08534).
- [33] S. Shalev-Shwartz and S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*, 1st ed. Cambridge University Press, May 2014, ISBN: 978-1-107-05713-5 978-1-107-29801-9. DOI: [10.1017/CBO9781107298019](https://doi.org/10.1017/CBO9781107298019). (visited on 12/01/2023).
- [34] J. Deng, W. Dong, R. Socher, L.-J. Li, Kai Li, and Li Fei-Fei, "ImageNet: A large-scale hierarchical image database", in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, Miami, FL: IEEE, Jun. 2009, pp. 248–255, ISBN: 978-1-4244-3992-8. DOI: [10.1109/CVPR.2009.5206848](https://doi.org/10.1109/CVPR.2009.5206848). (visited on 12/01/2023).
- [35] J. Wasserthal, H.-C. Breit, M. T. Meyer, M. Pradella, D. Hinck, A. W. Sauter, T. Heyne, D. T. Boll, J. Cyriac, S. Yang, M. Bach, and M. Segeroth, "TotalSegmentator: Robust Segmentation of 104 Anatomic Structures in CT Images", *Radiology: Artificial Intelligence*, vol. 5, no. 5, e230024, Sep. 2023, ISSN: 2638-6100. DOI: [10.1148/ryai.230024](https://doi.org/10.1148/ryai.230024). (visited on 11/27/2023).
- [36] MONAI Consortium, *MONAI: Medical Open Network for AI*, Zenodo, Oct. 2023. DOI: [10.5281/ZENODO.4323058](https://doi.org/10.5281/ZENODO.4323058). (visited on 11/28/2023).
- [37] A. Myronenko, D. Yang, Y. He, and D. Xu, *Automated 3D Segmentation of Kidneys and Tumors in MICCAI KiTS 2023 Challenge*, Oct. 2023. arXiv: [2310.04110 \[cs\]](https://arxiv.org/abs/2310.04110). (visited on 11/28/2023).

- [38] A. Torralba and A. A. Efros, "Unbiased look at dataset bias", in *CVPR 2011*, Colorado Springs, CO, USA: IEEE, Jun. 2011, pp. 1521–1528, ISBN: 978-1-4577-0394-2. DOI: [10.1109/CVPR.2011.5995347](https://doi.org/10.1109/CVPR.2011.5995347). (visited on 11/28/2023).
- [39] L. Liu, W. Lin, L. Wu, Y. Yu, and M. Y. Yang, *Unsupervised Deep Domain Adaptation for Pedestrian Detection*, Feb. 2018. arXiv: [1802.03269 \[cs\]](https://arxiv.org/abs/1802.03269). (visited on 11/28/2023).
- [40] Y. Ganin and V. Lempitsky, "Unsupervised domain adaptation by backpropagation", in *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*, ser. ICML'15, Lille, France: JMLR.org, 2015, pp. 1180–1189.
- [41] R. Bermudez-Chacon, P. Marquez-Neila, M. Salzmann, and P. Fua, "A domain-adaptive two-stream U-Net for electron microscopy image segmentation", in *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, Washington, DC: IEEE, Apr. 2018, pp. 400–404, ISBN: 978-1-5386-3636-7. DOI: [10.1109/ISBI.2018.8363602](https://doi.org/10.1109/ISBI.2018.8363602). (visited on 11/28/2023).
- [42] M. Noroozi and P. Favaro, "Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles", in *Computer Vision – ECCV 2016*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds., vol. 9910, Cham: Springer International Publishing, 2016, pp. 69–84, ISBN: 978-3-319-46465-7 978-3-319-46466-4. DOI: [10.1007/978-3-319-46466-4_5](https://doi.org/10.1007/978-3-319-46466-4_5). (visited on 02/16/2022).
- [43] A. N. Carr, Q. Berthet, M. Blondel, O. Teboul, and N. Zeghidour, *Shuffle to Learn: Self-supervised learning from permutations via differentiable ranking*, 2021.
- [44] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A Simple Framework for Contrastive Learning of Visual Representations", *arXiv:2002.05709 [cs, stat]*, Jun. 2020. arXiv: [2002.05709 \[cs, stat\]](https://arxiv.org/abs/2002.05709). (visited on 02/10/2022).
