

PREDICTING BRAIN CONNECTIVITY MAPPING USING RADIOMICS FEATURES IN ANATOMICAL MRI

LEVENTE ZSOLT NAGY

Thesis supervisor

ALFREDO VELLIDO ALCACENA (Department of Computer Science)

Thesis co-supervisor

ESTELA CAMARA MANCHA (Hospital Universitari de Bellvitge)

Degree

Master's Degree in Artificial Intelligence

Master's thesis

School of Engineering
Universitat Rovira i Virgili (URV)

Faculty of Mathematics
Universitat de Barcelona (UB)

Barcelona School of Informatics (FIB)
Universitat Politècnica de Catalunya (UPC) - BarcelonaTech

Abstract

This study explores alternatives to diffusion MRI for mapping brain connectivity, predicting fractional anisotropy and mean diffusivity, using radiomic features derived from T1 and T2 structural MRI images. This approach aims to significantly enhance the cost and time efficiency of data acquisition, eliminating the need for diffusion MRI and tractography. The research is centered on the basal ganglia, a region primarily affected by neurodegeneration in Huntington's disease, comparing its characteristics between control subjects and patients with the condition.

Contents

1	Introduction	8
1.1	Objectives	9
1.2	Motivation	10
1.3	State of the Art	10
2	Design	11
2.1	Preprocessing	11
2.1.1	Raw Data	11
2.1.2	Quality Control	14
2.1.3	Radiomics Features	15
2.1.4	Coordinates	16
2.1.5	Data Augmentation	16
2.1.6	Scaling and Normalization	16
2.1.7	Data Balancing	18
2.1.8	Clinical Data	19
2.1.9	Relative Connectivity	20
2.2	Evaluation	20
2.2.1	Train, Validation and Test Splits	20
2.2.2	Accuracy and Pearson Correlation	21
3	Experiments	23
3.1	Subcortical Segmentation	23
3.2	Methodology	24
3.2.1	Missing Records	25
3.2.2	Architecture Tuning	26
3.3	Diffusion Fractional Anisotropy Regression	26
3.4	Mean Diffusivity Regression	29
3.5	Relative Connectivity Segmentation	32
3.5.1	Exhaustive Sequential Backwards Feature Selection	36
3.5.2	Streamline Regression	38
4	Sustainability	39
4.1	Environmental Aspect	39
4.1.1	Developement	39
4.1.2	Production	40
4.2	Economic Aspect	40

4.2.1	Cost	40
4.2.2	Return	41
4.3	Social Aspect	41
4.3.1	Development and Collaborations	41
4.3.2	Inclusivity	42
4.4	Risks	42
5	Conclusions	43
5.1	Future Improvements	43
5.2	Project Future	44
Sources of Information		45
A	Software Design	46
A.1	Raw Data	46
A.2	Common Functions	48
A.3	Preprocessed Data	49
A.4	Data Generator	50
B	Software Implementation	54
B.1	Interpolation	54
B.2	Multithreading	54
B.3	Warp Pre-Computing	55
B.4	Data Normalization	55
C	Additional Figures	56
D	Additional Tables	60
E	Source Code	77

List of Notations & Abbreviations

MRI	magnetic resonance imaging	8
DTI	diffusion tensor imaging	8
FA	fractional anisotropy	8
MD	mean diffusivity	8
RD	radial diffusivity	8
ROI	region of interest	8
NN	neural network	15
FNN	feedforward neural network	23
CNN	convolutional neural network	15
FCNN	fully convolutional neural network	14
NIfTI	neuroimaging informatics technology initiative	11
FMRIB	functional magnetic resonance imaging of the brain	
FSL	FMRIB software library	47
FNIRT	FMRIB's nonlinear image registration tool	8
GLCM	gray level co-occurrence matrix	16
GLSZM	gray level size zone matrix	16
GLRLM	gray level run length matrix	16
NGTDM	neighbouring gray tone difference matrix	16
GLDM	gray level dependence matrix	16
cUHDRS	composite Unified Huntington's Disease Rating Scale	19
CAP	CAG Age Product	19
UML	unified modeling language	47

List of Figures

1.1	Basal Ganglia (ROI) & Cortical Targets	8
1.2	Connectivity Maps	9
2.1	Simple Model Overview	11
2.2	Basal Ganglia Subcortical Segmentation	13
2.3	Histogram: Firstorder Energy	17
2.4	Histogram: GLDM Small Dependence High Gray Level Emphasis	17
2.5	Histogram: NGTDM Busyness	18
2.6	Balance: Subcortical	19
2.7	Balance: Diffusion MD	19
2.8	Balance: Diffusion FA	19
2.9	Balance: Relative Connectivity (thresholded at 0.6 & binarized)	19
2.10	Distribution of Records in relation to Datapoints	21
3.1	Training Curve: Diffusion Fractional Anisotropy	27
3.2	Train Predictions: Diffusion Fractional Anisotropy	28
3.3	Validation Predictions: Diffusion Fractional Anisotropy	28
3.4	Test Predictions: Diffusion Fractional Anisotropy	29
3.5	Training Curve: Mean Diffusivity	30
3.6	Train Predictions: Mean Diffusivity	30
3.7	Validation Predictions: Mean Diffusivity	31
3.8	Test Predictions: Mean Diffusivity	31
3.9	Confusion Matrices (Precision): Relative Connectivity	33
3.10	Confusion Matrices (Recall): Relative Connectivity	34
3.11	Training Curve: Relative Connectivity	34
3.12	Train Predictions: Relative Connectivity	35
3.13	Validation Predictions: Relative Connectivity	35
3.14	Test Predictions: Relative Connectivity	36
A.1	Files: Raw	46
A.2	Class Diagram: Data Types	47
A.3	Class Diagram: Preprocessing	48
A.4	Class Diagram: Common	49
A.5	Files: Preprocessed	50
A.6	Class Diagram: Experimentation	51
B.1	Index Array Mapping	55

C.1	Slice: GLDM Small Dependence High Gray Level Emphasis	56
C.2	Slice: NGTDM Busyness	57
C.3	Training Curve: Subcortical	58
C.4	Train Predictions: Subcortical	58
C.5	Validation Predictions: Subcortical	59
C.6	Test Predictions: Subcortical	59

List of Tables

1.1	Regions Legend	9
2.1	Raw Data	12
2.2	Uniform Data	14
2.3	Radiomic Feature Types	16
3.1	Hyperparameters: Common	23
3.2	Hyperparameters: Subcortical	23
3.3	Hyperparameter Tuning: Subcortical	24
3.4	Missing Records	25
3.5	Feature Selection	38
4.1	Sustainability	40
4.2	Billing	41
A.1	Data Generator Properties	53
D.1	Voxel Based Radiomic Features	61
D.2	Shape Based Radiomic Features	61
D.3	Hyperparameter Tuning: Diffusion Fractional Anisotropy - Native T1	62
D.4	Hyperparameter Tuning: Diffusion Fractional Anisotropy - Native T1/T2	63
D.5	Hyperparameter Tuning: Diffusion Fractional Anisotropy - Normalized T1	64
D.6	Hyperparameter Tuning: Diffusion Fractional Anisotropy - Normalized T1/T2	65
D.7	Architecture Tuning: Diffusion Fractional Anisotropy	66
D.8	Hyperparameter Tuning: Mean Diffusivity - Native T1	67
D.9	Hyperparameter Tuning: Mean Diffusivity - Native T1/T2	68
D.10	Hyperparameter Tuning: Mean Diffusivity - Normalized T1	69
D.11	Hyperparameter Tuning: Mean Diffusivity - Normalized T1/T2	70
D.12	Architecture Tuning: Mean Diffusivity	71
D.13	Hyperparameter Tuning: Relative Connectivity - Native T1	72
D.14	Hyperparameter Tuning: Relative Connectivity - Native T1/T2	73
D.15	Hyperparameter Tuning: Relative Connectivity - Normalized T1	74
D.16	Hyperparameter Tuning: Relative Connectivity - Normalized T1/T2	75
D.17	Architecture Tuning: Relative Connectivity	76

Introduction

Basal ganglia is a part of the human brain which is a group of subcortical nuclei responsible primarily for motor control, as well as other roles such as motor learning, executive functions and behaviors, and emotions. [1] Huntington's disease is a disorder that causes the progressive degeneration of the basal nuclei. [2]

Hospital de Bellvitge provided an excellent dataset of anatomical magnetic resonance imaging (MRI) and diffusion tensor imaging (DTI) records of 32 control and 38 Huntington patient records of T1 and T1/T2 MRI images with isotropic voxels of 1 millimeter resolution and DTI fractional anisotropy (FA), mean diffusivity (MD) and radial diffusivity (RD) images with isotropic voxels of 2 millimeter resolution. Furthermore this dataset also contains the mask for the basal ganglia, which will also be referenced as the region of interest (ROI). Masks for the 7 main cortical regions of the brain, which will also be referenced as the target regions: Limbic, Executive, Rostral-Motor, Caudal-Motor, Parietal, Occipital and Temporal are also included in the dataset. Tractography was performed on the DTI images to figure out which parts of the ROI are connected to which cortical target, in a similar manner to how it was done in this paper [3]; where the relative connectivity maps are representing the ratio of the number of streamlines to each cortical target. Furthermore, the raw streamline images are also available, where there are a maximum of 5000 streamlines starting from each voxel in the ROI. The subcortical segmentation of the Basal Ganglia is also available, for the Caudate, Putamen and Accumbens on the control records. And lastly FMRIB's nonlinear image registration tool (FNIRT) warp fields were also provided for converting the records into normalized space.

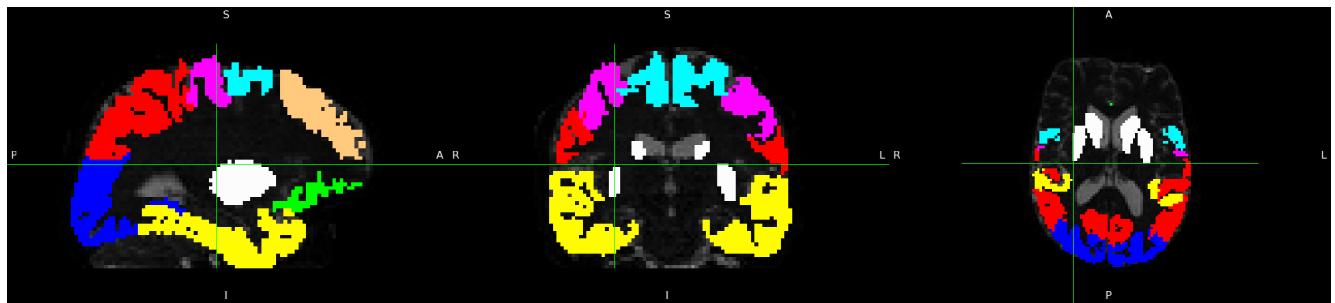


Figure 1.1: Basal Ganglia (ROI) & Cortical Targets

Color	Region
White	Basal Ganglia (ROI)
Green	Limbic
Brown	Executive
Light Blue	Rostral-Motor
Purple	Caudal-Motor
Red	Parietal
Blue	Occipital
Yellow	Temporal

Table 1.1: Regions Legend

Furthermore, for both the ROI and cortical targets, the dataset distinguishes between the right and left hemispheres of the brain. Thus there are actually 2 ROIs and $2 \cdot 7 = 14$ target regions.

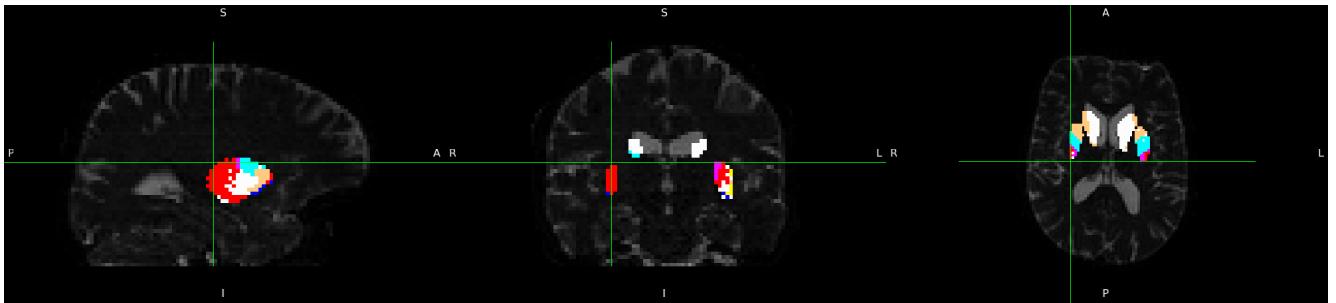


Figure 1.2: Connectivity Maps

1.1 Objectives

The end goal is to predict the relative connectivity of the Basal Ganglia to the cortical targets, from the radiomics features of the T1 and T1/T2 images.

This being a very complex problem, there is the possibility that the correlation between the connectivity of the brain and the T1, T1/T2 images are too weak to be mapped on this dataset. As from a data science perspective, 70 records are not much. But from a medical perspective it is substantial as it is very hard to collect uniform, clean data, with permissions to use it for research.

A simpler task leading up to the complex end goal, is a model for the simple segmentation of the Basal Ganglia for the subcortical regions Caudate, Putamen and Accumbens. In order to confirm that the radiomics texture of the T1 and T1/T2 images of this dataset are correlated to the segmentation of the Basal Ganglia. This problem is inherently connected to the main goal, as the relative connectivity does obey certain anatomical restrictions, and the subcortical segmentation of the Basal Ganglia is confirmed to be related to the relative connectivity. Thus if this simpler prediction fails, it is almost certain that the complex end goal will fail as well.

Another intermediate task, is a model for predicting FA and MD images. This is also related to the main goal, as these images are computed from the DTI images, the same image that the relative connectivity is computed from. But it is inherently simpler, not needing to perform complex algorithms like tractography.

