

Automatic Left Ventricle Segmentation in Ultrasound

Ezequiel Ortiz

2018

Contents

1	Abstract	1
2	Acknowledgements	1
3	Intro	2
3.1	Motivation	2
3.2	Metrics	2
3.3	Ejection Fraction	3
3.4	Left Ventricular Mass	5
3.5	Cardiac Imaging	6
3.6	Ultrasound	9
3.7	Echocardiogram	10
3.8	Manual Segmentation	10
3.9	Automatic Segmentation	12
3.10	Deep Learning	14
3.11	CNN	16
3.12	Unet	17
4	Method	18
4.1	Programming Environment	18
4.2	python modules	19
4.3	Dataset	22
4.3.1	dataset preprocessing	23
4.4	training	27
4.4.1	training params	28
4.5	class organisation	31
4.6	data generation	31
4.7	Validation	32
4.7.1	K Fold Validation	32
4.7.2	Generated LVEF and LVM	32
4.7.3	ROC Curve	33
5	Results	33
6	Discussion	34
7	Further Research	35

1 Abstract

This report outlines our investigation into the performance of deep learning methods on semantic segmentation of the left ventricle blood pool and myocardium from 2D echocardiogram images. We found that a convolutional neural network configured in a unet architecture was effective at segmenting the blood pool but further work is required to effectively segment the myocardium.

2 Acknowledgements

I would like to thank Esther and Andy for guiding me through this project. I did every step of this for the first time and learned more than I ever have. Thank you for your expertise and time. Thank you Bram for helping me understand echocardiograms and the potential space for algorithms such as mine in the NHS. I would also like to thank Nick, Irina, James, and Eric for helping me with deep learning python and bash and most of all for making me feel welcome in Beckett House. Thank you Mom and Dad for supporting me though this project and degree. Thank you Lauren for being at my side I can't picture any of this without you.

3 Intro

3.1 Motivation

Cardiovascular disease is responsible for 17.9 million deaths, or 44% of all noncommunicable deaths per annum [1]. Detection and monitoring of heart disease is tantamount to treatment and management of symptoms. Quantification of the interior volume of the left ventricle and thickness of muscle is extremely important due to how illustrative these measurements are to the overall function of the ventricle. [2] The most prevalent method to assess these metrics is cardiac imaging. While there are many imaging modalities, ultrasound has been cemented as a cornerstone of cardiac imaging due to the non-invasive procedure low cost and real-time image acquisition. Interpretation of ultrasound images is time intensive and can be challenging even for trained personnel. The automatic segmentation of the left ventricle blood pool and myocardium would reduce the need for expert evaluation and inter-operator variance. Access to more accurate clinical data with reduced wait times would better patient outcomes and reduce monetary strain on hospitals.

3.2 Metrics

Quantification of the performance of the heart is crucial for diagnosis of cardiac pathology. Due to the complexity of the heart and cardiovascular system as a whole the amount that can go wrong is staggering. We want the metrics that we collect to be illustrative of how the heart is functioning,

reproducible, and easy to understand....

The two most common metrics that are collected during echocardiograms are left ventricular ejection fraction(LVEF) and left ventricular mass(LVM). [3] These quantifications in conjunction with the visual information of the echocardiogram are potent in the assessment of how effectively the heart is pumping blood and the health of the myocardium

3.3 Ejection Fraction

LVEF is actually the combination of several ventricle volume measurements that are made at different points within the cardiac cycle. LVEF quantifies the effectiveness of the heart as it is the ratio of end diastolic volume(EDV) to end systolic volume(ESV), indicating how much blood the heart is pumping per beat.

EF it is calculated according to EQ 2. EQ. 1

$$SV = EDV - ESV$$

EQ. 2

$$EF = \frac{SV}{EDV} * 100$$

Pinning down what is normal LVEF is not trivial as it is very patient specific. Everyone is different and an LVEF measurement that is normal for one patient might be unhealthy for another. LVEF can also change depending on the physical state of the patient when the measurement is made.

The patient's level of cardiovascular stress and whether they are supine or seated will all change the measured LVEF. M. E. PFISTERER et al set out to determine ranges of normal LVEF by assessing the mean LVEF within a population of 1200 patients using ultrasound. The lower limit of normal LVEF values for men has been estimated to be $51.1 \pm 4.2\%$ with the upper limit of normal being $76.6 \pm 3.8\%$. These values were assessed with echocardiogram while the patient was supine and at rest. LVEF was also assessed with the more accurate but invasive x-ray angiography; no significant differences were found. LVEF was for the most part independent of sex or age group, except for patients over 60 who had slightly elevated LVEF levels. Patients that were able to exercise then reached 85% of their maximal HR before their LVEF was measured. On average, Men had an LVEF increase of 10.5% while women only saw an increase of 5.3%. Young patient LVEF increased about 6% more than older patients [4]. These findings show that there is no bright line to determine what is normal or abnormal in LVEF.

The American Society of Echocardiography has published general guidelines for LVEF classification. Reduced LVEF should be classified as equal to or below 40% with borderline being 41%-49% and normal function occupying the 50%-70% range. []

As many variables affect LVEF its diagnostic value comes from its use in conjunction with other metrics. However it is an excellent indicator for change in cardiac function because it is easy and consistent to measure. It is another piece of data about the function of the heart that can be moni-

tored over time. Changes in LVEF can indicate common heart ailments like hypertrophic cardiomyopathy heart failure.