- [45] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning", *arXiv preprint arXiv:1911.05722*, 2019. arXiv: [1911.05722](https://arxiv.org/abs/1911.05722).
- [46] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, *Generative Adversarial Networks*, Jun. 2014. arXiv: [1406.2661 \[cs, stat\]](https://arxiv.org/abs/1406.2661). (visited on 12/01/2023).
- [47] H.-C. Shin, N. A. Tenenholtz, J. K. Rogers, C. G. Schwarz, M. L. Senjem, J. L. Gunter, K. P. Andriole, and M. Michalski, "Medical Image Synthesis for Data Augmentation and Anonymization Using Generative Adversarial Networks", in *Simulation and Synthesis in Medical Imaging*, A. Gooya, O. Goksel, I. Oguz, and N. Burgos, Eds., vol. 11037, Cham: Springer International Publishing, 2018, pp. 1–11, ISBN: 978-3-030-00535-1 978-3-030-00536-8. DOI: [10.1007/978-3-030-00536-8_1](https://doi.org/10.1007/978-3-030-00536-8_1). (visited on 12/01/2023).
- [48] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models", *arXiv preprint arxiv:2006.11239*, 2020. arXiv: [2006.11239](https://arxiv.org/abs/2006.11239).
- [49] F. Khader, G. Müller-Franzes, S. Tayebi Arasteh, T. Han, C. Haarburger, M. Schulze-Hagen, P. Schad, S. Engelhardt, B. Baefler, S. Foersch, J. Stegmaier, C. Kuhl, S. Nebelung, J. N. Kather, and D. Truhn, "Denoising diffusion probabilistic models for 3D medical image generation", *Scientific Reports*, vol. 13, no. 1, p. 7303, May 2023, ISSN: 2045-2322. DOI: [10.1038/s41598-023-34341-2](https://doi.org/10.1038/s41598-023-34341-2). (visited on 12/01/2023).
- [50] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*, ser. Adaptive Computation and Machine Learning. Cambridge, Massachusetts: The MIT Press, 2016, ISBN: 978-0-262-03561-3.

- [51] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation", in *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, N. Navab, J. Hornegger, W. M. Wells, and A. F. Frangi, Eds., vol. 9351, Cham: Springer International Publishing, 2015, pp. 234–241, ISBN: 978-3-319-24573-7 978-3-319-24574-4. DOI: [10.1007/978-3-319-24574-4_28](https://doi.org/10.1007/978-3-319-24574-4_28). (visited on 08/23/2021).
- [52] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition", in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA: IEEE, Jun. 2016, pp. 770–778, ISBN: 978-1-4673-8851-1. DOI: [10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90). (visited on 08/23/2021).
- [53] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation", in *Computer Vision – ECCV 2018*, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds., vol. 11211, Cham: Springer International Publishing, 2018, pp. 833–851, ISBN: 978-3-030-01233-5 978-3-030-01234-2. DOI: [10.1007/978-3-030-01234-2_49](https://doi.org/10.1007/978-3-030-01234-2_49). (visited on 08/23/2021).
- [54] Q. Liu, X. Hong, W. Ke, Z. Chen, and B. Zou, "DDNet: Cartesian-polar Dual-domain Network for the Joint Optic Disc and Cup Segmentation", *arXiv:1904.08773 [cs]*, Apr. 2019. arXiv: [1904.08773 \[cs\]](https://arxiv.org/abs/1904.08773). (visited on 04/01/2021).
- [55] H. Salehinejad, S. Valaee, T. Dowdell, and J. Barfett, "Image Augmentation Using Radial Transform for Training Deep Neural Networks", in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Calgary, AB: IEEE, Apr. 2018, pp. 3016–3020, ISBN: 978-1-5386-4658-8. DOI: [10.1109/ICASSP.2018.8462241](https://doi.org/10.1109/ICASSP.2018.8462241). (visited on 08/23/2021).