The biggest obstacle of this project is the preprocessing of the data, as there are many variations and hyperparameters that can be tuned. An exhaustive search definitely will not be viable, thus

the preprocessing and model will needed to be tuned in a waterfall like manner, making educated guesses and comparing model performances across different tries. The main metric to measure model performance, will be the accuracy of the label prediction across voxels, as it should be comparable between all approaches. And pearson correlation will be used as the metric to evaluate the FA and MD regression predictions.

1.2 Motivation

The motivation for predicting the connectivity maps from the T1 and T1/T2 MRI images, is skipping the time and resource consuming process performing DTI and tractography.

1.3 State of the Art

Design

In order to understand some of the following design choices, it makes sense to establish it early that the model will be operating on extracted voxel based features and non-voxel based features, and will predict on a voxel by voxel level.

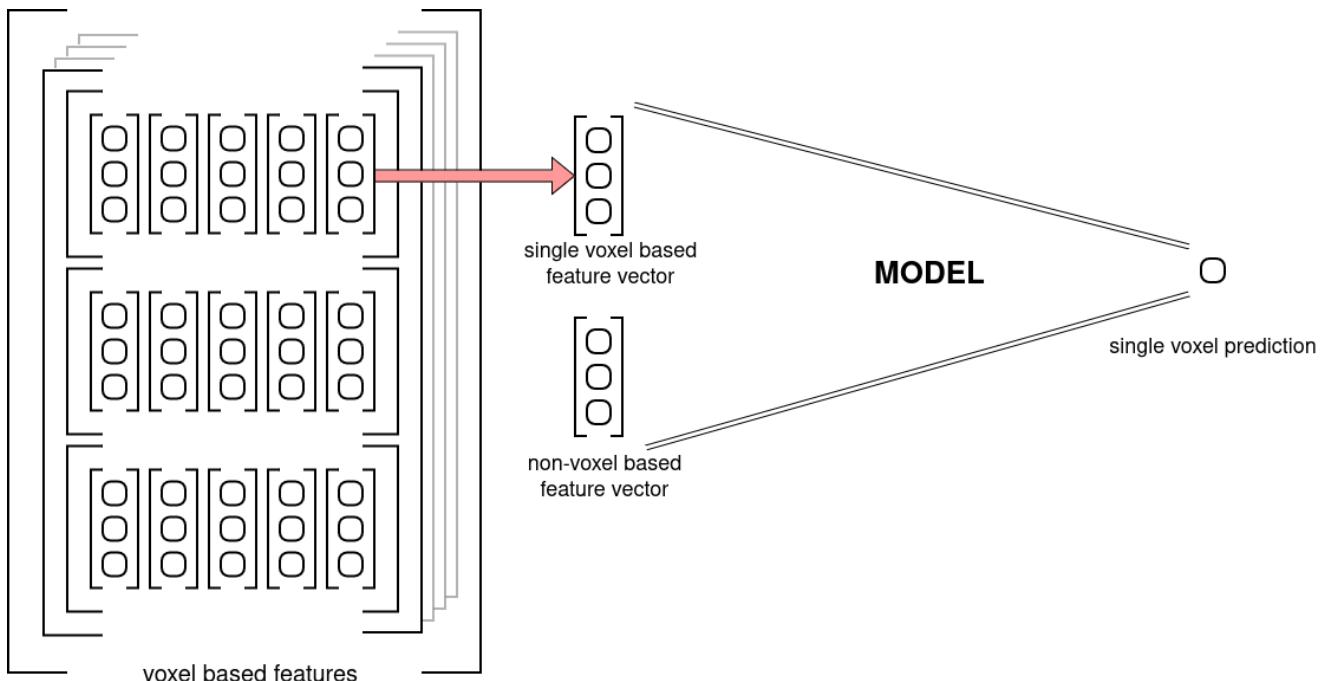


Figure 2.1: Simple Model Overview

This report will reference to the control/patient spatial data as '**Record**' ("voxel based features" in Figure 2.1) and will reference to the individual feature vectors as '**Datapoint**' ("single voxel based feature vector" and "non-voxel based feature vector" in Figure 2.1). This logical differentiation is needed, as the model only operates on Datapoints and has no global context available, while some preprocessing and evaluation logic should happen on a Record level.

2.1 Preprocessing

2.1.1 Raw Data

All provided records are in the neuroimaging informatics technology initiative (NIfTI) format, first these are need to be understood and parsed. This format stores the raw output of the MRI record, and additionally an affine transformation matrix used for aligning different spaces.

2.1.1.a Available Data

The following records will be preprocessed and read, even if not all of them are going to be used later on it helps providing the largest possible flexibility.

Data	Shape	Range	Type	Space	Reference
DTI	(118, 118, 60, 74)	[0, 4096]	uint	diffusion	diffusion
Diffusion FA	(118, 118, 60)	[0, 2]	float	diffusion	diffusion_fa
Diffusion MD	(118, 118, 60)	[0, 0.01]	float	diffusion	diffusion_md
Diffusion RD	(118, 118, 60)	[0, 0.01]	float	diffusion	diffusion_rd
T1	(208, 256, 256)	[0, 1000]	float	t1	t1
T1/T2	(208, 256, 256)	[0, 1]	float	d_aligned	t1t2
Cortical Targets	(118, 118, 60, 14)	{0, 1}	bool	diffusion	targets
Relative Connectivity	(118, 118, 60, 14)	[0, 1]	float	diffusion	connectivity
Streamline Image	(118, 118, 60, 14)	[0, 5000]	uint	diffusion	streamline
ROI Mask (Basal Ganglia)	(118, 118, 60, 2)	{0, 1}	bool	diffusion	mask_basal & roi
Brain Mask	(208, 256, 256)	{0, 1}	bool	t1	mask_brain
Basal Ganglia Segmentation	(208, 256, 256)	[0, 58]	uint	t1	basal_seg

Table 2.1: Raw Data

2.1.1.b Brain Mask

The provided dataset did not apply the brain masks for the T1 images out of the box so it can be done with a simple element wise multiplication of the T1 image and T1 mask.

2.1.1.c Registration

The process of aligning different records into the same native space is called "registration". The provided dataset comes with 2 (3) different spaces, earlier referenced to as t1 and diffusion (and d_aligned). Most of the data are in diffusion space, thus it is logical to register the rest into the same space. After manual inspection, only 15 records required registration. Out of which 3 only required a tiny translation, and the rest 12 needed a complete affine registration.

The image T1/T2 is the odd one out, as it is inherently in a different space from diffusion (due to them being different resolution). But they are aligned into diffusion space. Although they do not need to be registered, this has to be taken into account later on.

2.1.1.d Normalization

The process of warping each brain into a common space is called "normalization". Applying the FNIRT warp fields are more or less straight forward, as two warp fields are provided, one for the diffusion space and one for the T1 space. Note that this process inherently contains the benefits of registration, as it is warping the different images into a common brain shape and space. This also paves the direction of future experiments, as it opens the door to working in either native and normalized space.

The only encountered obstacle was with the T1/T2 image. As it is aligned in diffusion space, but FNIRT convention ignores the affine transformation of the NIfTI format, thus making it's registration useless as the raw data of the t1t2 has nothing to do with the raw diffusion data

(due to them being different resolution). The solution is to apply an affine matrix to t1t2's raw data which transforms it into t1's raw data space, after which the t1's FNIRT warp field can be applied to the t1t2 image. This affine transformation matrix can be easily calculated from the already given matrices. Let A denote T1/T2's affine matrix and B denote T1's affine matrix (after registration), thus the matrix which transforms the T1/T2 into T1 space is $M = A \cdot B^{-1}$.

2.1.1.e Basal Ganglia Segmentation

As the tractography of the brain is performed on the diffusion image, it inherently means that the connectivity maps and the roi are in diffusion space. But the basal ganglia's subcortical segmentation is in T1 space. This means that even if they are registered in the same space, they will not have a pixel perfect union due to the different resolutions.

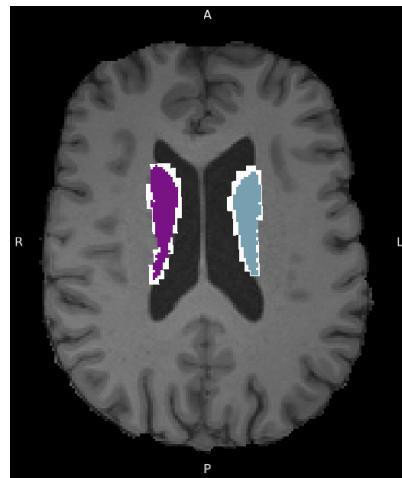


Figure 2.2: Basal Ganglia Subcortical Segmentation

The figure above visualizes the alignment of the Caudate subcortical region, where the white (larger) region is the Basal Ganglia mask from the diffusion space and the colored (smaller) regions are the Basal Ganglia segmentation from T1 space.

In order to keep the data consistent, mapping the segmentation to the Basal Ganglia mask can be done by assigning the same label for each voxel in the basal ganglia as the label of the closest voxel in the subcortical segmentation.

2.1.1.f N-Dim Array

The used NIfTI format stores the raw voxel space and the affine transformation matrix separately, in order to not loose data in the process of interpolating voxels when applying the transformation. But in order to consistently compare voxel data across different spaces (even if they are registered in the same space), the transformation needs to be applied, computing the interpolated voxels in the common space, bringing them into the same raw format of matching X, Y and Z dimensions, and discarding the stored affine matrices.

By default the native anatomical space's origin is near the center of mass of the brain, between the ears. This makes sense for medical professionals, when working with MRI records, but data-structure wise an array is indexed from 0. Meaning after applying the transformation to the voxel space, the yielded array will only contain one quadrant of the record as the rest are clipped in the

negative regions. Thus the space is also needed to be translated with the negative vector of the transformed space's bounding box's lower end.

The translation value can be calculated by calculating the boundaries of the transformed space's bounding box. Get all 8 corners of the voxel space and apply the transformation matrix to all of them. Then get the min-max coordinates along X, Y and Z from the 8 transformed vectors, yielding the lower and upper bounds of the transformed space's bounding box.

It is very important to use the same translation value across different spaces to properly align them in the native space. For example let D and T denote a diffusion and t1 records and M_D and M_T denote their respective transformation matrices. Let T_D and T_T denote their respective translation values. In order to properly align them we need to apply $A_D = (M_D \cdot \textcolor{red}{T}_D)$ matrix and $A_T = (M_T \cdot \textcolor{red}{T}_D)$ matrix to D and T respectively, with matching $\textcolor{red}{T}_D$ translation values.

The last issue is the missaligned shape of the dimensions of the T1 and diffusion records. This can be simply fixed by truncating the excess along each dimension.

2.1.1.g Uniform Shape

After aligning the data into the same space per record, it is still very likely that the individual records do not have a uniform shape. This is due to them being in native space, some records will contain a smaller volume brain, some will contain a larger, they will not be the same.

Due to the per-voxel based prediction model architecture this is not a problem, but fixing this for being able to use the data in a spatial model like a fully convolutional neural network (FCNN) can be simply solved by adding padding to the records in order to match their shapes.

Data	Volumes	Range	Type
diffusion	74	[0, 4096]	float16
diffusion_fa	1	[0, 2]	float16
diffusion_md	1	[0, 0.01]	float16
diffusion_rd	1	[0, 0.01]	float16
t1	1	[0, 1000]	float16
t1t1	1	[0, 1]	float16
targets	14	{0, 1}	bool
connectivity	14	[0, 1]	float16
streamline	14	[0, 5000]	float16
mask_basal	2	{0, 1}	bool
mask_brain	1	{0, 1}	bool
basal_seg	6	{0, 1}	bool

Table 2.2: Uniform Data

2.1.2 Quality Control

Having a low count of records means that if there are even just a few outliers, it can heavily affect the end result. Thus all data were manually inspected to make sure they are as clean as possible.

2.1.2.a Mismatched Data

Looking through the diffusion, diffusion_fa, diffusion_md and diffusion_rd images, 2 records' FA, MD and RD images were seemingly from completely different patients. Thus the FA, MD and RD images were omitted for 2 records.

2.1.2.b Garbled Data

Looking through the subcortical segmentation of the Basal Ganglia revealed that 1 record had a garbled segmentation. Thus, said basal_seg image was omitted for 1 record.

And one record had a garbled T1 FNIRT warp field. Said record was entirely omitted from the normalized set of records.

2.1.2.c Missing Data

Looking through the relative connectivity and streamline images, 3 records were missing these images, said 3 records were completely omitted, as these records are effectively missing the labels.

And the t1t2 images were missing for 10 records, but these were not omitted completely as the t1 images were present for these records, thus experiments only concerning the t1 can have a bit more available data.

2.1.3 Radiomics Features

Although the term is not strictly defined, radiomics generally aims to extract quantitative, and ideally reproducible, information from diagnostic images, including complex patterns that are difficult to recognize or quantify by the human eye. [4] Using these features is key, as there are not nearly enough data for neural network (NN) based features extraction such as a convolutional neural network (CNN).

Extracting the voxel based radiomic features has two main parameters to tune, the bin width and the kernel width. Where the binning parameter(s) influence how the intensity values of the image are binned, and the kernel size influences the size of the 'sliding window' similar to a convolution.

The two approaches for binning are absolute discretization and relative discretization. Where in the prior one, a fixed bin width is chosen and in the latter one, a fixed number of bins are chosen and the bin width scales relatively according to the min-max voxel values. This study found that "The absolute discretization consistently provided statistically significantly more reproducible features than the relative discretization." [5] Relying on this information, the obvious choice to start with is the absolute discretization.

The bin width and the kernel width will be tuned in later experiments. And possibly features calculated with different settings will be concatenated and used simultaneously for better results. The used default values will be 25 and 5 for the bin and kernel widths respectively.

