

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

October 28, 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	VI
List of Tables	VIII
List of Abbreviations	IX
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
Bibliography	29

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 up-sampled 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

List of Tables

2.1	Approximate Hounsfield values of various tissues [1], [2].	4
-----	--	---

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.

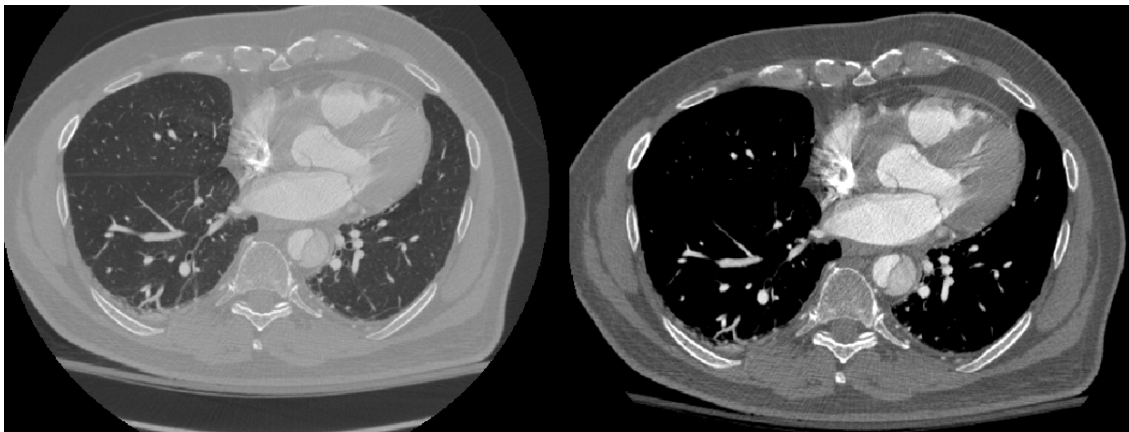


Figure 2.1: A cardiac CTA in its full range (left) and windowed (right). [5]

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 2.2](#).

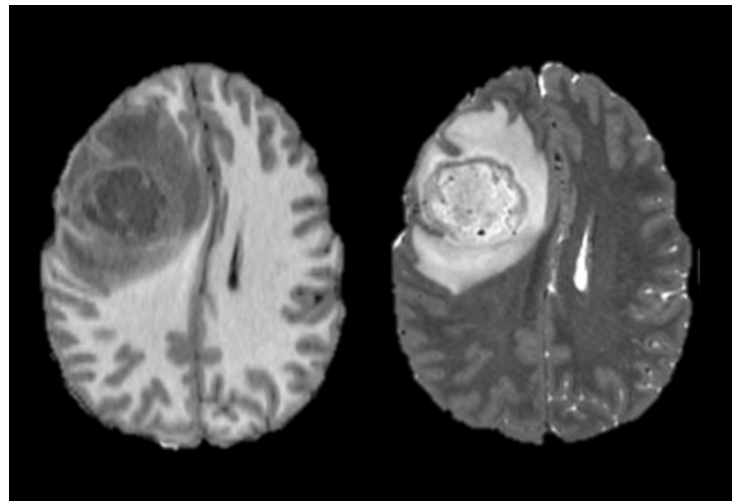


Figure 2.2: An example T1-weighted MRI (left) and a T2-weighted MRI (right) showing a diffuse glioma. [6]

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, downsampled, 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.

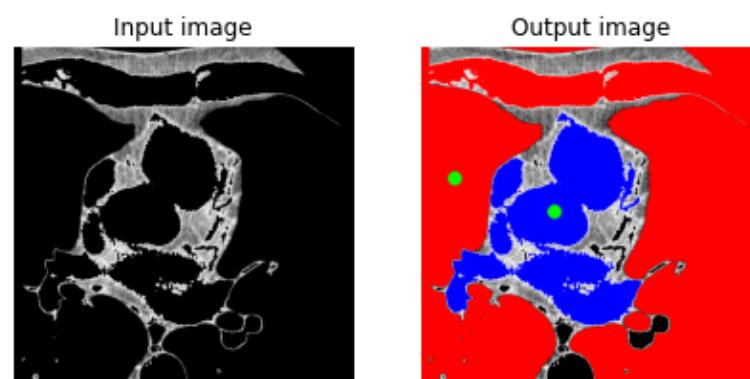


Figure 2.3: A demonstration of region growing for delineating the internal and external areas of the pericardium on a CT slice, set to the adipose tissue intensity range. The left image is the original input, and the right image depicts the segmented outcome with the heart exterior in red and the interior 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.