- [56] J. Kim, W. Jung, H. Kim, and J. Lee, "CyCNN: A Rotation Invariant CNN using Polar Mapping and Cylindrical Convolution Layers", *arXiv:2007.10588 [cs, eess]*, Jul. 2020. arXiv: [2007.10588 \[cs, eess\]](https://arxiv.org/abs/2007.10588). (visited on 04/08/2021).
- [57] B.-S. Kim, J.-Y. Sun, S.-W. Kim, M.-C. Kang, and S.-J. Ko, "CNN-based UGS method using Cartesian-to-polar coordinate transformation", *Electronics Letters*, vol. 54, no. 23, pp. 1321–1322, Nov. 2018, ISSN: 0013-5194, 1350-911X. DOI: [10.1049/el.2018.5051](https://doi.org/10.1049/el.2018.5051). (visited on 04/01/2021).
- [58] C. Esteves, C. Allen-Blanchette, X. Zhou, and K. Daniilidis, "Polar Transformer Networks", *arXiv:1709.01889 [cs]*, Feb. 2018. arXiv: [1709.01889 \[cs\]](https://arxiv.org/abs/1709.01889). (visited on 04/01/2021).
- [59] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu, "Spatial Transformer Networks", *arXiv:1506.02025 [cs]*, Feb. 2016. arXiv: [1506.02025 \[cs\]](https://arxiv.org/abs/1506.02025). (visited on 04/26/2021).
- [60] A. Newell, K. Yang, and J. Deng, "Stacked Hourglass Networks for Human Pose Estimation", in *Computer Vision – ECCV 2016*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds., vol. 9912, Cham: Springer International Publishing, 2016, pp. 483–499, ISBN: 978-3-319-46483-1 978-3-319-46484-8. DOI: [10.1007/978-3-319-46484-8_29](https://doi.org/10.1007/978-3-319-46484-8_29). (visited on 08/23/2021).
- [61] A. Buslaev, V. I. Iglovikov, E. Khvedchenya, A. Parinov, M. Druzhinin, and A. A. Kalinin, "Albumentations: Fast and flexible image augmentations", *Information-an International Interdisciplinary Journal*, vol. 11, no. 2, 2020, ISSN: 2078-2489. DOI: [10.3390/info11020125](https://doi.org/10.3390/info11020125).

- [62] J. Bernal, F. J. Sánchez, G. Fernández-Esparrach, D. Gil, C. Rodríguez, and F. Vilar-iño, "WM-DOVA maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians", *Computerized Medical Imaging and Graphics*, vol. 43, pp. 99–111, Jul. 2015, ISSN: 08956111. DOI: [10.1016/j.compmedimag.2015.02.007](https://doi.org/10.1016/j.compmedimag.2015.02.007). (visited on 04/13/2021).
- [63] D. Jha, M. A. Riegler, D. Johansen, P. Halvorsen, and H. D. Johansen, "DoubleUNet: A Deep Convolutional Neural Network for Medical Image Segmentation", in *2020 IEEE 33rd International Symposium on Computer-Based Medical Systems (CBMS)*, Rochester, MN, USA: IEEE, Jul. 2020, pp. 558–564, ISBN: 978-1-72819-429-5. DOI: [10.1109/CBMS49503.2020.00111](https://doi.org/10.1109/CBMS49503.2020.00111). (visited on 08/23/2021).
- [64] P. Bilic, P. F. Christ, E. Vorontsov, G. Chlebus, H. Chen, Q. Dou, C.-W. Fu, X. Han, P.-A. Heng, J. Hesser, S. Kadoury, T. Konopczynski, M. Le, C. Li, X. Li, J. Lipková, J. Lowengrub, H. Meine, J. H. Moltz, C. Pal, M. Piraud, X. Qi, J. Qi, M. Rempfler, K. Roth, A. Schenk, A. Sekuboyina, E. Vorontsov, P. Zhou, C. Hülsemeyer, M. Beetz, F. Ettlinger, F. Gruen, G. Kaassis, F. Lohöfer, R. Braren, J. Holch, F. Hofmann, W. Sommer, V. Heinemann, C. Jacobs, G. E. H. Mamani, B. van Ginneken, G. Chartrand, A. Tang, M. Drozdzał, A. Ben-Cohen, E. Klang, M. M. Amitai, E. Konen, H. Greenspan, J. Moreau, A. Hostettler, L. Soler, R. Vivanti, A. Szekeskin, N. Lev-Cohain, J. Sosna, L. Joskowicz, and B. H. Menze, "The Liver Tumor Segmentation Benchmark (LiTS)", *arXiv:1901.04056 [cs]*, Jan. 2019. arXiv: [1901.04056 \[cs\]](https://arxiv.org/abs/1901.04056). (visited on 04/13/2021).