The following types of radiomic features will be used:

Feature Type	Number of Features
first order	18
gray level co-occurrence matrix (GLCM)	23
gray level size zone matrix (GLSZM)	16
gray level run length matrix (GLRLM)	16
neighbouring gray tone difference matrix (NGTDM)	5
gray level dependence matrix (GLDM)	14
3D shape	17

Table 2.3: Radiomic Feature Types

2.1.3.a Voxel Based

The 92 features in Table D.1 will be calculated voxel based. Shape features do not make sense to calculate voxel based as it would just describe the shape of the used kernel, which is constant and independent from the input image.

2.1.3.b Non-Voxel Based

However, the additional shape features in Table D.2 do make sense for the non-voxel based features. As it can be computed for each target region, both hemispheres of the ROI and the entire brain.

2.1.4 Coordinates

One additional input that can be included in the experiments is the coordinates. Although this approach only makes sense in normalized space, where the images from different records are aligned. This theoretically would allow the model to learn certain anatomical markers based on the location of the voxel, adding a type of global context to the input of the model.

Furthermore, this approach can be adopted to the native space, by constructing the normalized coordinate map and then 'de-normalizing' them with an inverse FNIRT warp field.

2.1.5 Data Augmentation

The only data augmentation that makes sense involves applying small rotation values to the input images in their native space before calculating radiomic features. Applying transformations to the already extracted features is illogical, as interpolating between voxels in feature space is unlikely to yield the same results as computing features after transforming the input images. In summary, any spatial data transformations should be performed upstream. Furthermore, data augmentation only makes sense in native space, as by definition such transformations would make the normalized image pointless.

2.1.6 Scaling and Normalization

As the extracted features have very different ranges, it makes sense to follow the standard practice of scaling the data to a fixed range. Inspecting the histogram of some of the radiomic features reveals that most of them follow a bell curve with moderate standard deviation, such as Figure 2.3 (Firstorder Energy).

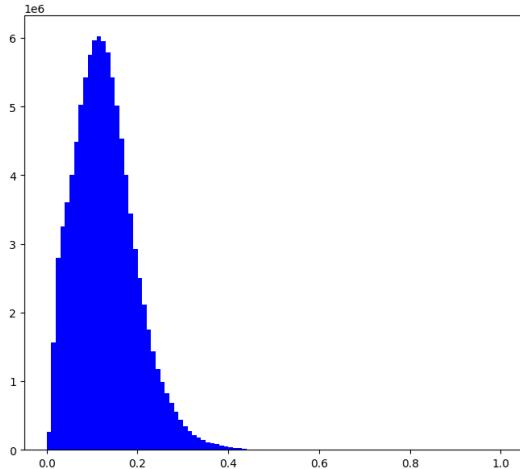


Figure 2.3: Histogram: Firstorder Energy

However, some other features like Figure 2.4 (GLDM Small Dependence High Gray Level Emphasis) and Figure 2.5 (NGTDM Busyness) have a very skewed distribution, the latter one being the most extreme case. This skewing can be mitigated by applying logarithm to the offending features.

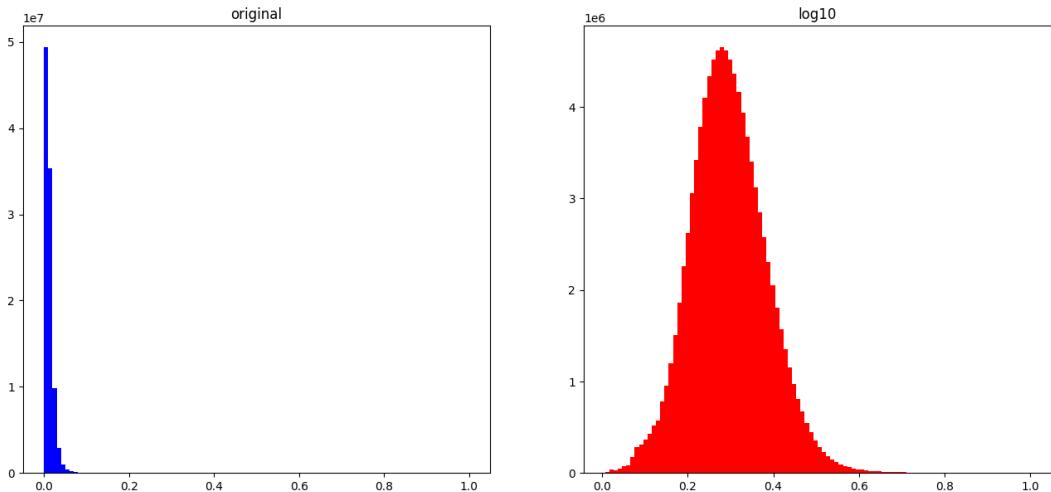


Figure 2.4: Histogram: GLDM Small Dependence High Gray Level Emphasis

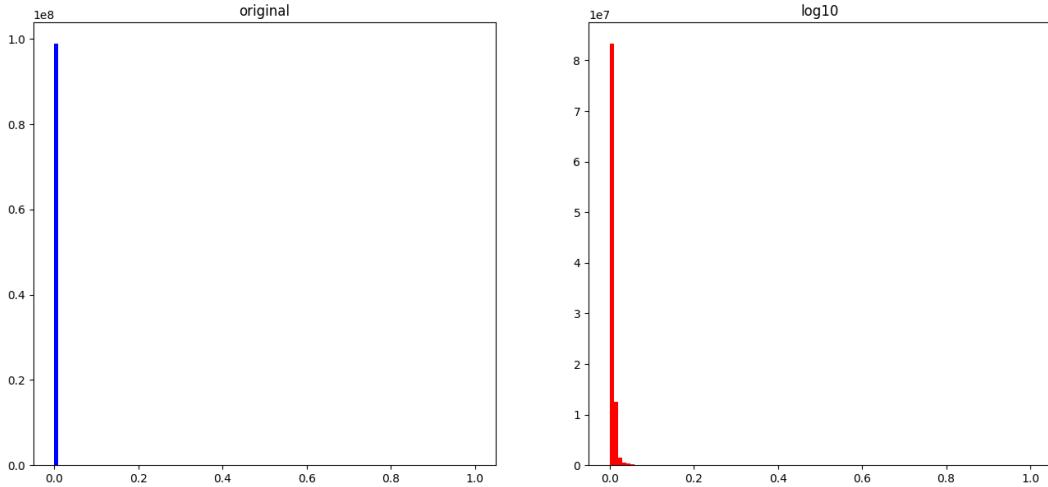


Figure 2.5: Histogram: NGTDM Busyness

Besides the standard benefits of making the optimization process more stable and efficient, and reducing the sensitivity to outliers. It also have some less evident benefits.

Although it is very subtle, but storing these records in float16 inherently loses some information. This loss is not a problem for the features that have a healthy distribution, but in the more extreme cases it can cause compression artifacts visible even to the naked eye, such as the very subtle loss of detail in Figure C.1. And in the most extreme case it can even render the entire feature useless like in Figure C.2. While the normalized features have no problem storing this fine detail in float16.

This makes the system much more robust from a practical perspective; as depending on the hardware, some GPUs are much more efficient at computing in float16. And it also halves the memory and storage requirements, as in float32 a sinlge MRI image of 92 volumes (for the 92 features) takes up around 1GB of space.

Selecting which features need normalization is done programmatically, and the exact selection criteria is detailed in Appendix B.4.

2.1.7 Data Balancing

Working with highly unbalanced data can be challenging, and balancing it does not necessarily going to help the model's generalization capability. Thus, a method for partially balancing the data will be used, where the bins of the unbalanced data will be up-sampled by a ratio of the difference of the number of datapoints in the bin (compared to the bin with the maximum number of datapoints). Figure 2.6 demonstrates how a ratio 1 means perfectly balanced data, 0 means unbalanced data. And how the ratios in between are approximately preserving the shape of the distribution and partially balance the data.

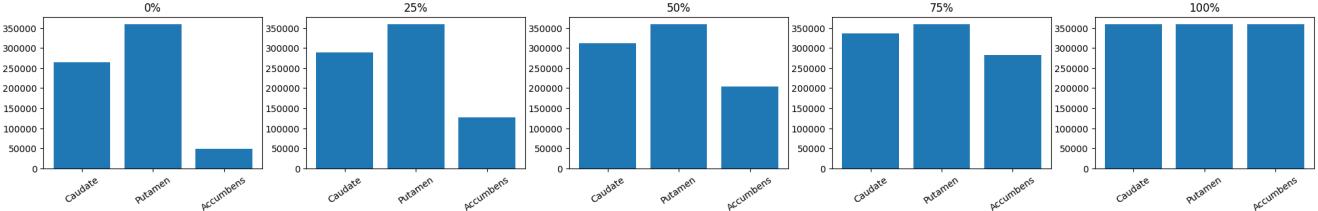


Figure 2.6: Balance: Subcortical

For the diffusion_md and diffusion_fa, which are regression problems and have continuous labels, binning can be used to create artificial groups which can be balanced.

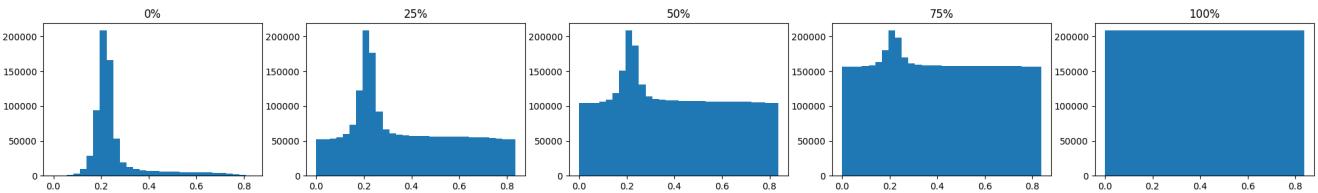


Figure 2.7: Balance: Diffusion MD

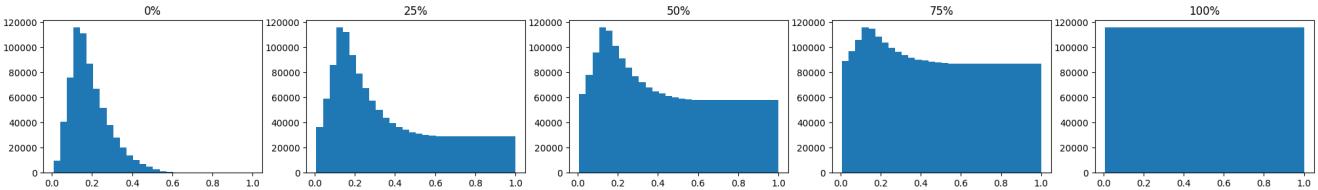


Figure 2.8: Balance: Diffusion FA

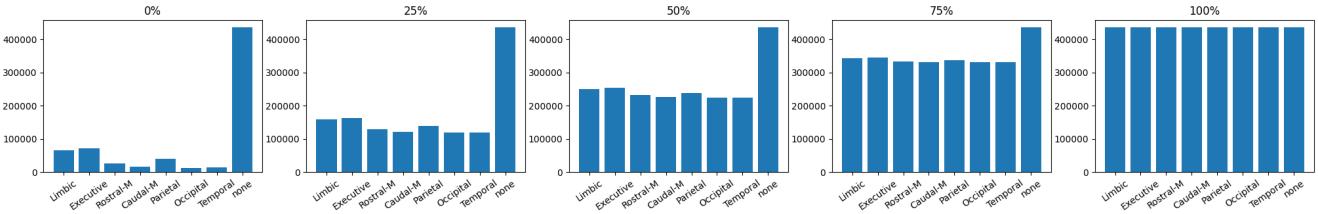


Figure 2.9: Balance: Relative Connectivity (thresholded at 0.6 & binarized)

2.1.8 Clinical Data

There are additional clinical data available for the Patient records. Disease severity can be characterized in terms of CAG Age Product (CAP) score. Providing a measure of cumulative exposure to the mutant HTT gene. [6] This widely accepted and used CAP score is available for all patients.

Another, newer metric for characterizing disease severity is the composite Unified Huntington's Disease Rating Scale (cUHDRS) [7]. This is calculated from 4 other basic metrics: Total Functional Capacity, Total Motor Score, Symbol Digit Modalities Test and Stroop Word Reading. These are available for most patients, with a handful exceptions.

And there are a total of 91 available clinical features, with relatively a lot of missing data on some of these features. There are 8 additional patients available in the clinical data. These can be

used to aid the data imputation for the missing values, and can be omitted afterwards, as these have no corresponding MRI records.

All clinical features were scaled the range of 0-1 with min-max scaling (per feature). And euclidean distance was used for the following imputation process. The imputation strategy itself consisted of 2 steps, first the few missing cUHDRS values were imputed from the CAP score. And then the remaining features were imputed from the combined CAP and cUHDRS values.

2.1.9 Relative Connectivity

Relative Connectivity describes the ratio of the number of streamlines going into each cortical target. This means that the Streamline record can be converted into Relative Connectivity by simply dividing by the total number of streamlines in each voxel. With the additional inbetween step of filtering some noise, by thresholding the streamlines at 250 (5% of the total 5000 streamlines), meaning any voxel which has less than 250 streamlines to a target region are set to zero.

Then the relative connectivity could be converted into a label, by picking the label of the cortical target to the highest connection per voxel. However this would yield very noisy labels, as these ratios can be quite balanced between the multiple cortical targets, for example voxels with ratios of 0.31/0.29/0.3/0.1. To mitigate this, the relative connectivity can be thresholded at a value higher than 0.5, meaning a label can only be picked for a voxel if at least half of the connections are going to a single target.

But this also means that there can be voxels without labels. This can be dealt with introducing an artificial 'Not Connected' label for these voxels.