Athletes can often test for LVEF below normal even though their hearts are healthy and functional [5].

And LVEF can be within normal values while the function of the heart is impaired. Heart failure preserved ejection fraction (HFpEF) and heart failure reduced EF are conditions where the heart is failing to supply the body with enough blood LV systolic func LV diastolic func This is because EF only depends EDV and ESV and does not take into account the time duration or the geometry of the contraction. Measuring LVEF is measured with cardiac imaging techniques. [6]

3.4 Left Ventricular Mass

LVM is a measure of how thick the myocardium of the left ventricle is. It is a strong predictor for cardiovascular events [?], as it is an excellent indicator of the health of the myocardium.

From 2D echocardiography images LVM is calculated by first segmenting the myocardium from the echo image.(2-4CH?) From this segmentation, the long as LVM should be measured at the end of diastole, when the left ventricle is relaxed and full of blood. In order to calculate the mass the volume of the myocardium must be measured. To approximate the volume with 2D echo methods the Area Length Method This and other 2D LVV methods assume that the thickness of the wall is relatively uniform so hearts with

regional thickness changes like those afflicted with HCM can lead to inaccurate measurements. We are only attempting to segment the myocardium from 2 and 4CH images as we do not have any short axis images, so we will not be calculating LVM directly. Hypertrophic cardiomyopathy (HCM), as indicated via elevated LVM(how elevated?) occurs when the left ventricle myocardium thickens to a point where heart function is impaired. There are elevated risks of dilated cardiomyopathy in athletes as the natural thickening of the ventricle walls due to exercise can mask the condition.

Dilated cardiomyopathy (DCM) as indicated by lower than normal LVM occurs when the left ventricle myocardium thins and stretches which compromises the strength and function of the left ventricle.

In conjunction, LVEF and LVM can help inform physicians on the efficiency of the heart and the reasons behind that level of function. There are many methods both invasive and non-invasive that can be used to measure LVEF. Non-invasive measurement techniques are always preferred if they offer similar performance to invasive methods. Imaging the heart has the ability of providing detailed structural information of the heart.

3.5 Cardiac Imaging

The most common imaging modalities used image the heart are Magnetic resonance imaging (MRI), x-ray computed tomography (CT), and ultrasound (US).

MRI employs strong magnetic fields and quantum mechanical properties

of particles to get signal from the presence of hydrogen ions i.e. protons. The behaviour of protons within the context of an MRI scan is highly dependent on the structure and magnetic properties of the tissue surrounding each proton. However, the scanner is not getting signal from individual protons but groups of protons within each voxel of the image. MRI is usually non-invasive but can require the use of contrast agents. Contrast agents are usually used in cardiac MRI to detect scar tissue in the myocardium. Scar can form when the blood supply to regions of the heart becomes restricted or halted. As scar is not electricity active like healthy myocardium, its presence will impair hearts ability to pump as well as increase the likelihood of electrical pathologies like ventricular tachycardia and fabulation. As infarcted tissue by definition has limited or absent blood supply, the perfusion of contrast agents through the myocardium will highlight damaged regions. MRI offers unparalleled contrast between different biological tissues. Once the heart has been scanned, it is relatively easy to calculate the interior volume of the left ventricle. For example, region growing methods are usually used to directly calculate the volume from the 3D scan. This approach is very accurate as no assumptions are made about the shape of the ventricle.

However, scanners are usually in short supply, require skilled operators, are expensive to run and maintain, and are slow to acquire images. When the subject of the image is moving, long image acquisition time can lead to motion artefacts. Patients must hold their breath while their heart is being imaged to avoid introducing such artefacts. Patients that are unable

to hold their breath, or those whose heart moves while it pumps may be ill-suited for MRI. This long image acquisition time also prevents the cardiac cycle from being captured in real time. The CMR images are captured at equivalent points within the cardiac cycle using the ECG to time each capture point. As CMR image acquisition gets faster this will be less of a problem, but for now the CMR images are an average heart over many cycles. Metrics gathered from CMR images are generally easier to calculate and more accurate but eliminate beat to beat variation. Additionally, due to the extremely strong magnetic fields present within the scanner, those with electrical or ferrous metal implants will not be able to be scanned. Some medical implants are labelled as MRI conditional, meaning that they are safe to be scanned, with some limitations to the scan. However, due to the complex restrictions placed on scanning these devices, the increased expense of MRI conditional technology and the relative rarity of such devices, most clinicians elect to use an alternative imaging modalities. [7]

CT imaging involves taking many 2D x-ray images from different angles then using computed tomography algorithms to generate 3D images. CT imaging is not generally used for cardiac imaging as soft tissue does not provide sufficient contrast. The high energy x-rays used in CT are not heavily attenuated by low z elements found in soft tissue. Functional imaging of the heart with CT must utilize contrast agents [5]. High z compounds are injected into the bloodstream and subsequently imaged to give images of the blood pool within the LV CT imaging has excellent spacial and temporal resolution.

Fast image acquisition means minimal motion artefacts when compared to MRI. However, the radiation dose and use of contrast agents make cardiac CT rare.