- [65] N. C. F. Codella, D. Gutman, M. E. Celebi, B. Helba, M. A. Marchetti, S. W. Dusza, A. Kalloo, K. Liopyris, N. Mishra, H. Kittler, and A. Halpern, "Skin lesion analysis toward melanoma detection: A challenge at the 2017 International symposium on biomedical imaging (ISBI), hosted by the international skin imaging collaboration (ISIC)", in *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, Washington, DC: IEEE, Apr. 2018, pp. 168–172, ISBN: 978-1-5386-3636-7. DOI: [10.1109/ISBI.2018.8363547](https://doi.org/10.1109/ISBI.2018.8363547). (visited on 08/23/2021).
- [66] P. Tschandl, C. Rosendahl, and H. Kittler, "The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions", *Scientific Data*, vol. 5, no. 1, p. 180161, Dec. 2018, ISSN: 2052-4463. DOI: [10.1038/sdata.2018.161](https://doi.org/10.1038/sdata.2018.161). (visited on 04/13/2021).
- [67] É. Rodrigues, F. Morais, N. Morais, L. Conci, L. Neto, and A. Conci, "A novel approach for the automated segmentation and volume quantification of cardiac fats on computed tomography", *Computer Methods and Programs in Biomedicine*, vol. 123, pp. 109–128, Jan. 2016, ISSN: 01692607. DOI: [10.1016/j.cmpb.2015.09.017](https://doi.org/10.1016/j.cmpb.2015.09.017). (visited on 04/13/2021).
- [68] D.-P. Fan, G.-P. Ji, T. Zhou, G. Chen, H. Fu, J. Shen, and L. Shao, "PraNet: Parallel Reverse Attention Network for Polyp Segmentation", in *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*, A. L. Martel, P. Abolmaesumi, D. Stoyanov, D. Mateus, M. A. Zuluaga, S. K. Zhou, D. Racoceanu, and L. Joskowicz, Eds., vol. 12266, Cham: Springer International Publishing, 2020, pp. 263–273, ISBN: 978-3-030-59724-5 978-3-030-59725-2. DOI: [10.1007/978-3-030-59725-2_26](https://doi.org/10.1007/978-3-030-59725-2_26). (visited on 08/23/2021).
- [69] N. K. Tomar, D. Jha, M. A. Riegler, H. D. Johansen, D. Johansen, J. Rittscher, P. Halvorsen, and S. Ali, "FANet: A Feedback Attention Network for Improved Biomedical Image Segmentation", *arXiv:2103.17235 [cs, eess]*, Mar. 2021. arXiv: [2103.17235 \[cs, eess\]](https://arxiv.org/abs/2103.17235). (visited on 04/19/2021).

- [70] J. M. J. Valanarasu, V. A. Sindagi, I. Hacihaliloglu, and V. M. Patel, "KiU-Net: Over-complete Convolutional Architectures for Biomedical Image and Volumetric Segmentation", *arXiv:2010.01663 [cs, eess]*, Oct. 2020. arXiv: [2010.01663 \[cs, eess\]](https://arxiv.org/abs/2010.01663). (visited on 04/22/2021).
- [71] Q. Zhang, J. Zhou, B. Zhang, W. Jia, and E. Wu, "Automatic Epicardial Fat Segmentation and Quantification of CT Scans Using Dual U-Nets With a Morphological Processing Layer", *IEEE Access*, vol. 8, pp. 128 032–128 041, 2020, ISSN: 2169-3536. DOI: [10.1109/ACCESS.2020.3008190](https://doi.org/10.1109/ACCESS.2020.3008190). (visited on 04/22/2021).