To achieve the best results and filtering, 0.6 was chosen for the thresholding value, as it also filters potential 50-50 situations and only allows labeling strong connections.

2.2 Evaluation

2.2.1 Train, Validation and Test Splits

There are 2 important aspects when splitting the data into Train/Validation/Test groups. In order to truly validate the model's generalization capability, the split must happen on a record level and not on a datapoint level. This means that our model can only learn on certain records, and it can be validated on records that it never seen before, not even partially. This has the consequence of that the split will not follow the defined ratio on a datapoint level, as it could happen that by pure chance the train split contains records with larger volumes, resulting in having a bit more datapoints than the validation split.

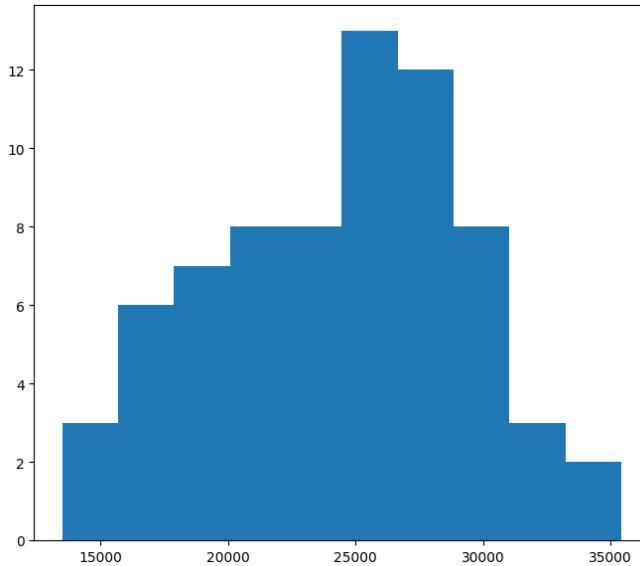


Figure 2.10: Distribution of Records in relation to Datapoints

In practice the lower end of datapoint count per record is around half of the higher end. Figure 2.10 shows the distribution of both Control and Patient records and both Left and Right datapoints. During experimentation, with 0.8 Train split and 0.5 Validation/Test split, the datapoint ratios stayed in the range of 0.8 ± 0.02 and 0.5 ± 0.06 for the two split ratios.

Furthermore to avoid introducing bias in the case of experiments with mixed Control and Patient records, the ratio of Controls/Patients must be constant across the different splits. This is required as Controls and Patients can have vast differences due to neurodegeneration. An extra caveat is having another ratio that must also be kept constant for the same reason, which is the symptomatic and asymptomatic patients, as they can also have vast differences due to different stages of neurodegeneration.

2.2.2 Accuracy and Pearson Correlation

There will be 3 groups of metrics for evaluating each model. First is the 'raw' (will also be referenced as 'train') metric group, which is computed on the datapoints that were extracted with the same hyperparameters as the datapoints during the training process. Meaning that this metric group best reflects how the model performs on the different splits, such as if the model was trained with a balancing of 0.5, all splits will be balanced the same way and the metrics will be computed on a datapoint level (the same way as how it is naturally computed in the loss function).

The second and third groups are for comparing model performances in between models and are for practical evaluation. The difference is that the metrics in this case are computed for each record, and then averaged out inbetween records. This means that it is computed on a record level instead of a datapoint level, resulting in the elimination of potential bias coming from the deviation from the number of datapoints per records. It also means that these metric groups will inherently ignore data balancing, as it operates on a record level.

And the 2nd metric group is computed in native space, while the 3rd is computed in normalized space. This means that if the model operates in native space, the normalized metrics will be computed by predicting the datapoints for each record, then the spaital record is reconstructed from the datapoints and warped to normalized space, and then the datapoints are extracted from

the normalized spatial prediction, and compared against the normalized labels (This process would be computationally quite expensive, so the implementation does not follow this exact logic, but numerically it is doing the same; more information on this in Appendix B.3). This way the models can have comparable metrics even if they operate in different spaces.

Experiments

The following hyperparameters were constant during all of the experiments:

Hyperparameter	Value
Train Split	0.8
Validation/Test Split	0.5
Model Type	feedforward neural network (FNN)
Optimizer	Adam

Table 3.1: Hyperparameters: Common

3.1 Subcortical Segmentation

This simple problem did not need a lot of tuning, as it was working very well almost from the start. The following set of hyperparameters were constant during these experiments:

Hyperparameter	Value
Control/Huntington Datapoints	Control Only
Left/Right Hemisphere Datapoints	Both
Space	Native
Image	T1
Scaling/Normalization	Normalized Voxel Based Features
Hidden Layers	$1024 \rightarrow 512 \rightarrow 256 \rightarrow 128$
Loss	Categorical Crossentropy
Activation	Sigmoid (softmax for the output layer)
Learning Rate	0.001
Batch Size	10000
Early Stopping Patience	7

Table 3.2: Hyperparameters: Subcortical

The reasoning behind the initial choices of these parameters are straight forward. The T1 image and native space were chosen, because those are the simplest to acquire in practice. Thus if the model is doing great on those, there is no need for more complicated inputs. Including both hemispheres would hopefully result in a model which can generalize better. Only using control datapoints should translate into less variance between the general characteristics of the datapoints, as it does not contain patients with neurodegeneration. The number and sizes of the hidden layers were chosen based on the potential size of the input layer, which should range from 92 (single set of voxel based features) up to $\sim 1000 - 2000$ including many different kernel sizes and non-voxel based features as well. Categorical crossentropy loss function and the output layer's softmax activation function are standard practices for a classification problem. The Sigmoid activation function should work fine without having to deal with exploding gradients and dieing

relu problems. Learning rate is the default learning rate of the Adam optimizer in TensorFlow. And a batch size of 10,000 seems appropriate for a train split size of 1,000,000 datapoints. And the early stopping patience of 7 epochs should also be good enough to prevent overfitting and stop the training in time, but it will be evaluated based on the learning curves and the accuracy of the model.

The used metric for evaluating model performance on the Train/Validation/Test splits is Accuracy. The 'k' and 'b' notations stand for kernel and bin, where k5 means a kernel width of 5mm, b25 means an absolute bin size of 25, and b10r means relative binning with 10 bins. And in the case of multiple kernel sizes denoted by a dash, it naturally only means odd kernel sizes. Tuning the rest of the hyperparameters were done in the following experiments:

	Experiment	Train	Val	Test	Input Layer
1.	Voxel Features k5_b25	68.9	69.1	72.4	92
2.	<i>Voxel Features k5_b25</i> Non-Voxel Features of Target Regions b25	73.2	68.9	72.5	1576
3.	<i>Voxel Features k5_b25</i> Non-Voxel Features of ROI b25	75	74.3	78.5	304
4.	<i>Voxel Features k5_b25</i> <i>Non-Voxel Features of ROI b25</i> Non-Voxel Features of Brain b25	74.5	70.4	70.3	410
5.	<i>Voxel Features k5_b25</i> Non-Voxel Features of ROI b10 b25 b50 b75	71.2	70.6	74.3	856
6.	Voxel Features k5_b25 - k21_b25 <i>Non-Voxel Features of ROI b25</i>	94.5	94.1	95.1	1040
7.	Voxel Features k5_b25 - k21_b25	95	94.6	93.7	828
8.	<i>Voxel Features k5_b25 - k21_b25</i> <i>Non-Voxel Features of ROI b25</i> Balance Ratio 0.5	94.9	94.4	94.9	1040
9.	<i>Voxel Features k5_b25 - k21_b25</i> <i>Non-Voxel Features of ROI b25</i> Balance Ratio 1	95.9	95.5	95.9	1040

Table 3.3: Hyperparameter Tuning: Subcortical

The best performing model was in experiment number 9, where it achieved a 96% accuracy with practically no overfitting. The biggest improvement during the experiments was to include many different kernel sizes for the voxel based features. The additional non-voxel based features of the ROI yielded a small improvement. And balancing the data yielded a marginal improvement, by reducing overfitting.

Examples of the true/predicted records can be found in Figures C.4 C.5 C.6. And the loss training curves can be found in Figure C.3.

3.2 Methodology

The experimentation from this point on, will be divided into 4 main groups:

- Native - T1

- Native - T1/T2
- Normalized - T1
- Normalized - T1/T2

The same set of core experiments will be run for all 4 groups, and some additional experiments will be run per group, depending on how they perform. The experiments will consider the following aspects:

- Single/Many Different Kernel Sizes for Voxel Based Features
- Additional Non-Voxel Based Features
 - Single/Many Different Bin Sizes
- Control/Patient/Both Records
- Left/Right/Both Hemisphere Datapoints
- Additional Clinical Features for Patient Records
- Additional Coordinate Map Features
- Scaled Voxel Based Features (not normalized)
- Different Bin Sizes for Voxel Based Features
- Different Balance Ratios

3.2.1 Missing Records

In order to be completely fair when comparing model performances, only records should be used which are available for all 4 groups of experiments. In practice the following records were missing:

Record	Missing Amount
Normalized	1
T1/T2	10
Diffusion FA & MD	2

Table 3.4: Missing Records

This meant that for the Diffusion FA & MD experiments there were a total of 13 records omitted, yielding 57 records in total, out of which 29 are Control and 28 are Patient records. And for the Relative Connectivity experiment, 11 records were omitted, yielding 59 records in total, out of which 30 are Control and 29 are Patient records.

As additional experiments for the groups with more available data (such as T1, where 10 more records could be included), these records can be appended to the train split on the best performing model, feasibly increasing model performance.

3.2.2 Architecture Tuning

For the best performing model, the architecture will be further tuned, considering the following aspects:

- Number of Layers and Layer Sizes
- Activation Function
- Batch Size
- Learning Rate
- Dropout Normalization
- Early Stopping Patience

3.3 Diffusion Fractional Anisotropy Regression

All numerical results of the experiments can be found in Tables D.3 - D.7. The baseline starting experiment is trying to predict the FA from a single set of voxel based radiomic features, with a kernel size of 5. And with the same starting hyperparameters (Table 3.2) that were also used in the subcortical segmentation (with exception of the used loss function, which is Mean Squared Error instead of the Categorical Crossentropy).

The next few experiments were trying to determine how does each set (target regions, roi, and entire brain) of non-voxel based features affect the model performance. Between the 4 different group of experiments (Native-Normalized & T1-T1/T2) the observations were more or less consistent, with the final consensus being that the inclusion of the entire brain's non-voxel based features are yielding the best results, with an improvement of 5-10% better correlation compared to the baseline.

Including many different bin sized non-voxel based features worsened the model performance by 0-3%.

The biggest improvement was the inclusion of many different kernel sized voxel-based features, with an improvement of 10-15%. And surprisingly after removing the non-voxel based features, it further improved the performance of the T1 experiments by 1-2%, while worsening the T1/T2 experiments by 0-1%. Running this experiment on the Patient records, resulted in the models performing even worse with the additional non-voxel based features by 8-10%.

The experiments consistently showed the model performing much better on the Control records, compared to the Patient records, with much less overfitting and better correlation by 5-10%.

The inclusion of the clinical features were behaving inconsistently between the 4 groups of experiments. For the native T1, including the CAP and cUHDRS features marginally improved the model performance, and for the normalized T1/T2 it improved model performance by 4-5%. While for the native T1/T2 and normalized T1, it worsened the model performance by 5-10%. The overall Patient records even with the best performing clinical features, were still performing worse than the Control records.

As expected, mixing Control and Patient records were performed worse than Control records only, but only with 1-5% correlation.

Including coordinates, did not affect the T1 models' performance, but it did marginally increase the T1/T2 models' performance by 1-2%.

Only using min-max scaling, and not normalizing the datapoints, resulted in marginally worse performance.

Increasing the bin size for the voxel based radiomic features marginally decreased the model performance.

Balancing the data was a bit inconsistent between the groups of experiments, but the balance ratio of 1 usually resulted in a marginally worse, and a balance ratio of 0.5 resulted in a marginally better performance.

Re-including the 10 extra T1 records as part of the training split for the T1 experiment, only resulted in a marginal improvement for the native space, and a 2% improvement for the normalized space.

After combining all of the best configurations, the best performing model was the T1 normalized model, with Control records only, and re-included T1 records, without any additional non-voxel based features. It reached a final correlation of **84.4/84.6/82.8** for the train/val/test splits in native space, and **84.6/84.9/82.9** in normalized space.

After tuning the model architecture, by searching different layer sizes and numbers, activation functions, dropout normalization, adjusting learning rate and batch size, it only increased the model's overfitting, without any actual benefits.