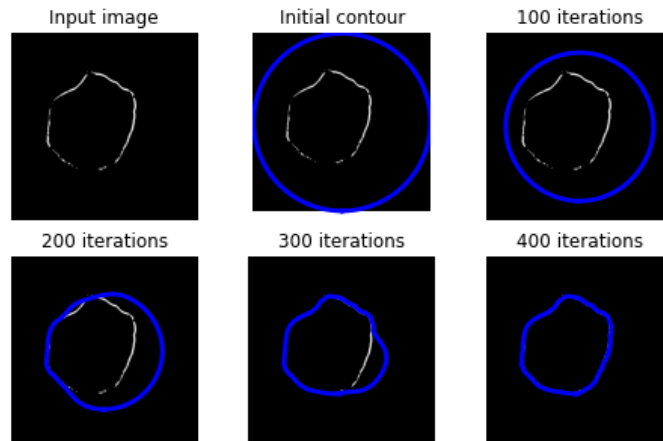


Figure 2.4: A demonstration of employing active contours to complete the absent segments of the pericardium line, displayed in white. The contour, illustrated in blue, starts as a complete circle surrounding the image. With every iteration, the contour adapts more closely to the image’s shape. [4]

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:

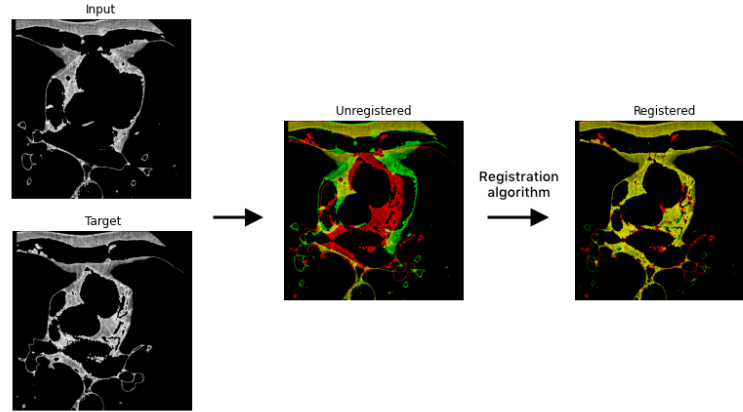


Figure 2.5: A schematic representation of the registration procedure. Initially, input and target images are chosen. Throughout the registration phase, the input image (shown in green) undergoes deformation to align with the fixed target image (shown in red). [4]

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

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.