3.6 Ultrasound

Ultrasound as used for echocardiograms is the core of cardiac imaging. High frequency sound waves are sent into the body, reflected refracted and scattered, and picked up by the same transducer that produced the original sound pulse. These sound waves are reflected most strongly across boundaries of differing acoustic impedance. In the images in our dataset this is illustrated with a distinct boundary between the myocardium and the blood pool. However, imaging the heart with ultrasound does present some problems. The ribcage encases the heart and a transducer that can fit in-between the ribs must be used as bone is effectively opaque to the sound pulses. When the heart is imaged across the ribcage the scan is called a trans thoracic echocardiogram(TTE). Even with correct positioning there can still be some shadowing from the ribs present in some images. Imaging around the ribs can be circumvented with transesophageal echocardiograms(TEE). This type of scan involves placing the ultrasound probe down the patients oesophagus and imaging the heart from within the ribcage. While TEE does provide higher quality images, it is usually reserved for acquiring detailed information about the atria or valves of the heart. TTEs are much more common than TEEs due to the invasive TEE procedure and marginal increase in detail for routine

echocardiograms.

It is the most widely used and prevalent modality to assess the structure and function of the heart. Ultrasound is robust as different scanner modes allow for 2D and 3D image and video acquisition. Doppler imaging allows for blood flow analysis and the use of contrast agents can provide increased levels of contrast. The relatively small size of ultrasound scanners when compared to MR and CT make ultrasound machines much more flexible. Unlike the other modalities where the patient is physically placed within the scanner, the ultrasound probe is positioned around the patient. In most echocardiograms, both a 2CH view and 4CH view is taken. As these image views are approximately orthogonal, information from both can be combined for a better understanding of the hearts 3 dimensional shape. 2D 3D 2CH 4CH

3.7 Echocardiogram

Echocardiograms are integral to cardiac evaluation. They are often the first type of cardiac scan that patients undergo. The procedure involves ECG Qualitative approach How the patient feels

3.8 Manual Segmentation

Once the echocardiogram has been completed, the images must be analysed by a trained professional in order to calculate meaningful metrics. Both

LVM and LVEF require segmentation, where some regions of the image are identified and marked as part of the myocardium or blood pool. In the case of LVEF, in order to minimise the assumptions being made about the shape of the ventricle, both the 2CH and 4CH views are segmented.

The American Society of Echocardiography and European Association of Cardiovascular imaging have published a set of guidelines for LV chamber quantification that are as follows. As most ultrasound machines that are used for echocardiograms come with software for segmenting blood pools, these steps should be considered within that context. First, frames of diastole and systole must be identified. This will be facilitated by the use of the single lead ECG that is present on the images. The QRS complex indicates the beginning of ventricular depolarisation or systole, while the T wave indicates the beginning of repolarization or diastole.

Then a line will be drawn from the mitral valve annulus, a fibrous ring at the base of the valve that shows up bright in most echoes, along the endocardial border and to the other side of the annulus. The beginning and end of the line must be connected with a straight line that ignores the position of the valve. Trabeculations are regions on the epicardium where blood becomes trapped and calcified within the myocardium. [?] They can look like myocardium within echo images and should not be counted as such when determining LV volume.

From the straight line across the valve, the longest orthogonal line is then drawn that connects the valve line to the epicardium at the apex of the heart.

Now the bi-plane method of disk summation, or Simpson’s Bi-plane, can be used to calculate the LV volume.

$$Volumeofdisk = pi(d1/2)*(d2/2)*heightofdisk \quad Volumeofheart = summationofthedisks \quad (1)$$

This method is widely used but suffers from endocardial dropout, user variation, and incompatibility with some deformed ventricles. [?]

Our procedure mirrored the steps laid out in the guidelines, but differed slightly. As we quickly found out, most of the hearts within our data were pathological so the ECG shape was not very useful. In these cases we simply tried to find frames that looked like they contained the largest and smallest LV volumes. We also struggled with endocardial dropout as highlighted in the guidelines.

3.9 Automatic Segmentation

While most every medical professional can take heart rate and blood pressure by hand, these tasks are now relegated to machines. [?] In this case we find that it is best practice to take manual blood pressure and heart rate by hand because of the patient contact [?]. Humans will pick up on variation within the heart rate that could be more insightful than the digital readout of automatic systems. We also don’t patients to feel as if their care has been relegated to machines. As manual segmentation is prone to interoperator variance [?] and time spent manually segmenting images is time spent

away from patients, an automatic procedure would be beneficial. While 2D echocardiograms are useful, the dialogue with patients to determine how they view their physical abilities is just as valuable to determining whether more scans should be taken. Qualitative judgement of echocardiograms would be expected if the cardiologist is given ample information to make a decision. For example, if a patient has an LVEF that looks to be normal and they report no issues associated with reduced EF or heart function, the cardiologist would not need to manually calculate LVEF. LVEF is well suited to automation because the metric is time consuming and variable compared to automatic methods and the metric is useful but almost never the sole reason to make a diagnosis or prognosis description, it is an indicator.

While cardiologists are experts at diagnosing what is wrong with a patient's heart, they synthesise an array of information that is comprised of much more than just LVEF and LVM. Many are able to make their diagnoses simply by looking at the echoes and making a qualitative judgement. When qualitative measurements are made variance older methods Stat shape models Kalmann filters are

Classically, neural networks had been used to classify images. When an image is fed into the network, the network assigns a label to classify the contents of the image. This task requires a transformation of the spacial information of the image into feature information, or information about what the image contains. Automatic semantic segmentation is the logical next step, where we not only want to know what is in an image, but what pixels

contain our object.