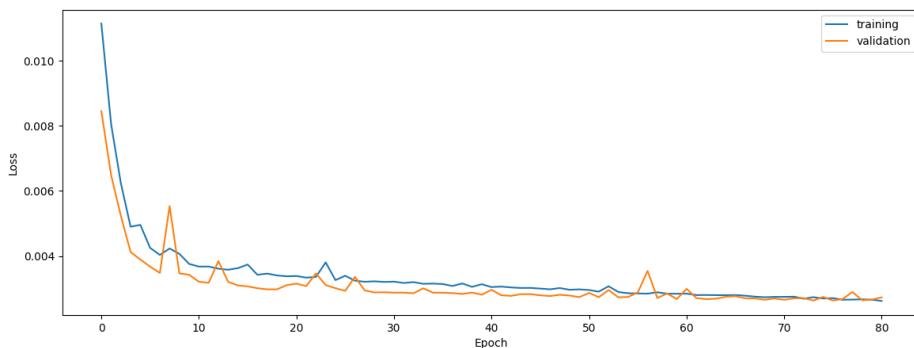


Figure 3.1: Training Curve: Diffusion Fractional Anisotropy

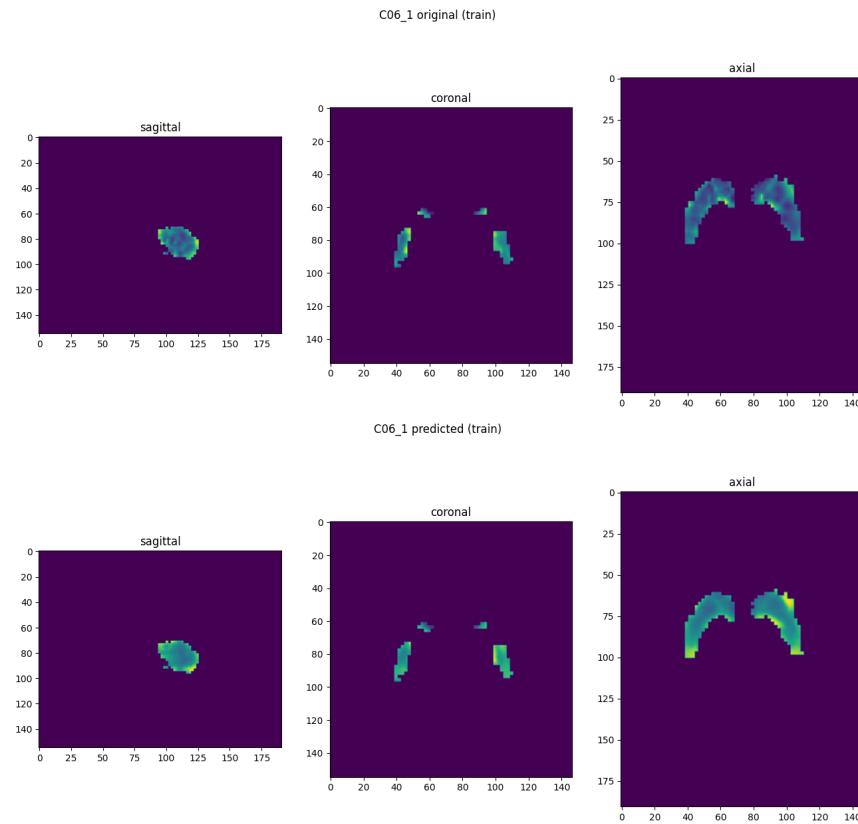


Figure 3.2: Train Predictions: Diffusion Fractional Anisotropy

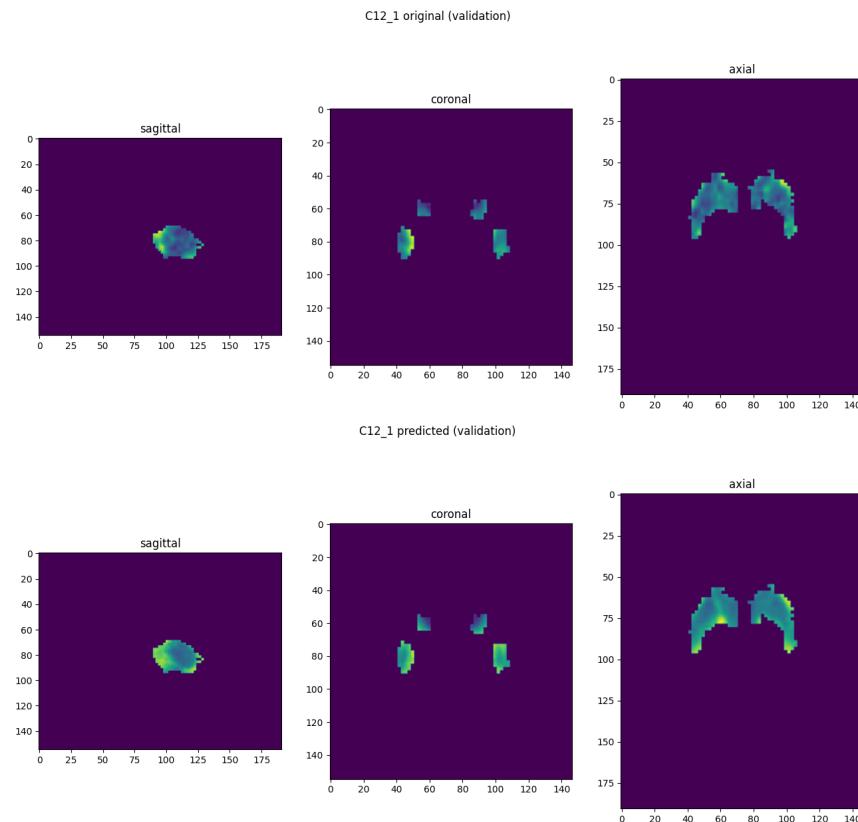


Figure 3.3: Validation Predictions: Diffusion Fractional Anisotropy

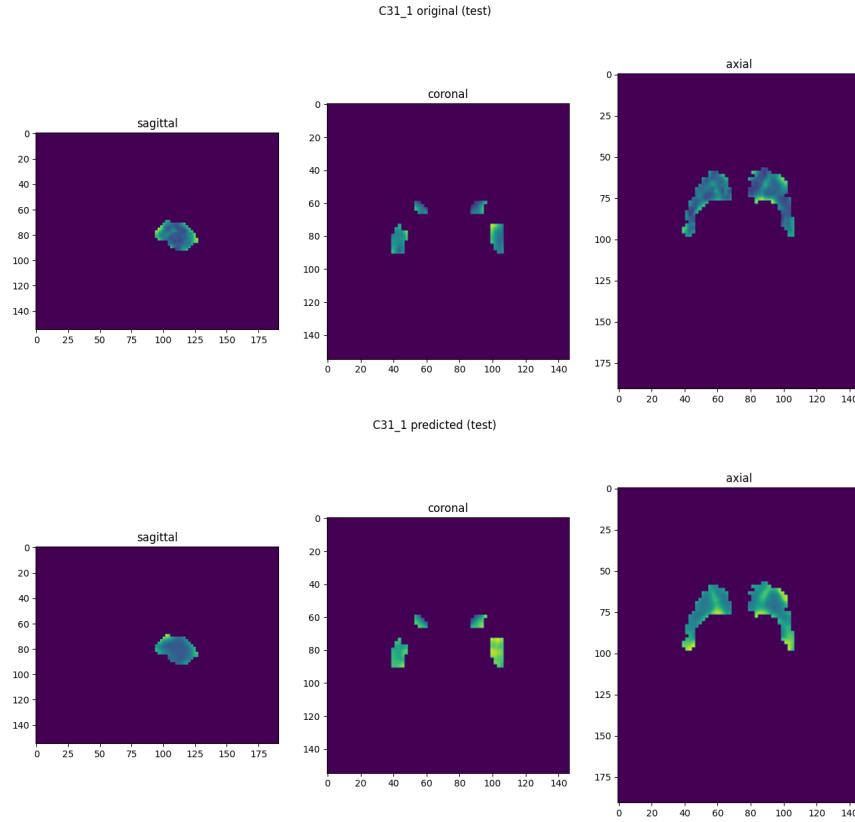


Figure 3.4: Test Predictions: Diffusion Fractional Anisotropy

3.4 Mean Diffusivity Regression

A similar set of experiments were run for predicting MD with the numerical results available in Tables D.8 - D.12. But these experiments were performing very well out of the box, even the baseline Native T1 experiment with a single set of voxel based features resulted in a correlation of 94% without any overfitting.

No significant observation can be made here, besides the Patient records performing marginally worse.

The best performing model was Native T1, on the Control records only, with the additional non-voxel based features of the entire brain, and many different voxel based kernel sizes. It reached a final correlation of **94.7/95.5/95.1** for the train/val/test splits in native space, and **95.4/95.7/96.3** in normalized space.

After tuning the model architecture, by searching different layer sizes and numbers, activation functions, dropout normalization, adjusting learning rate and batch size, it could not increase the model performance, not even on the train split. Indicating that this is the absolute best this model can do.