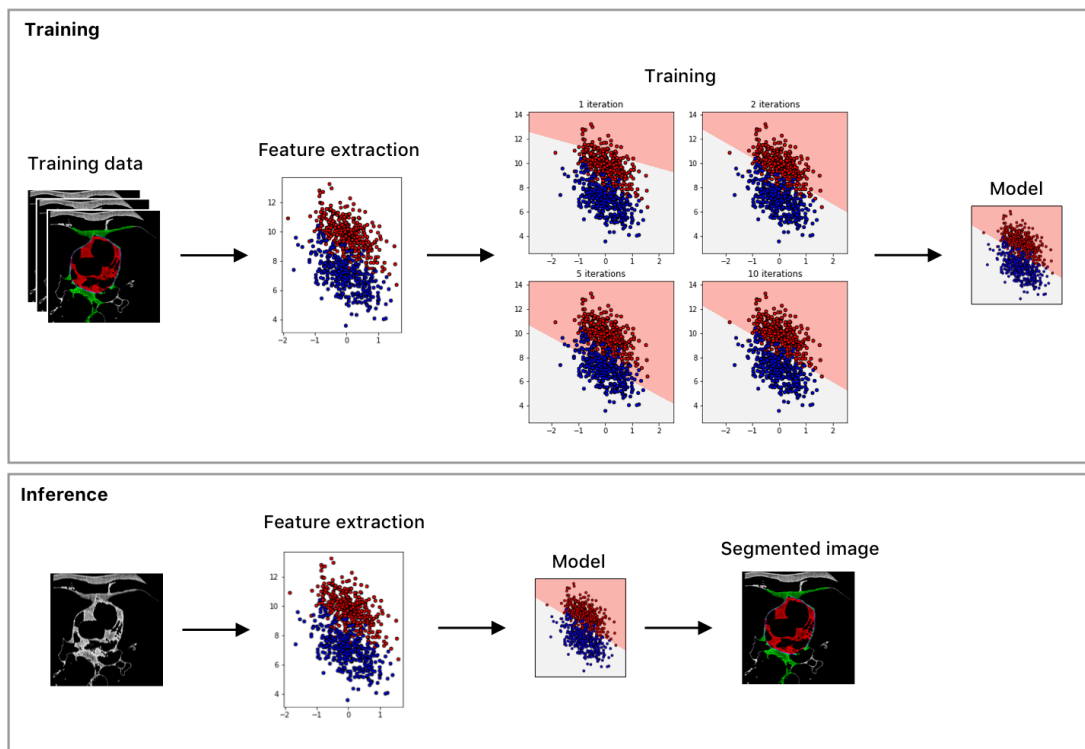


Figure 2.6: A schematic of a supervised linear classifier in a machine learning workflow. The upper section illustrates the training phase. Here, features are color-coded according to their known class from training data, depicted in red and blue. The parameters of the decision boundary, which demarcates 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 $x = [x_0 \ x_1 \ \cdots \ x_{n-1}]$ and the weights column vector w is calculated and added with the bias vector b :

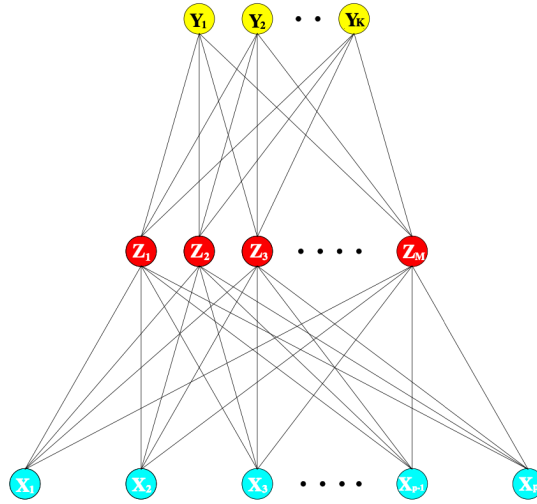


Figure 2.7: A small neural network with three layers, wherein each layer is connected to every neuron in the next layer. [16]

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

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:

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

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

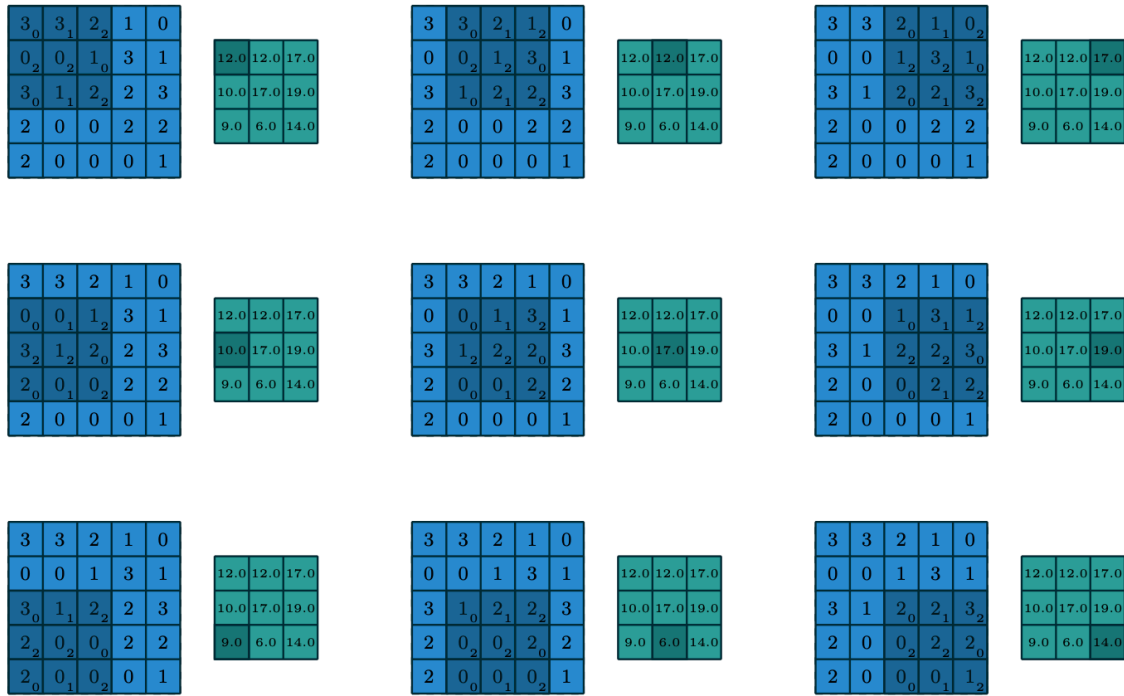


Figure 2.8: A visual depiction of a convolution procedure step-by-step. In each step, the kernel slides over the image (shown in blue). The overlapping elements between the kernel and the image are multiplied and then summed to produce a value in the resultant image (shown in green). [18]

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](#).

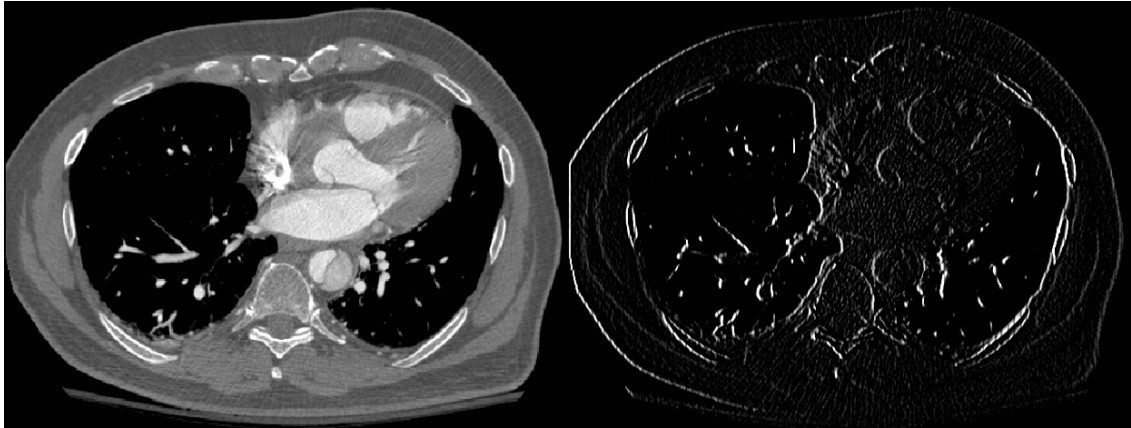


Figure 2.9: An example of an input image (left) convolved with the Prewitt operator (right). Note that vertical edges are accentuated in the convolution result.

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:

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

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:

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

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

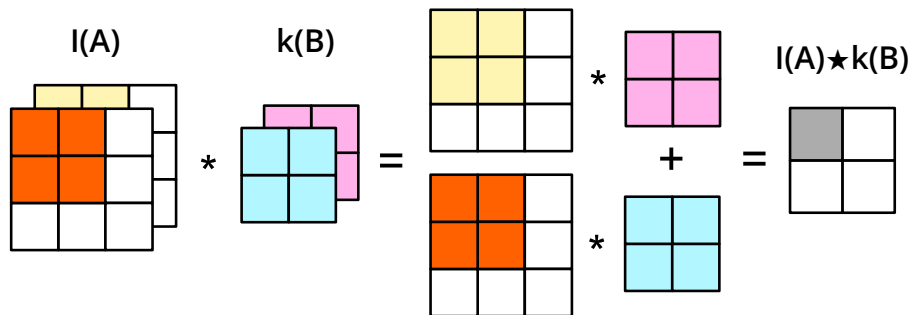


Figure 2.10: A view of one step of a single convolution operation inside a CNN layer. The layer performs multiple convolutions, each with a different kernel that has an equal number of channels as the input image. In each step, the whole kernel slides over the width and height of the image, and the overlapping 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](#).