3.10 Deep Learning

in theory less assumptions than stat shape models can be trained to better robustness as computational power increases can be used more and more and gets more powerful easy to continue improving networks

Neural networks are computational structures that mimic the function of biological neurons. Networks with sufficient complexity can model any mathematical function. The functional unit of neural networks is the neuron. Individual neurons, like the whole, will each receive an input and provide an output based on an internal activation function. Activation functions can vary in sophistication from simple step functions that provide a binary output based on a hard coded threshold, to sigmoid functions that organically map the input to a range of values between zero and one. The selection of activation functions will depend on the desired behaviour of the network.

One of the simplest configurations of neurons into a network is the perception. We will examine a hand written digit classifier multi-layer perceptron(mlp) to determine how neural networks make decisions. The structure of most networks consist of layers. In mlps, there is an input layer, hidden layers, and the output layer. The input layer in the case of the mnist digit image will simply be a vectorized image. The hidden layers contain the neurons that make the classification decision. The output layer will have a neuron for every class within the dataset so in our digit case we will have

ten output neurons(0-9). The output layer will also be a softmax layer. The sum of all neuron values in a softmax layer must sum to one. For a given input each pixel value of our image will get sent to every neuron in the first of the hidden layers. These layers map the image input to a decision of what digit the network thinks it is seeing. If every neuron within the first hidden layer is getting input from every input layer neuron, then the output of every hidden layer neuron will be the same. In order to prevent this, we need to introduce the ideas of weights. Weights signify the strength of connections between neurons. If a neuron has been strongly associated with a three, in our digit classification example, then there will be neurons that fire strongly when the image provided is a three. The weights in any neural network determine the flow of information through the network and the final decision. In this classification example, each digit from zero to nine has different characteristics about it that make it easily recognisable to our brain as well as a trained neural network. A fully trained network is structurally identical to its untrained counterpart. It is the weights that dictate the strength of the connections and by extension the pathways responsible for the pseudo cognition.

The difference between the mlp and a deep network is simply the number of layers. The mlp by convention is limited to one hidden layer and any network with more than one hidden layer is considered deep. Deep networks have seen much more use than mlps recently as computer hardware advancements have allowed the training and implementation of networks with many

more than two hidden layers. While the mlp simply took the pixel intensity levels as input and subsequently made decisions based purely on pixel intensity, more and more deep networks are utilising convolutions to extract relevant features from the input images. Convolution involves generating a new image from the input image. Each pixel of this new image is the result of an element-wise multiplication of a region in the original image and a filter. For example, a 3 by 3 averaging mask averages the intensities of a pixel and its eight neighbours to give the intensity of the output pixel. This convolution is carried out repeatedly as the filter is swept across the input image.

3.11 CNN

Padding is as simple as adding pixels to the edges of the input image to increase the resolution. Convolution naturally reduces the resolution of the output due to the thickness of the mask. There are many types of mask that can be used in convolutions and they produce a wide variety of output images that highlight different features of the original image. Sobel and Laplacian filters make edges bright and flat surfaces dark. Averaging filters can mitigate noisy images and produce blur. The size of mask can also be adjusted. We used an example of a 3x3 mask but larger masks such as 5x5 are also common. Masks are usually square with an odd width and height to ensure a middle pixel. All of the convolutions we have talked about are predetermined. A known mask is applied in order to produce

a specific output. In deep convolutional neural networks, the masks are analogous to the mlp weights in that they will be defined through training. The nature of optimizing the network to complete segmentations means that the filters within the unet will be optimized to better select and highlight the features within the image that are relevant to the segmentation task at hand. For example, the network trained to segment the ultrasound cone from out images will most likely have filters that tend to identify the boundary between the background and ultrasound data, while the network trained to segment ventricles should be looking to identify the relatively dark blood pool and bright myocardium.

3.12 Unet

The Unet architecture is the standard approach for segmenting both 2D and 3D biomedical imaging data. It is strongly reliant on image augmentation to increase the size and variance of comparatively small biomedical image datasets. The signature architecture is comprised of a downsampling pathway followed by upsampling to produce a mask that is the same dimension as the input image. As the input image is convolved and downsampled, spacial information is converted into feature information. What constitutes as feature information would not make much sense anywhere but within the workings of the neural network. The information that is distilled as the image progresses through the downsampling phase is directly related to the overall purpose of the network. The feature information in our case will be vital to

the network to generate predicted masks for a given image.

The subsequent expansion after the initial distillation of information gives the unet its "U" shape.

During the expansion phase, images from the contracting end are combined with the upsampled image. The network is "reminded" what the original image looks like with these connections.

The downsampling of the input image also allows us to use one filter size. Without downsampling we may need to use many filters of different sizes to capture differently sized features.