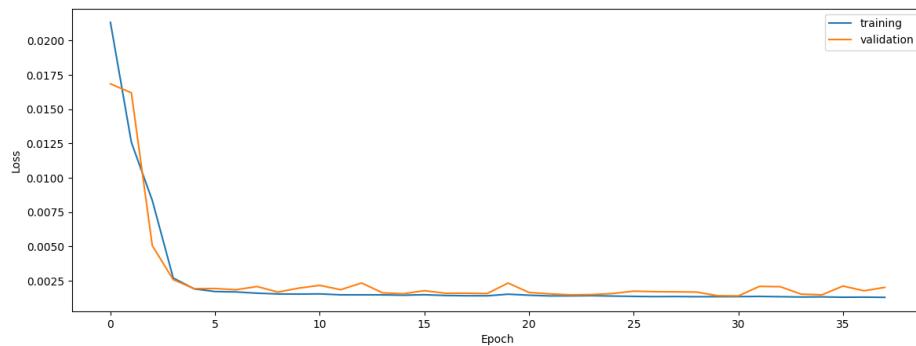


Figure 3.5: Training Curve: Mean Diffusivity

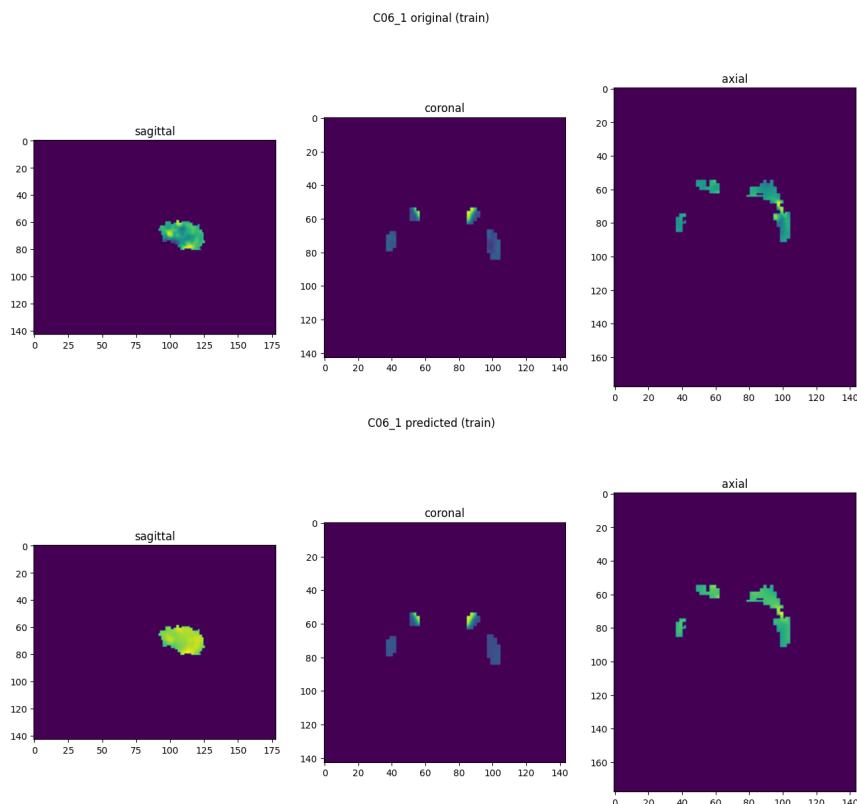


Figure 3.6: Train Predictions: Mean Diffusivity

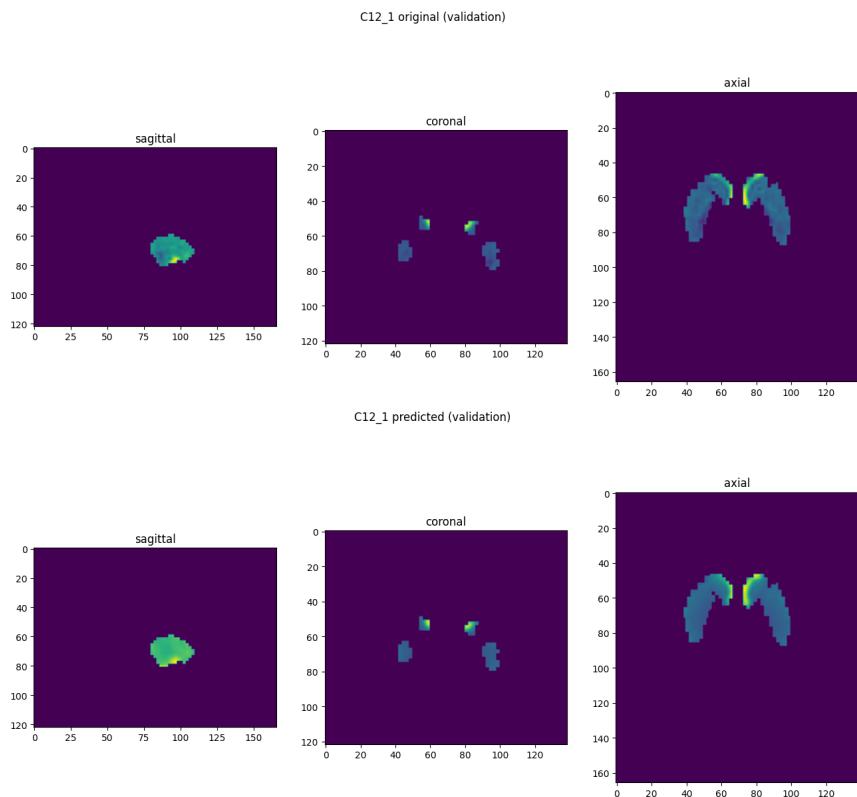


Figure 3.7: Validation Predictions: Mean Diffusivity

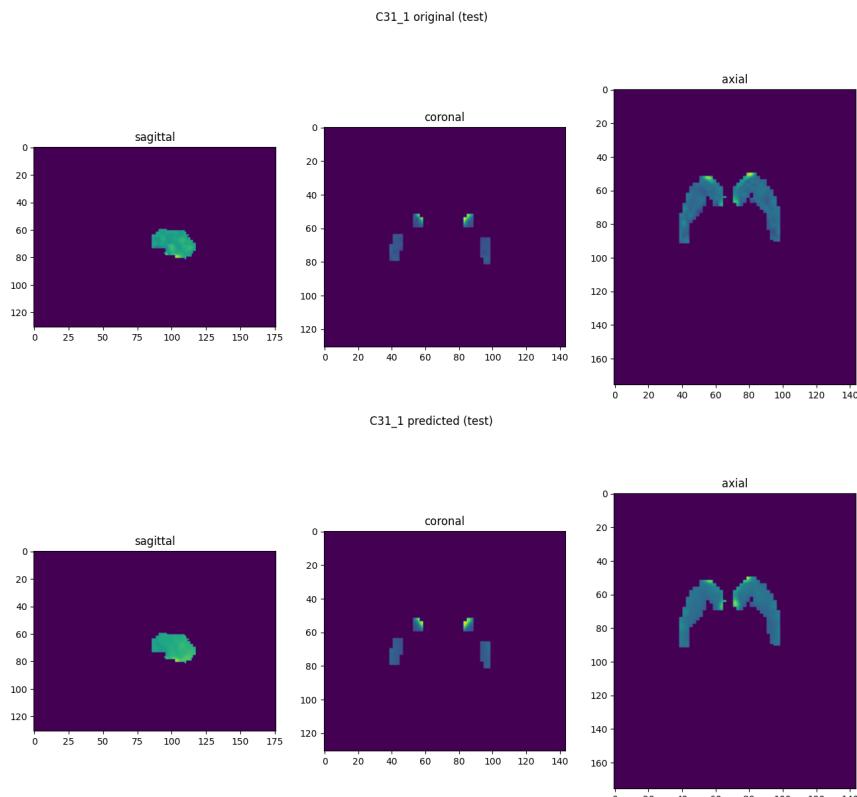


Figure 3.8: Test Predictions: Mean Diffusivity

3.5 Relative Connectivity Segmentation

All numerical results of the experiments can be found in Tables D.13 - D.17. The baseline starting experiment is trying to predict the Relative Connectivity (preprocessed with the method described in subsection 2.1.9) from a single set of voxel based radiomic features, with a kernel size of 5. And with the same starting hyperparameters (Table 3.2) that were also used in the subcortical segmentation.

The next few experiments were trying to determine how does each set (target regions, roi, and entire brain) of non-voxel based features affect the model performance. Between the 4 different group of experiments (Native-Normalized & T1-T1/T2) the observations were more or less consistent, with the final consensus being that the inclusion of non-voxel based features does not increase model performance.

The biggest improvement was the inclusion of many different kernel sized voxel-based features, with an improvement of 5-10%.

The experiments consistently showed the model performing much better on the Control records, compared to the Patient records, with much less overfitting and better accuracy by 2-5%.

The inclusion of the clinical features were behaving inconsistently between the 4 groups of experiments. For the normalized experiments, it mostly had no effect, or made it marginally worse. And for the native T1 experiments, including CAP and cUHDRS features marginally improved the model performance. And for the native T1/T2 it increased train/validation performance substantially, but yielded worse test performance, overfitting a lot.

Mixing Control and Patient records only performed marginally worse than using only Control records.

Including coordinates, consistently increased accuracy by 1-2%.

Only using min-max scaling, and not normalizing the datapoints, resulted in marginally worse performance.

Increasing the bin size for the voxel based radiomic features marginally decreased the model performance.

Balancing the data consistently decreased performance by quite a lot, depending on the balance ratio by 5-20%.

Re-including the 10 extra T1 records as part of the training split for the T1 experiments did not affect the model performance.

After combining all of the best configurations, the best performing model was the T1 normalized model, with Control records only, with the additional coordinate inputs, without any additional non-voxel based features. It reached a final accuracy of 72.6/72.1/73.3 for the train/val/test splits in native space, and 73.2/71.9/73 in normalized space.

After tuning the model architecture, by searching different layer sizes and numbers, activation functions, dropout normalization, adjusting learning rate and batch size, the only thing which marginally increased model performance was lowering the batch size to 10^3 and lowering the learning rate to 10^{-4} , yielding a final accuracy of **73.3/72.9/73.4** in native space, and **73.5/72.3/73.4** in normalized space.

As these numbers can be misleading due to the highly unbalanced data, and the best way to get more insight on how the model is performing, is by observing the confusion matrices in Figure 3.9 and Figure 3.10. Where matrices in the prior Figure 3.9 are normalized along the predicted label axis, effectively displaying the precision on the diagonals; and in the latter Figure 3.10 are

normalized along the true label axis, effectively displaying the recall on the diagonals. The first and most evident observation is that the unbalanced nature of the data is reflected on the confusion matrix, as the over represented 'not connected' datapoints have a much better precision and recall than the rest.

Also, the model is more effective at minimizing false positives, than false negatives, since it generally has a higher precision than recall for practically all labels (except the 'not connected').

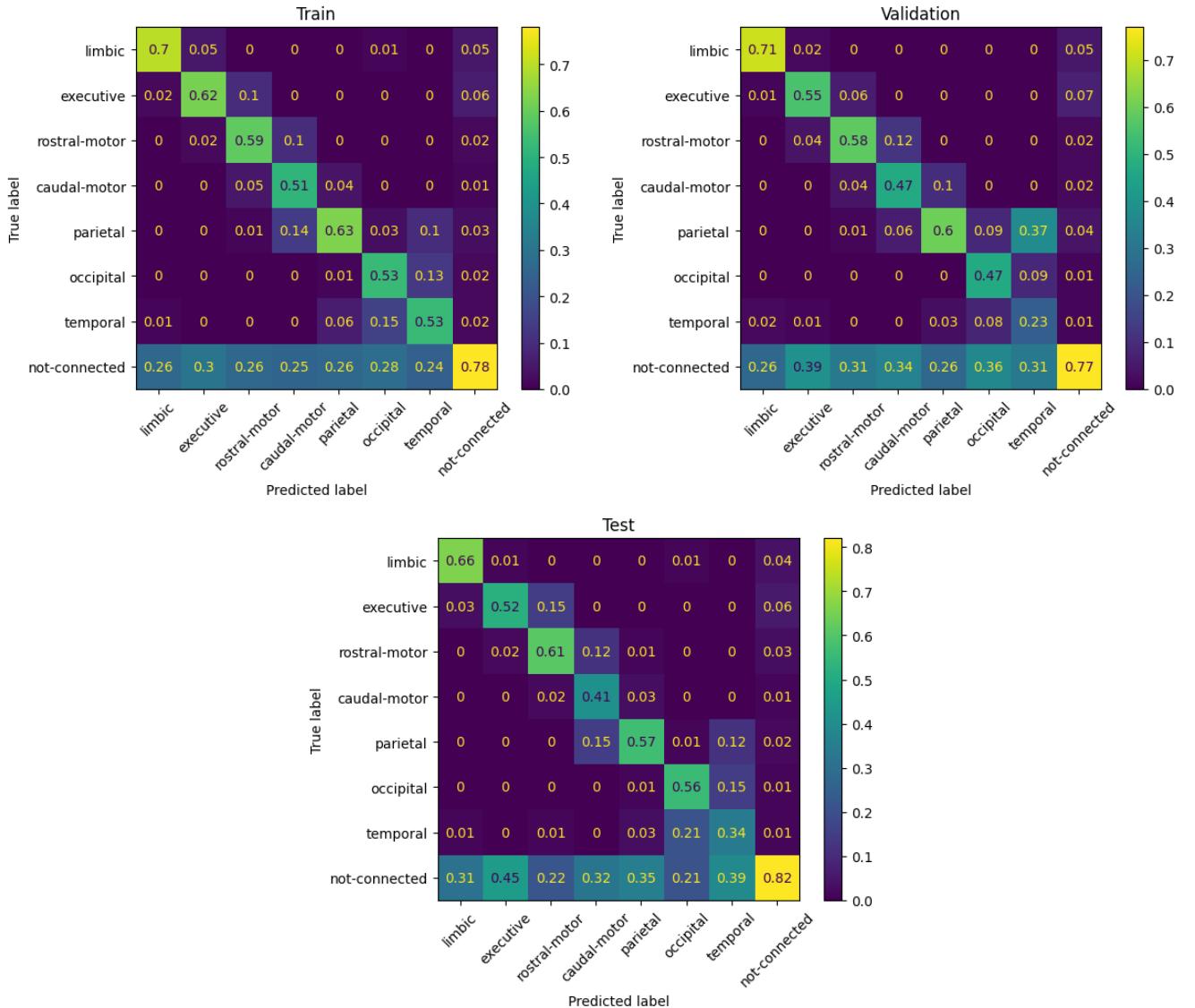


Figure 3.9: Confusion Matrices (Precision): Relative Connectivity

These matrices also tell that not all labels are performing equally. As the model clearly struggles with the recall of 'temporal' target region datapoints.