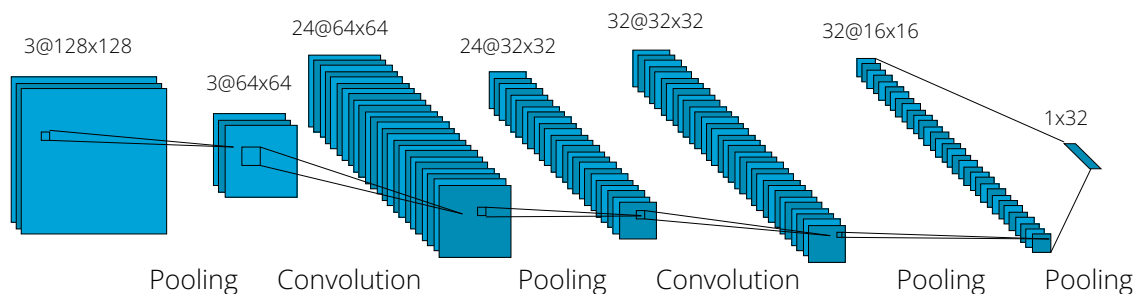


Figure 2.11: A typical architecture of a CNN encoder. The encoder consists of consecutive convolutional and pooling layers that gradually increase the feature map depth and decrease its width and height. The result is a map of features that tells the decoder what features are on the image but does not 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 upsampling.

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](#).

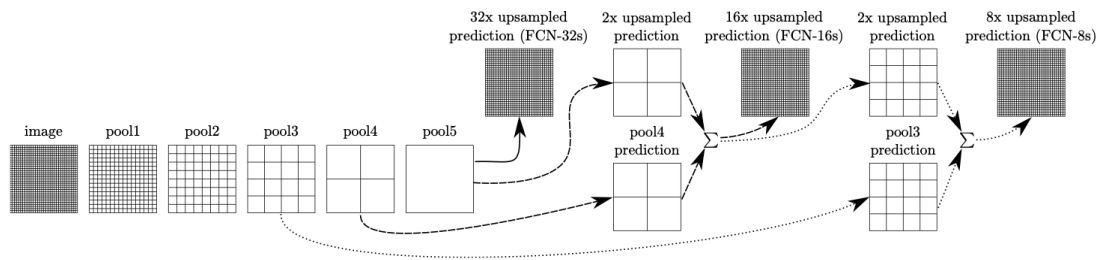


Figure 2.12: A diagram of how FCN forms predictions based on the output of different encoder layers. Encoder layers are shown on the left and the grid represents the coarseness of the feature map. The maps are combined at three different levels to produce three predictions. Each prediction is compared 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 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](#).

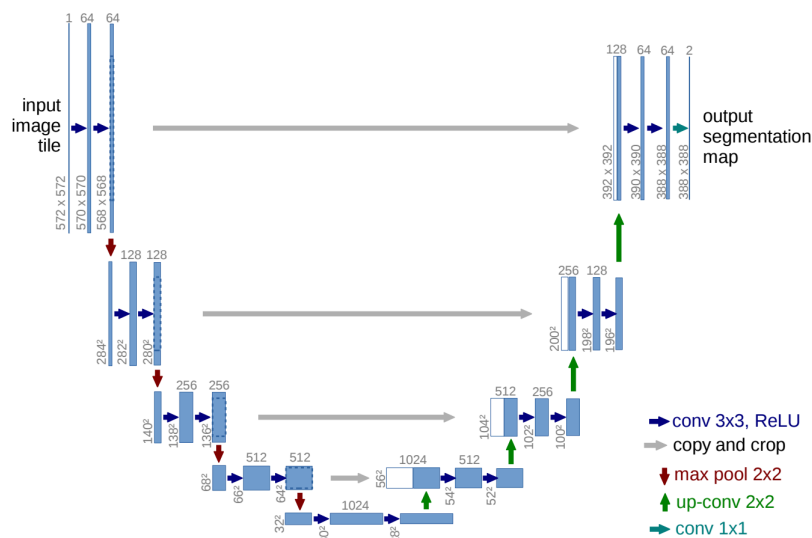


Figure 2.13: A diagram of the U-Net model. The output of each layer of the encoder is concatenated to the input of its corresponding layer in the decoder. [21]