The structure of the unet from the original paper can be see in figure x

4 Method

4.1 Programming Environment

During the first stages of the project, we worked within the Google Colab environment. Google colab is an online coding suite where python code can be run from a web browser. We chose colab because of the free cloud GPU that would be needed to train the network. Colab mimics jupyter notebooks functionality in that code can be split up and run in sections even within the same script. We were also unconstrained by a physical computer and code could be written and ran from any web browser given access to the project. However there were some downsides to colab that prompted the eventual switch. All of the data for training and validation had to be stored

in the google drive of the account that owns the project. Any changes that were made to the dataset had to be uploaded to drive, and we had to make sure that the dataset was the most recent iteration. The performance of the cloud GPU was not as fast as a local GPU. Some longer training regimes were halted on account of google colab becoming idle and ending our processes. This was more of a issue before we implemented model checkpoints that allowed us to save model progress as it was training. Google colab was an excellent starting point for the project as it removed lots of start up cost for a deep learning project such as this, but the downsides outweighed the benefits in the end.

The move to running our code and storing our data on a local machine came once we started getting promising results with the blood pool segmentation. As the blood pool was considerably more complicated to segment than the ultrasound cone, the need for more intense training prompted a move. We moved our project to a pc containing an Nvidia GTX 1060 GPU and running Ubuntu 18.04. After installing the relevant cuda drivers for the GPU and anaconda to manage our python libraries our training regimes were faster and could run for longer.

4.2 python modules

We chose Ubuntu as an operating system because it is free and open software. As the project is as much about working with data as it is programming, the Bash terminal has been an indispensable tool. While the learning curve

has been steep, we find ourselves settled into Linux using vim to write all of the code and report. Vim is extremely and extremely powerful and extensible terminal text editor. Tmux, a terminal multiplexer, has extended the capability of the terminal by allowing us to perform many tasks within one window. The power of bash scripts has allowed to perform monotonous manual tasks such as renaming every file in a directory with one line of code.

Python is the language of choice due to the many well maintained modules geared for deep learning.

We use anaconda to maintain our python programming environment. Anaconda gives us complete control over package versions to avoid dependency problems. For example, we have to use python version 3.6.8 instead of the more recent python 3.7.3 as tensorflow cannot yet run on python 3.7 at the time of writing. With anaconda we can have a tensorflow environment where we can use this older version of python and still have the ability to use the most recent releases by changing environments.

The core of project from a programming perspective is keras. Keras is a high level python module that wraps tensorflow. Tensorflow is the most ubiquitous deep learning library and converts python code into GPU optimised c++ code. The parallel processing power of GPUs is extremely well suited for the intensive image operations and extensive micropropagation steps that training CNN's involves. Keras is mostly used for quick prototyping projects. Its ease of use lends itself to quick implementations and not production code. Pure tensorflow is much more powerful, but significantly more complex than

keras, but we never found our work impeded by keras.

Numpy is an extremely popular python matrix library. It has an extensive set of methods that allow for a diverse set of optimised matrix operations. As our images were stored in numpy arrays, most every operation in the dataset preprocessing section relies on numpy functions. Keras is designed to work with datasets that are stored in numpy arrays due the optimizations built into the module.

As the images were stored in nifti files, we needed a way to properly read the files into memory. After looking at both simpleitk and nibabel, we settled on nibabel. We use nibabel to read the nifti files into memory from storage. We can then extract the image information, image header, and affine from the nifti file. The image information is then stored in numpy arrays. The image header contains the dimensions of the pixels in the image. If our end goal is to generate approximations of volume, we need to correspond pixels to physical distances. The affine of a nifti image contains information that relates the coordinates of the nifti image to world coordinates in RAS+ space. This has to do with radiological conventions that govern how medical data is interacted with and viewed in computers. We use affines to generate nifti files from our predicted masks. As each mask corresponds to a specific image all the information that we need to turn the numpy array that is returned from our network predictions into a nifti file to be stored on disk is stored in the image header and affine.

4.3 Dataset

The raw dataset that was provided to us consisted of pathological echocardiography images from 96 patients. Most patients had both a 2CH and 4CH view of their heart. The files are stored as zipped nifti images and are 2D greyscale videos of the heart through multiple diastole systole cycles. Constructing the final network that will segment cardiac ultrasound will be broken up into several successive steps.

1. Train a network to successful segment ultrasound cone from the image
2. Train a network on 2CH blood pool data
3. Train a network on 4CH blood pool data
4. Train a network on both 2CH and 4CH blood pool data
5. Depending on the success of combining 2 and 4CH data, train a network to segment the myocardium

The purpose of segmenting the ultrasound cone from the rest of the image to make successive training easier. Eliminating image data that is irrelevant to the segmentation of the heart saves the network from having to learn that the extraneous information is extraneous. The build up to the final network allows me to gradually build up the training dataset.

The images need to be converted into numpy arrays before they are used to train the network. The python library nilearn is used to convert the zipped

nifti files into numpy arrays. The images are then resized to a resolution of 512 by 512 and the intensities are remapped from 0-255 to 0-1. The images are then collectively stored in a 4th order tensor with dimensions (number of images, X dimension, Y dimension, color channel). This allows for easy manipulation with Keras

The masks go through the same resizing conversion from nifty to numpy array. As the mask images can contain cone or ventricle masks, one layer of the masks must be chosen for training. This is accomplished via a thresholding operation to eliminate either the cone or ventricle mask. Once the correct mask configuration has been chosen, the masks can be resized and stored in a fourth order tensor as well.