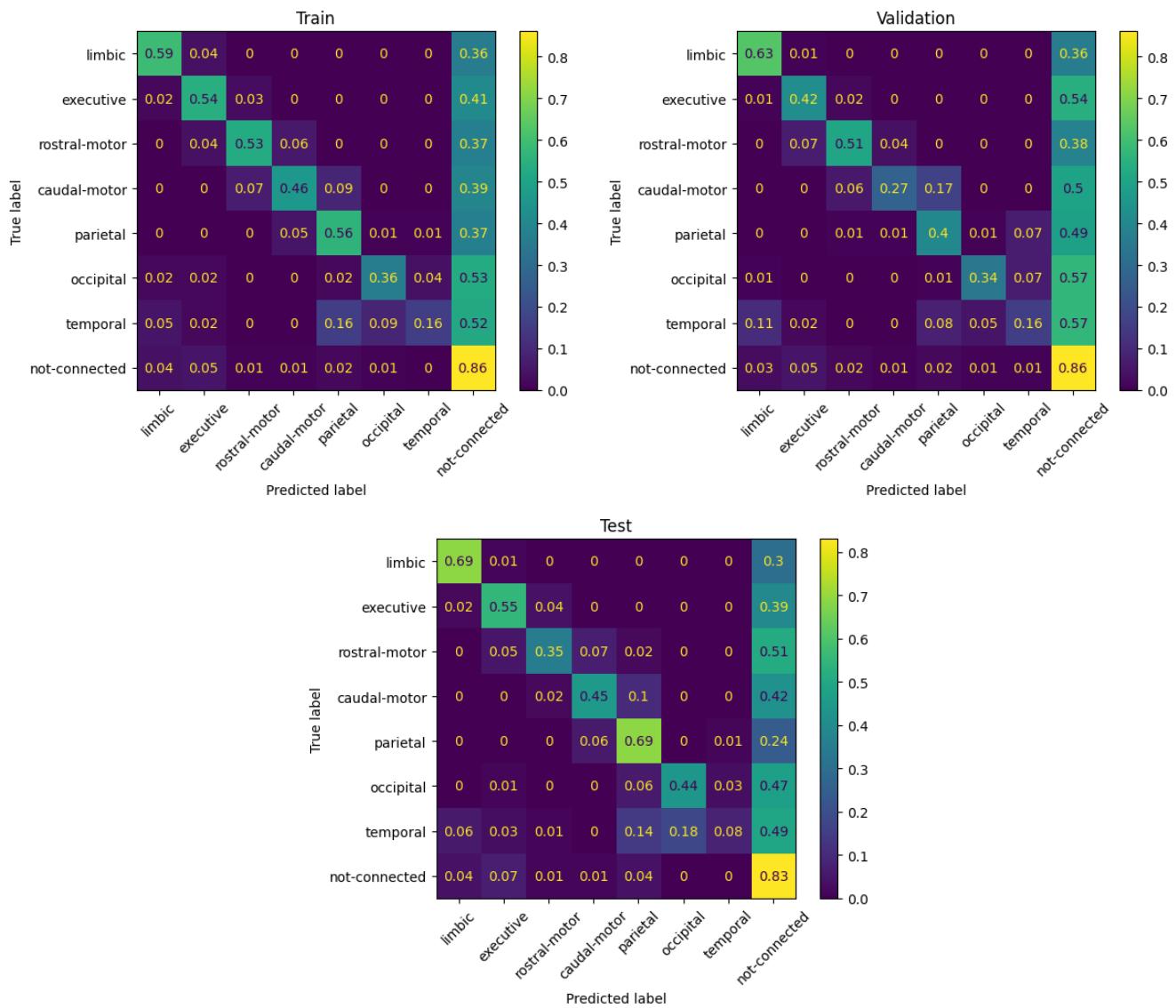


Figure 3.10: Confusion Matrices (Recall): Relative Connectivity

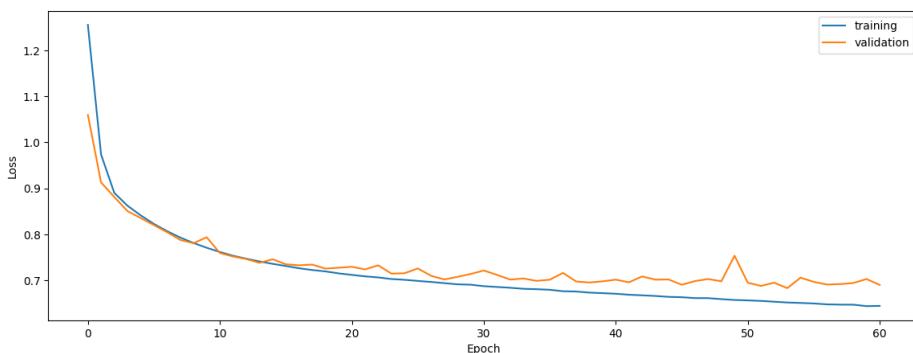


Figure 3.11: Training Curve: Relative Connectivity

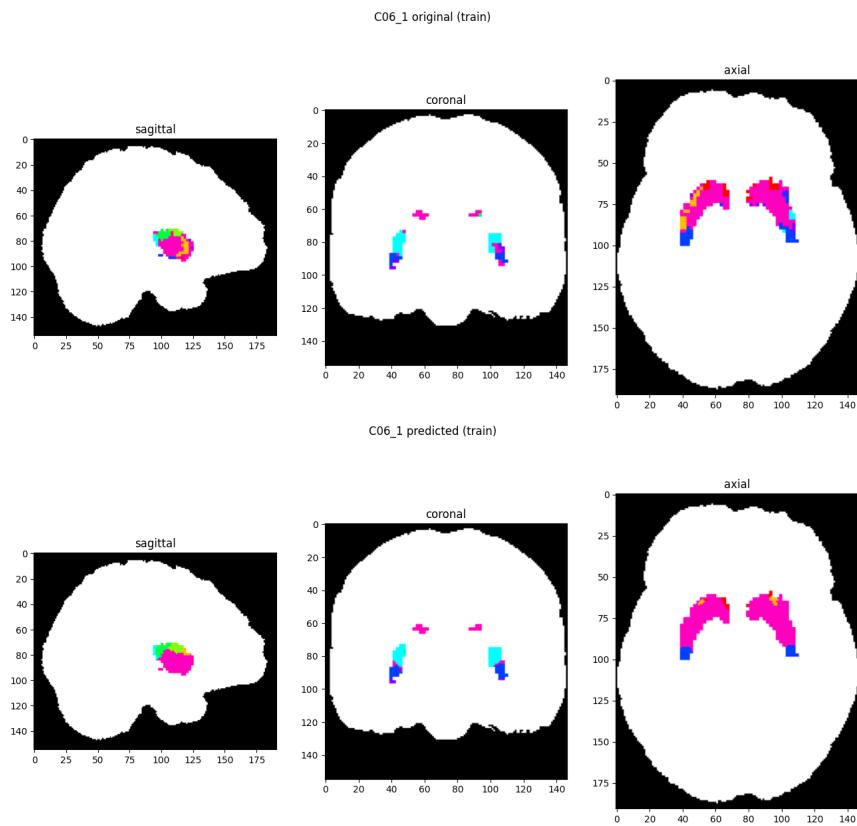


Figure 3.12: Train Predictions: Relative Connectivity

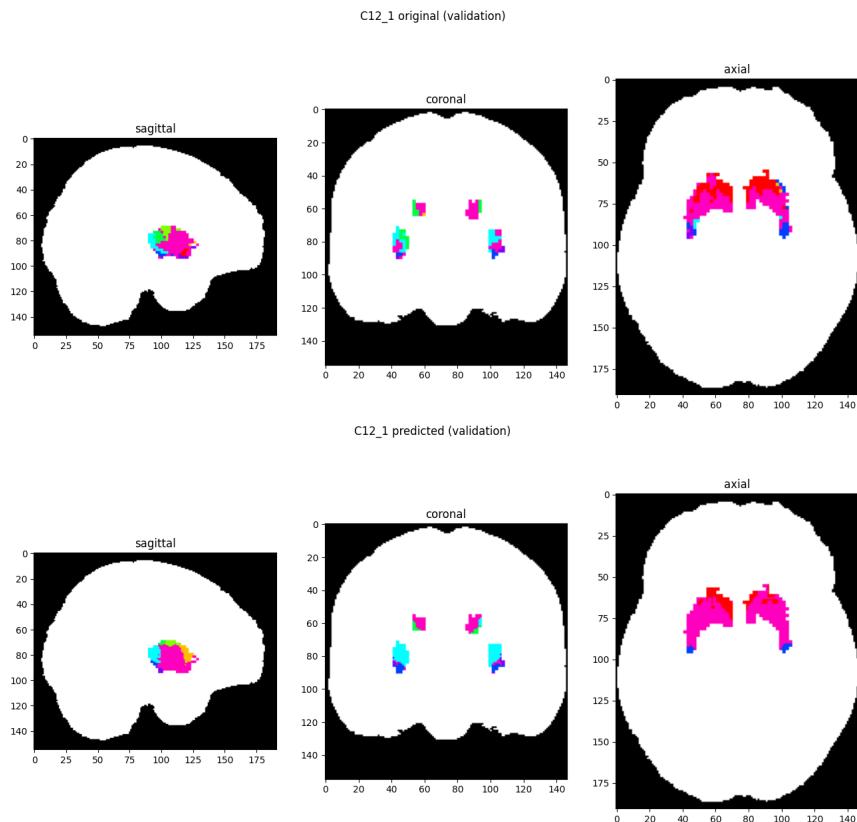


Figure 3.13: Validation Predictions: Relative Connectivity

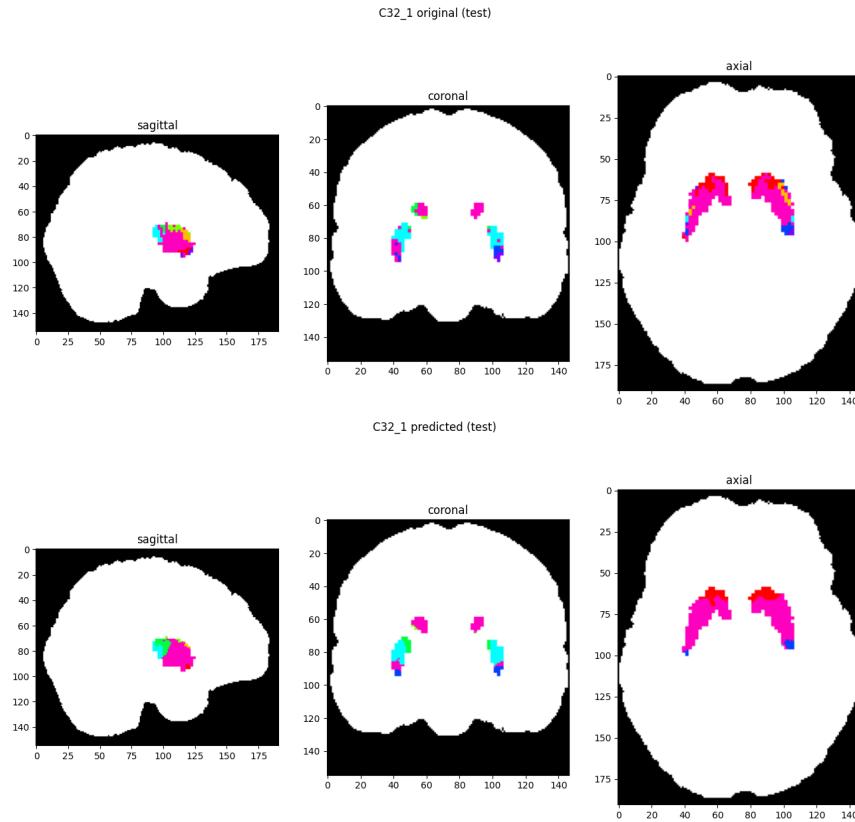


Figure 3.14: Test Predictions: Relative Connectivity

3.5.1 Exhaustive Sequential Backwards Feature Selection

Exhaustive sequential backwards feature selection, is about training the model iteratively by removing a single feature at a time, going through all features. After an iteration is done, the best performing model is chosen, and the corresponding feature is permanently removed before the next iteration, where the model goes through the remaining $n - 1$ features.

And the stopping criteria can be varied from setup to setup, but in this case the evaluation metric was chosen to be the validation raw accuracy, and a stopping point of where all runs in the iteration performed worse than the baseline (the model with all features included) by more than 2%.

Feature selection was only ran for the Relative Connectivity Segmentation problem, due to it being time consuming and computationally intensive. This was executed on a network cluster, where the cluster head would give out tasks (a task being a model to train with a set of features to exclude), and the workers returning the validation accuracy.

There were 2 logistical oversights that were not mitigated in time. The first mistake being that the execution of the feature selection was rushed due to the longevity of the task, and the rest of the hyperparameters (besides the excluded features) were chosen before finishing the basic experimentation. And the second being an indirect consequence to it being rushed, and it was the usage of a buggy code that was responsible for balancing the data, as the bug was found after already executing the feature selection.

Due to these mistakes, the selected features turned out not to be very efficient with the final hyperparameters. The feature selection was not executed again due to time and resource limitations. More information detailing this mistake and oversight can be found in Future Improvements 5.1

subsection.

Nevertheless the expected result of re-running the feature selection would be something similar to the current sub-optimal result. Where around 1/3rd of the features could be excluded before stopping and a maximum increase of 2% for the validation accuracy. The first features to be excluded was predominantly from the firstorder feature class, followed by the GLCM feature class.

Iteration	Excluded Feature	Accuracy
0.	BASELINE	59.0
1.	firstorder_MeanAbsoluteDeviation	59.6
2.	firstorder_Entropy	59.9
3.	firstorder_Energy	58.7
4.	glszm_SmallAreaLowGrayLevelEmphasis	58.5
5.	firstorder_90Percentile	59.3
6.	glcm_Autocorrelation	59.5
7.	firstorder_Mean	59.4
8.	firstorder_Maximum	59.1
9.	glszm_LowGrayLevelZoneEmphasis	58.9
10.	glcm_Imc1	60.0
11.	firstorder_Median	60.1
12.	glszm_ZoneVariance	59.3
13.	firstorder_TotalEnergy	58.9
14.	firstorder_10Percentile	60.4
15.	firstorder_Minimum	59.1
16.	glszm_SmallAreaHighGrayLevelEmphasis	59.9
17.	firstorder_InterquartileRange	60.0
18.	firstorder_Kurtosis	59.6
19.	firstorder_RobustMeanAbsoluteDeviation	59.3
20.	firstorder_RootMeanSquared	59.8
21.	glcm_LargeDependenceEmphasis	58.9
22.	firstorder_Variance	59.1
23.	glcm_DifferenceAverage	59.4
24.	firstorder_Uniformity	58.9
25.	firstorder_Skewness	59.0
26.	firstorder_Range	58.6
27.	glcm_JointAverage	58.7
28.	glcm_ClusterProminence	59.2
29.	glcm_ClusterTendency	58.5
30.	glcm_ClusterShade	60.2
31.	glcm_Correlation	60.0
32.	glcm_DifferenceEntropy	58.1
33.	glrlm_RunVariance	59.3
34.	glcm_JointEnergy	59.6
35.	glcm_Contrast	60.5
36.	glcm_Idm	59.3
37.	glcm_Imc2	58.7
38.	glcm_DifferenceVariance	59.2

39.	glcm_Idmn	58.0
40.	glcm_MCC	57.9
41.	glcm_JointEntropy	56.6

Table 3.5: Feature Selection

3.5.2 Streamline Regression

Another approach would be to try predicting the raw streamline images, and then process the predictions the same way as the relative connectivity labels were computed in the first place. This method could result in a more robust solution with predicting a the upstream data, before preprocessing.

Training 7 expert models for the 7 streamline images (to each cortical target), with the same hyperparameters as the best performing model from the relative connectivity segmentation, yielded underwhelming results. With the computed accuracies being 65.1/65.3/63.8, and terrible precision/recall on all labels (except the 'not connected'), which is explained by the predicted labels mainly being 'not connected'.

This approach was not explored further in depth due to limited time and resources.

Sustainability

4.1 Environmental Aspect

4.1.1 Development

The biggest environmental impact of this project was the required computational power to train the models. During the entire duration of the project there approximately a total of 5000 models were trained.

Out of which the most demanding part was the exhaustive sequential backward feature selection, which in itself required the training of 3000 models. The return value of this computation is close to none, as it did not increase the model performance. The only marginal benefit is the model being able to reach the same performance on 2/3rd of the original feature count. This is marginal because the omitted features belong to different feature groups, and the biggest overhead for most feature groups is calculating the voxel-based matrix volume; meaning if not all features are omitted from a feature group, it only decreases the computational requirements marginally. In conclusion, this was not worth computing.

The 3000 models trained for the features selection were trained on a wide variety of GPUs in a network cluster. These GPUs were H100, P40, RTX3060, RTX3090, RTX2080, multiple T4s (Google Collab's), multiple P100s (Kaggle's). Mining the cluster head's log, reveals that around half of these models were trained on the H100 and the other half on the rest. The H100 was consuming around 300 Watts of power while training with a speed of 5s/epoch and taking 30-60 epochs on average to stop. The H100 consumed between $1500 \cdot 0.3\text{kW} \cdot 30 \cdot 5\text{s} \div 60 \div 60 = 18.75\text{kWh}$ on the lower and $2 \cdot 18.75 = 37.5\text{kWh}$ on the upper end. The other half is harder to estimate due to the many different GPUs used in the cluster, but a crude estimate is $1500 \cdot 0.2\text{kW} \cdot 30 \cdot 8\text{s} \div 60 \div 60 = 20\text{kWh}$ on the lower and $2 \cdot 20 = 40\text{kWh}$ on the upper end. Rounding the sum of them up to 80kWh and adding an extra 30% overhead for running the rest of the components of the servers/PCs, results in an upper boundary of 104kWh consumed during the feature selection.

The rest of the 2000 models were trained on the P40 and RTX3060, with preferring the P40 and probably training 3/4th of them. Doing the same estimations for the P40 yields $1500 \cdot 0.2\text{kW} \cdot 30 \cdot 7\text{s} \div 60 \div 60 = 17.5\text{kWh}$ on the lower and $2 \cdot 17.5 = 35\text{kWh}$ on the upper end. And for the RTX3060 yields $500 \cdot 0.15\text{kW} \cdot 30 \cdot 7\text{s} \div 60 \div 60 = 4.375\text{kWh}$ on the lower and $2 \cdot 4.375 = 8.75\text{kWh}$ on the upper end. Rounding the sum of them up to 45kWh and adding an extra 30% overhead for running the rest of the components of the server, results in an upper boundary of 58.5kWh consumed during the the rest of the trainings.