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

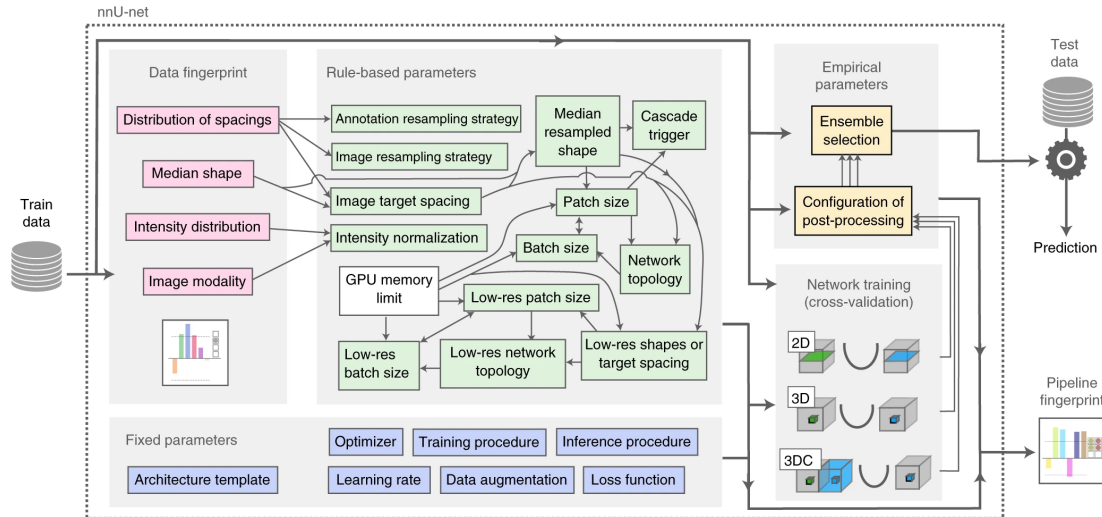


Figure 2.14: A diagram of the nnU-Net procedure of creating a training configuration. [22]

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 reimagined 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 reimagined 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.

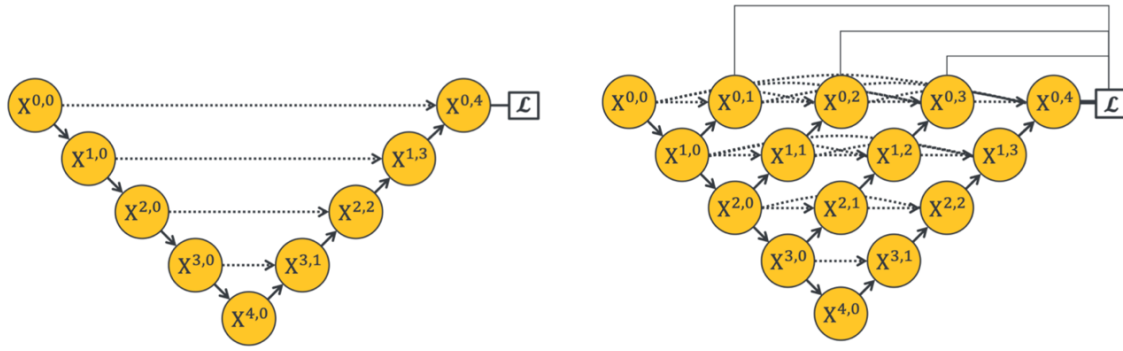


Figure 2.15: A comparison between U-Net (left) and U-Net++ (right). Each node in the graph represents a convolutional layer. The dashed arrows represent skip connections, while full arrows are downsampling or upsampling operations. [24]

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.

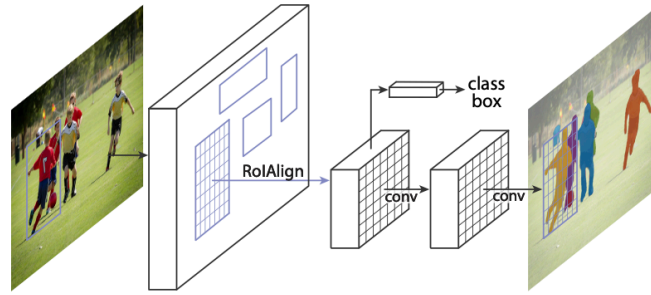


Figure 2.16: A diagram of the Mask R-CNN architecture. Two parallel decoder branches are used to achieve segmentation and object detection simultaneously. [17]

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.