4.3.1 dataset preprocessing

The data was presented to us in the form of 96 folders, each corresponding to a certain patient. The names of these folders were all unique as they came from different hospitals at different times and were acquired by different machines and operators. Most patients had both two and four chamber images, with some having multiple of either. the naming convention for the individual images was as follows:

US_2CH.nii.gz

US_4CH.nii.gz

and if there were multiple of one view a two is appended to the end of the

name as follows, *US_2CH2.nii.gz*

One of the first things that I did was to separate the dataset into two and four chamber images. However, the image files would need to be renamed so that they could be told apart. The name of the directory containing the images was appended to the image filenames so that there would be no conflict of names and the image view should still be identified by reading the name. An example of a filename in this form is as follows:

KCL_GC_001_US_2CH2.nii.gz

Once the images were somewhat sorted, manual segmentation could begin. Behind every competent neural network is a meticulously crafted dataset. To create a data point the blood pool and ultrasound cone from an individual frame of the nifti file must be segmented. To do this ITK-snap was used to both open the nifti file and produce the segmentations. In order to increase variance within the dataset and therefore robustness of the trained network, the chosen frames are at the general point of diastole and systole. Most all of the scans taken were of pathological hearts, so finding exact diastole and systole was challenging. Most easily distinguished features of the ECG within the echocardiogram were undistinguishable to our untrained eyes. We usually settled on two or three frames from each scan that looked sufficiently different. If more data is needed at a later time, additional frames could

be segmented. The product of the segmentation is a mask that contains information on both the boundaries of the ultrasound cone and blood pool of the left ventricle. The background of the mask is always zero. If a cone mask is present, the cone is always one with the ventricle mask valued at two. If the ventricle is the only mask present then it will have a value of one. TODO(image of mask with cone and vent and mask with only cone) As the ultrasound cone is quite simple to identify within the images, we only ended up producing 25 data points with the cone manually segmented.

Once a mask was completed it needed to be associated with the correct frame from the correct image. The naming convention that we settled on was to simply use the same name as the image with the correct frame number appended to the end of the name.

KCL_GC_001_US_2CH2_01.nii.gz

The above example refers to the mask associated with the first frame of image *KCL_GC_001_US_2CH2.nii.gz*. This limits us to two digits to determine the frame so we are limited to frames 0-99. Most images are well under 100 frames. Images that are longer usually have a full cardiac cycle within the first 100 frames.

In order to train a network on the images, they needed to be a standard resolution. Initially this resolution was to be 512 by 512. 512 is a power of two in order to avoid non integer resolutions as the images are downsampled

in the unet. The raw images in the dataset varied in resolution(TODO exact resets pls) with 512 by 512 being a rough mean. We later moved to a resolution of 800 by 800. This change was made based on two factors. We wanted to zero pad the images up to a resolution instead of cropping them or interpolating them to a smaller size. Interpolation would give a reproduction of our images at a different resolution while zero padding would encase the exact original image in zeros. Our zero padding scheme attempts to keep the original image in the center of the padding. While interpolating image data is common practice, the end goal of producing a measurement of LVEF in real world units places a higher degree of importance in preserving the pixel spacing and dimension information contained in the image headers. As the data generation script took shape, moving data from nifti to numpy array and back again without worrying about pixel spacing was a relief. Once the images and masks were a standard resolution the image intensities needed to be scaled from zero to one. This was achieved with scalar division by 255.

Now that the data has been standardised, we can construct the datasets for training. As some masks contain both ventricle and cone information, they must be processed to produce a separate ventricle and cone mask. At the initial round of training the cone dataset consisted of 25 image and mask pairs and the 2CH dataset consisted of 43 image and mask pairs.

4.4 training

Due to the small amount of data that we were working with along with memory restrictions, on the fly data augmentation was going to be crucial to the performance of the project. The premise was to apply random image transformations to an image and mask pair to generate a "new" piece of data from a base image. We decided to use the keras data generator class, which is capable of performing a multitude of rigid image transformations. While some non-rigid transformations could be useful, and are recommended in the unet paper, we needed them to be precise enough to only deform the ventricle as we wanted the cone to remain unreformed. We needed to tread the fine line of augmenting our data in a way that increases the variation within the dataset without going too far. Too much augmentation that is too drastic could lead to the network struggling to find any sort of pattern without an immense amount of training. Due to machine and time constraints training needed to be as efficient as possible. After some trial and error we settled on the following parameters to define our image and mask data generator:

1. rotation range of 45 degrees
2. width shift range of 10% of total image width
3. height shift range of 10% of total image height
4. shear range of 5 degrees
5. zoom range of 0.1

6. fill mode of nearest so most likely zero

We apply this same augmentation scheme to both the training and validation data to increase the size and variance of both sets. It is customary in machine learning, to take some data that could be used for training, and use it instead to test the trained model. Here we would give the trained model an image that it had never seen before and have it predict a mask. Comparison between the predicted and true masks gives us a better understanding at the effectiveness of the model. As the model is literally optimized to perform on the training data, using it to validate the network performance would not illustrate any ability to perform on real life data.

4.4.1 training params

We settled into training parameters early in the project. We needed to ensure that we would train long enough to get our validation loss sufficiently low, while minimizing training time. For the final cross validation training, we chose 1000 epochs. We knew that this was more than enough, so after one fold was trained, we noted the point of diminishing returns TODO and chose x epochs for the final run.

batch_size Batches are groups of images to be trained on each step. When we attempted to have a batch size greater than one our machine ran out of memory and the training script crashed. This could be due to machine restrictions, but is more likely the large image resolution and size of unet.