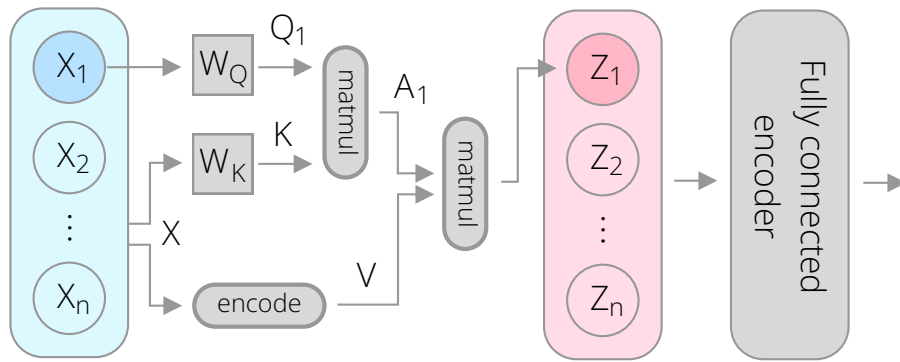


Figure 2.17: A visualization of a single encoder layer in a transformer network. This shows the encoding process for one element of the input sequence.

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 `citattnAllYouNeed`, 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

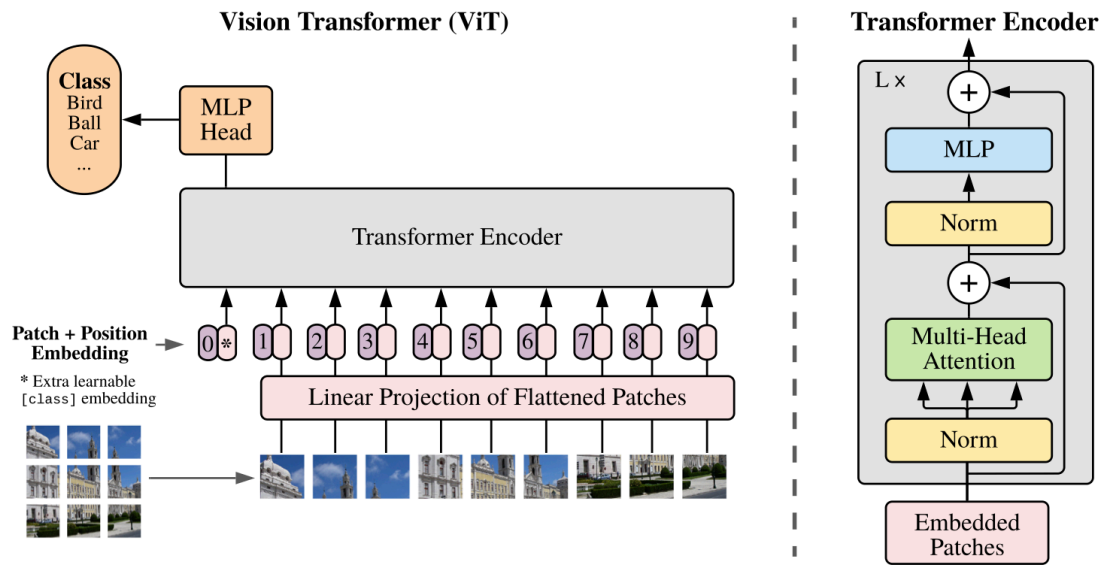


Figure 2.18: The ViT architecture for image classification. [29]

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.

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](https://arxiv.org/abs/2109.00356) [cs, eess]. (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](https://arxiv.org/abs/2012.15029) [eess.IV].
- [8] J. Irvin, P. Rajpurkar, M. Ko, Y. Yu, S. Ciurea-Illcus, C. Chute, H. Marklund, B. Haghighi, 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](https://arxiv.org/abs/1901.07031) [cs, eess]. (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. Malvey, 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, L. 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).