`steps_per_epoch` The number of steps per epoch determines how batches are trained on per epoch. Ideally every image is trained on for every epoch. As our batch size was 1, we needed our steps per epoch to be equal to the number of training images, which in our case was 80.

`optimizer` The optimizer is the most important aspect of the network learning. This algorithm performs stochastic gradient descent with an adaptive learning rate. Instead of performing computationally expensive gradient descent on every single batch, the stochastic aspect randomly selects images to optimize for. This ensures that the network is being optimized on a set of images that reflect the variance of the entire dataset. The other major aspect to the adam optimizer is the dynamic learning rate. If the stochastic gradient descent step tells us in what direction we want to tweak our weights, the learning rate defines how much we change our weights. If our learning rate is too high, our steps will be too large and our network will not be able to settle into a configuration that minimizes our loss. If our learning rate is too low, it will take too long for our network to reach a loss minimum. Adam handles this problem by dynamically changing the learning rate based on the parameters of the problem. If a certain feature is very common, for example the apex of the heart in most scans is usually a source of good signal, then adam will reduce the learning rate for this feature. While for sparser features would require a larger learning rate to ensure the network learns about them even though they are rare. An example of a rare feature could be the presence of hypertrophic myocardium. This means that there is no longer

one global learning rate but a learning rate that is associated with different features within the training data. The other defining feature of the adam algorithm is the root mean square propagation(RMSprop). RMSprop also adapts the per feature learning rates based on the average of recent learning rates. This gives adam good robustness on noisy sparse datasets [?].

loss To track the loss of the network we decided to use binary cross entropy(BCE) and dice loss. Binary cross entropy is. Dice loss measures the overlap between the predicted label and the ground truth. For the two classes in our segmentation task, we use the 2 class dice loss as calculated by EQ The dice coefficient ranges from 0-1, with 0 being no overlap and 1 being a perfect match. To convert the dice coefficient into dice loss we simply subtract it from 1 to invert the metric so that we can minimise it

To combine the two metrics we simply add them together and minimize the sum as our loss function.

saving In order to monitor the progress of our network and save training progress we utilized a number of keras callbacks. Callbacks are simply functions that can be called during training. We used the CSVLogger, ModelCheckpoint and LearningRateScheduler. The CSVLogger is run after every epoch and saves developer defined metrics the epoch number, training and validation BCE dice loss, and training and validation dice loss. ModelCheckpoint saves the models

4.5 class organisation

Once we completed our first successful cone segmentations within google colab, it was clear that we needed to move to a more object oriented approach. We decided to store relevant functions for data preprocessing within a module. We wanted to instantiate a unet object from a predefined parameters object that could then train, make predictions, and load previously trained models.

4.6 data generation

After constructing the cone and 2CH datasets, we saw a need to streamline the process. We were spending a considerable amount of time entering file-names and clicking on images to open in itk-snap for the segmentation. We then decided to write a script that would iterate over every nifti file in a given directory and:

1. open the nifti in itk-snap for the user to view the image
2. prom the user to input which frames would be segmented
3. give the image to a partially trained network
4. post process the generated mask and convert it into nifti file
5. open the image with the mask in itk-snap for user refinement

As the names of the masks were automatically generated based on the user inputted frame and original image name the user no longer had to type in

the name of the masks. We also defined the directory where the masks and images were to be saved within the script so the user did not have to specify that within itk-snap either. This script proved challenging due to the movement of image data between np and nifti but the overall time savings outweighed the development time.

4.7 Validation

A network is only as good as it can be proved to be.

4.7.1 K Fold Validation

To assess the performance of our model we will both assess the performance with a K fold cross validation scheme and compare the accuracy of the predicted LVEF and LVM. K fold cross validation involves a number of steps.

1. Split the dataset into K equal groups
2. Select one group for the test set and have the rest be training data
3. train K models so that every data point gets to be in the test set

In our case we chose K to be 5. This fit well with our final blood pool dataset size of 100 image mask pairs as it is cleanly divisible by 5.

4.7.2 Generated LVEF and LVM

Predicting medically relevant metrics from 2D echocardiogram images is an additional method of validation. While calculating LVEF and LVM to high

precision is left to 3D imaging techniques, that does not make 2D quantitative analysis irrelevant. The ease of acquiring and working with 2D echoes make them perfect for early detection of abnormalities. We w

4.7.3 ROC Curve

As the network returns a probability map of the ventricle and not a binary mask, we must threshold the masks before we compare them to ground truth. In order to assess the best threshold value we will plot the networks true positive rate over its false positive rate across a range of thresholds.

5 Results

To assess the performance of our network under the 5 fold cross validation scheme, we tested every image in each fold's test set to get an average dice. Our 5 fold cross validation scheme generated the following metrics

The training history for one of our folds is as follows.

To test the efficivness of our model we constructed an ROC curve for each cross validation model..

The reason LVEF only requires the segmentation of two frames is because of the cost of making those segmentations. As we have a network that can segment every frame, we can plot the volume of the LV through time. This is a result that would have been arduous to perform without our automatic methods.

6 Discussion

Overall we are quite pleased with our network’s performance on segmenting the left ventricle. The cross validated ROC curve shows excellent auc metrics on account of the high true positive rate and low false positive rate. The auc, which is related to the accuracy of the test is very close to 1 for every fold. This shows that the model doesn’t under or overestimate the area of the blood pool.

The performance on segmenting the myocardium was not on par with the blood pool results.

While our method blood pool segmentation was successful, the myocardium segmentation was not. We constructed fewer data points for the myocardium because there were many echo images that did not contain the full epicardium. The thickness of the myocardium was also much more difficult to determine. We extrapolated heavily from regions with more defined thickness. Short axis images would allow us to not only be more sure about the thickness but also actually calculate the LVM of the heart. But, if short axis 2CH and 4CH views are all being collected, we recommend moving from 2D echos to 3D.

To construct our myocardium dataset, we manually segmented the myocardium, however most methods to calculate the LVM manually subtract the blood pool volume from the volume of the epicardium. Considering the effectiveness of our model on the bloodpool, perhaps the simpler segmentation

of the epicardium would have been a better route to go down.

The history of training showed that our model converged about halfway through our training regime. We avoided overfitting by only saving the model that improved the validation loss.

As we have

7 Further Research

If another team or person were to continue exploring deep learning on 2D echocardiogram data, this is where we recommend they spend their efforts.

While we stated the need for short axis images, in order to determine the LVM, we recomend 3D echo methods. 3D images remove geometric assumptions and would lead to more robust algorithms, even when irregular hearts are imaged.

When we were manually segmenting the ventricle and myocardium from the echo images, in many cases it was exceedingly difficult to discern where the walls of the heart were. We also consulted a practising cardiologist, who corroborated this saying that in some cases the segmentation can be difficult. But what humans can do that this network cannot, is watch the video to contextualize a certain frame. In order to allow a network accomplish this, an LSTM may be employed. LTSM's employ a recurrent nature where the network would watch the entire echocardiogram video before attempting to segment the individual frames. The idea being that the network will learn

general information about the behaviour of the heart through time and use that information to segment individual frames.

Before neural networks are put into practice in hospitals, doctors and legislators need to be sure that networks can be trusted. One possible solution to this would be to give networks the ability to explain why they make the decisions that they do and how sure they are about their answer. A practical extension of this idea for our network could be the deletion of regions within the image that inform its decisions. There is research being done... We would also need to compare the performance of these networks to humans. We make several claims that this network will reduce interoperator variance, this makes sense as the nature of learning the dataset averages the biases that went into making it. Perhaps we need to involve a larger group of people in the manual segmentation to increase the variance within the segmentations and allow the network to learn how more people segment echoes.

longitudinal strain computer resources spec around batch size architecture We only investigated a unet architecture. While the unet is standard for biomedical segmentation, perhaps some other architectures like densenet [?] could be explored. For our data augmentation, we only performed rigid transformations. Our thinking behind this was the network would probably use the edges of the ultrasound cone to find the ventricle and elastic deformations could disrupt this. Elastic deformations that are constrained to the interior of the cone should be explored.

8 Conclusion

References

- [1] W. H. Org, “World health statistics 2018,” WHO, Tech. Rep., 2018. [Online]. Available: http://origin.who.int/gho/publications/world_health_statistics/2018/en/
- [2] F. M. B. M. R. F. R. B. D. M. F. A. F. M. E. F. M. P. A. P. M. M. H. P. M. M. J. R. M. J. S. M. J. S. S. M. F. S. D. S. M. K. T. S. M. F. M. S. J. S. M. F. Roberto M. Lang, MD and M. William J. Stewart, “Recommendations for chamber quantification: A report from the american society of echocardiographys guidelines and standards committee and the chamber quantification writing group, developed in conjunction with the european association of echocardiography, a branch of the european society of cardiology,” *Journal of the American Society of Echocardiography*, vol. 18, pp. 1440–1463, 2005.
- [3] F. F. L. P. B. M. P. F. V. M.-A. P. F. J. A. M. M. A. A. M. M. L. E. M. P. F. A. F. M. F. E. F. M. F. S. A. G. M. T. K. M. P. P. L. M. P. F. D. M. M. P. M. H. P. M. F. E. R. R. M. P. L. R. M. F. K. T. S. M. F. W. T. M. Roberto M. Lang, MD and P. F. Jens-Uwe Voigt, MD, “Recommendations for cardiac chamber quantification by echocardiography in adults: An update from the american society of echocardiography

- and the european association of cardiovascular imaging,” *Journal of the American Society of Echocardiology*, vol. 29, pp. 277–314, 2016.
- [4] M. E. PFISTERER, A. BATTLER, and B. L. ZARET, “Range of normal values for left and right ventricular ejection fraction at rest and during exercise assessed by radionuclide angiocardiology,” *European Heart Journal*, vol. 6, no. 8, pp. 647–655, 08 1985. [Online]. Available: <https://dx.doi.org/10.1093/oxfordjournals.eurheartj.a061916>
- [5] T. H. Marwick, “Ejection fraction pros and cons: Jacc state-of-the-art review,” *Journal of the American College of Cardiology*, vol. 72, no. 19, pp. 2360 – 2379, 2018. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0735109718383542>
- [6] N. H. M. Arrow, N. K. McAlister, and K. Buttoo, “Understanding cardiac ”echo” reports,” *Canadian Family Physician*, vol. 52, pp. 869–874, 2006. [Online]. Available: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1781094/>
- [7] M. B. F. Duc H Do, MD; Noel G. Boyle, “Mri in patients with implanted devices: Current controversies,” *JACC*, 2016.