Feature extraction was running on a CPU, with the entire server consuming 200 Watts during the process. The entire duration while feature extraction was running is around 2 weeks. The total power consumption of the feature extraction was $0.2\text{kW} \cdot 14\text{d} \cdot 24\text{h} = 67.2\text{kWh}$.

In summary this project had the following environmental impact, calculating with the average of 0.25€/kWh for the price of the electricity [8] and 0.2kgCO₂e/kWh for the emission [9]:

Part	Consumption	Price	Emission
feature selection	104kWh	26€	20.8kgCO ₂ e
rest of the experiments	58.5kWh	14.6€	11.7kgCO ₂ e
feature extraction	67.2kWh	16.8€	13.4kgCO ₂ e
Total (+10% overhead; rounded up)	260kWh	65€	52kgCO ₂ e

Table 4.1: Sustainability

For reference, this number is a bit less than the average Spanish household's monthly power consumption of 324kWh. [10]

4.1.2 Production

This project, whether implemented in a production, clinical or research setting, has the potential to significantly enhance sustainability in its respective fields. Acquiring DTI data typically requires up to 40 minutes, depending on the specific parameters, whereas obtaining a T1 image takes only about 8 minutes, making data acquisition a considerable time-saving step. Additionally, DTI data demands extensive and labor-intensive preprocessing. **Accelerating clinical or research workflows by a factor of 5** not only saves time but also increases the capacity for data collection or patient processing by the same factor.

Synthesizing a single record of relative connectivity or diffusion FA/MD takes only a few seconds on a GPU and the radiomic feature extraction takes around 20 minutes beforehand. This translates into the approximate energy demand of $200W \cdot (5s \div 60 + 20m) \div 60 = 67Wh$ of synthetizing a record. It can be argued that the time saved with the data acquisition is lost with the time it takes to extract the radiomic features, but neither the doctor (or MRI operator) and patient needs to be present during the computation of the feature extraction. Assuming an MRI machine has a consumption of 25-70kW [11] and calculating with an average of 50kW, adds up to the power requirement of $50kW \cdot 40m \div 60 = 33.3kWh$ per DTI record. Adding up the total energy requirements of T1 data acquisition and record synthesis is $50kW \cdot 8m \div 60 + 67Wh \div 1000 = 6.7kWh$ per record. This means that **this method could be 5 times more sustainable than the traditional approach**. And it would return the energy consumption of the development phase in the data acquisition of $260 \div (33.3 - 6.7) = 10$ records, this is naturally a minuscule number compared to how many MRI records are being done in the world.

4.2 Economic Aspect

4.2.1 Cost

The electricity cost during the development phase amounted to 65€. Approximately half of this cost was covered by colleagues and their networks, who contributed computational resources during the feature selection process. Additional support came from free services provided by Google and Kaggle.

Realistically, billing this project would consist of three components: the cost of labor, the price of the server used for development, and the electricity consumed. With the electricity cost already detailed, the server cost and labor are interconnected. The project could have been executed on a less powerful machine, though this would have required more time and more careful focus on

software design. Throughout the development, the primary limiting factors were RAM and disk storage. Significant effort was invested in developing efficient data structures and optimizing the memory footprint to address these constraints.

In conclusion, the minimum hardware which this project could be feasibly executed on is 32GBs of RAM, Intel Core i5-12600K (or similar), Nvidia RTX3060 (or similar), and 256GBs of storage. The current cost (new parts ordered from Amazon) of such server (with 'cheap' consumer grade hardware) would be around 300€ for the GPU, 150€ for the CPU, 50€ for the RAM, 25€ for the SSD, and 150€ for the motherboard, adding up to 675€. The main limiting factors in this configuration are the RAM, which was utilized up to 128GBs during development; and disk storage, which was utilized up to 300GBs. The extra RAM and storage would cost an additional 200€ and 25€ totalling at a new sum of 900€.

The development of this project required approximately 500 hours of work. Calculating with the average hourly rate of a software engineer in Spain [12] of 19€, leads to a total labor cost of $19\text{€} \cdot 500\text{h} = 9500\text{€}$. Therefore, the total billing for the project would be:

Part	Price
Cumulative Salary	9500€
Server Components	900€
Electricity	65€
Total	10465€

Table 4.2: Billing

4.2.2 Return

This project by all means, is a proof of concept, checking the viability of this approach. More information on the exact conclusions are in chapter 5, but the potential return value of this project is huge.

As stated in the previous subsection 4.1.2, this project has the potential to accelerate the clinical and/or research workflows by a factor of 5. The economic implications are applicable for both clinical and research workflows. In the clinical workflow, **this could increase the number of patients processed by a factor of 5**, with minimal extra cost. Especially compared to the alternative of buying 4 additional MRI machines, and hiring operators for them. In a research workflow it probably has less of an economic impact, but more of a logistical impact. As it could simplify the data acquisition, thus **allowing researchers to have up to 5 times more data acquired** with minimal extra resources invested. But even more importantly it opens the door of processing past anatomical MRI records, as they are much more common than dMRI records, virtually **increasing the available data for researchers by several magnitudes** (depending on the exact application of course).

4.3 Social Aspect

4.3.1 Development and Collaborations

The development of the project mainly required collaboration from Estela Camara Mancha (external thesis supervisor and project originator) and Alfredo Vellido Alcacena (internal thesis su-

pervisor). Estela being the project originator, this project demanded her time the most, in the neighborhood of 40-50 hours **TODO**. Alfredo mostly helping with the formalities and miscellaneous nuances, had a demand of 10-15 hours **TODO** of his time. Additionally Estella's colleague Vasiliki Bikou also spent a good chunk of her time on catching me up to speed with the contextual information of the project. And lastly my good friends Botond Lovasz, Andris Gyori and Daniel Csepregi-Horvath donated computational power, and required a few hours of their time to set up the worker nodes on their hardware. I especially want to thank Botond for using his connections and giving me access to the H100 state of the art tensor core GPU.

4.3.2 Inclusivity

The most obvious consideration from age, gender, sex, and cultural diversity; is sex, as males and females have slightly different brain structures [13]. For the used control records, there are 17 male and 15 female records, which is a relatively even split; and for the patient records, there are 25 male and 13 female records, which is a bit unbalanced. This theoretically could negatively impact the under represented group if there are truly any differences between the two sexes that matter from the model's perspective. However in the case of the project due to the very limited number of records, it was not an option to discard 1/3rd of the patients, to have the sexes completely balanced. And it would also need further experimentation to determine if this truly impacts this project or not.

This potential difference was overlooked during most of the project's lifetime, but it should have been included as another constant ratio to be kept during the train/validation/test splitting, which was covered in subsection 2.2.1. Due to limited time and resources, the models will not be re-trained with this new rule in mind, and this is something that will be included in the Future Improvements 5.1 subsection.

The rest of the considerations are less pressing and hard to take into account, like the brain structure of different ethnicities is an under researched area, plus this dataset does not contain any information regarding this aspect. The only other consideration which can be related is age, but this is partially accounted for as part of the constant symptomatic/asymptomatic ratio. Because being symptomatic is closely related to the CAP score (covered in subsection 2.1.8), which is directly related to age. Thus symptomatic/asymptomatic is indirectly related to age.

4.4 Risks

Environmental, Economic and Social risks are not really applicable to this project, as it is a foundational proof of concept research project.

Conclusions

This foundational project has huge potential thanks to the possible return and applications of the used approach. Due time and resource limitations, the experiments only focused on the Basal Ganglia, so by no means any of the performance metrics prove that a generalized solution is viable. However the FA and MD predictions are really promising, and would be interesting to see how would it perform on the entire brain, and not just the ROI. On the other hand the Relative Connectivity predictions are not even nearly precise enough to call them usable, but admittedly the labels themselves will inherently contain quite a bit of noise, due to the very sensitive process extracting the labels with tractography.

This project only delved into a tiny fraction of the endless sea of possible preprocessing approaches, and it only experimented with some very basic model architectures. Nevertheless, the viability of this approach and hypothesis is not disproven, but it definitely needs some new ideas implemented to make it work.

5.1 Future Improvements

As mentioned in subsection 3.5.1, there was a serious oversight during the execution of the exhaustive sequential backwards feature selection. At the time it was believed that the model was performing better on balanced data, due to a bug in the code. Thus, all of these models are performing much worse than the models during the experimentation. The bug was fixed, but due to limited time and resources the feature selection was not executed again.

Further investigating the sex imbalance issue discovered in subsection 4.3.2, reveals that with the used seed for the experimentation is not terribly imbalanced for the most part, but there are definitely room for improvement by enforcing a constant ratio. The male/female ratio for the control record splits are 0.49/0.67/0.67 (train/validation/test), and for the patient records are 0.62/0.5/1.

Besides these known mistakes, there are some issues that are hard to quantify. For example due to the different imaging of the anatomical T1 record and the DTI record, even after a perfect affine registration they can have tiny misalignments, which could only be resolved by non-affine warping. It is hard to even estimate the impact of these tiny misalignments, but it is probably not negligible.

Additionally, one aspect was not investigated thoroughly during this project and that being the kernel size, binning parameters, and feature class relationships, during the voxel based feature extraction. For the sake of simplicity, only a unified binning method was used for all kernel sizes and all feature classes. But logically the binning parameters could be optimized for each kernel size and feature class. For example logically the GLRLM feature class could yield fundamentally different features even just by adjusting the bin size a tiny bit at large kernel sizes.

Selecting the most efficient kernel size combinations were not investigated thoroughly either. Due to limited resources, the maximum kernel size used was 21, but there were no reason to stop here (besides to be able to finish this project in time).

Ultimately, even with the 'few' simple aspects that were taken into consideration during this

project, the hyperparameter space is gigantic and some compromises were had to be made. After investing over 500 hours into this project, the educated guesses for these potential improvements are:

- The feature selection should have a marginal improvement on the model performance.
- The gender imbalance should not have a measurable impact on the model performance on this scale (with less than 70 available records in total).
- The non-affine registration should have at least a marginal improvement on the model performance, with questionable reward/effort ratio, as this would be a relatively huge added effort in the preprocessing pipeline, and a lots of extra added points of failure.
- The binning parameter optimization for the kernel sizes and feature classes would probably have a big impact on the model performance.

And lastly, all experiments should be repeated at least a few times, with different splits and seeds. And the experiments should be evaluated with the means and medians of the model performance indicators, across different runs.

5.2 Project Future

Ultimately there are two paths to continue this project. First, the FA and MD experiments could be expanded to the entire brain and a generalized model should be developed. The other being the refinement of the Relative Connectivity experiments, as it definitely need improvements and new ideas.

Sources of Information

- [1] José L Lanciego, Natasha Luquin, and José Obeso. “Functional neuroanatomy of the basal ganglia”. In: *Cold Spring Harbor perspectives in medicine* (2012). URL: <https://doi.org/10.1101/cshperspect.a009621>.
- [2] Olivia C Matz and Muhammad Spoter. “The Effect of Huntington’s Disease on the Basal Nuclei”. In: *Cureus* (2022). URL: <https://doi.org/10.7759/cureus.24473>.
- [3] Hyungyou Park et al. “Aberrant cortico-striatal white matter connectivity and associated subregional microstructure of the striatum in obsessive-compulsive disorder”. In: *Molecular Psychiatry* (2022). URL: <https://doi.org/10.1038/s41380-022-01588-6>.
- [4] Marius E Mayerhoefer et al. “Introduction to radiomics”. In: *Journal of Nuclear Medicine* 61.4 (2020), pp. 488–495. URL: <https://jnm.snmjournals.org/content/jnumed/61/4/488.full.pdf>.
- [5] Loïc Duron et al. “Gray-level discretization impacts reproducible MRI radiomics texture features”. In: *PLoS One* (2019). URL: <https://doi.org/10.1371/journal.pone.0213459>.
- [6] *Special Characteristics of HD Data*. 2022. URL: <https://enroll-hd.org/for-researchers/analyzing-data/special-characteristics-of-hd-data-2>.
- [7] Dylan Trundell et al. “Defining Clinically Meaningful Change on the Composite Unified Huntington’s Disease Rating Scale”. In: *Neurology* 92.15_supplement (2019), P1.8–043. URL: https://doi.org/10.1212/WNL.92.15_supplement.P1.8-043.
- [8] *Spain - Household electricity prices*. URL: <https://countryeconomy.com/energy-and-environment/electricity-price-household/spain> (visited on 12/28/2024).
- [9] B.W. Ang and Bin Su. “Carbon emission intensity in electricity production: A global analysis”. In: *Energy Policy* 94 (2016), pp. 56–63. ISSN: 0301-4215. URL: <https://doi.org/10.1016/j.enpol.2016.03.038>.
- [10] *Electricity Consumption per Dwelling*. URL: <https://www.odyssee-mure.eu/publications/efficiency-by-sector/households/electricity-consumption-dwelling.html> (visited on 12/28/2024).
- [11] *MRI and Sustainability*. URL: <https://www.siemens-healthineers.com/perspectives/MRI-reducing-energy-consumption> (visited on 12/28/2024).
- [12] *Average Software Engineering Salaries by Country*. URL: <https://codesubmit.io/blog/software-engineer-salary-by-country> (visited on 12/28/2024).
- [13] Daphna Joel. “Male or Female? Brains are Intersex”. In: *Frontiers in Integrative Neuroscience* 5 (2011). ISSN: 1662-5145. URL: <https://doi.org/10.3389/fnint.2011.00057>.

Software Design

The software design was mainly guided by the provided format of the raw data. Which was from a few different sources, as different people were working with different data at the Hospital, and they did not have a unified collection. Thus the following documentation of the software design will be structured going from the raw data, to the preprocessed data, and then the model itself.

A.1 Raw Data

The raw data is scattered amongst many files, are registered in different spaces, does not have masks applied, along with other challenges detailed in Section 2.1.

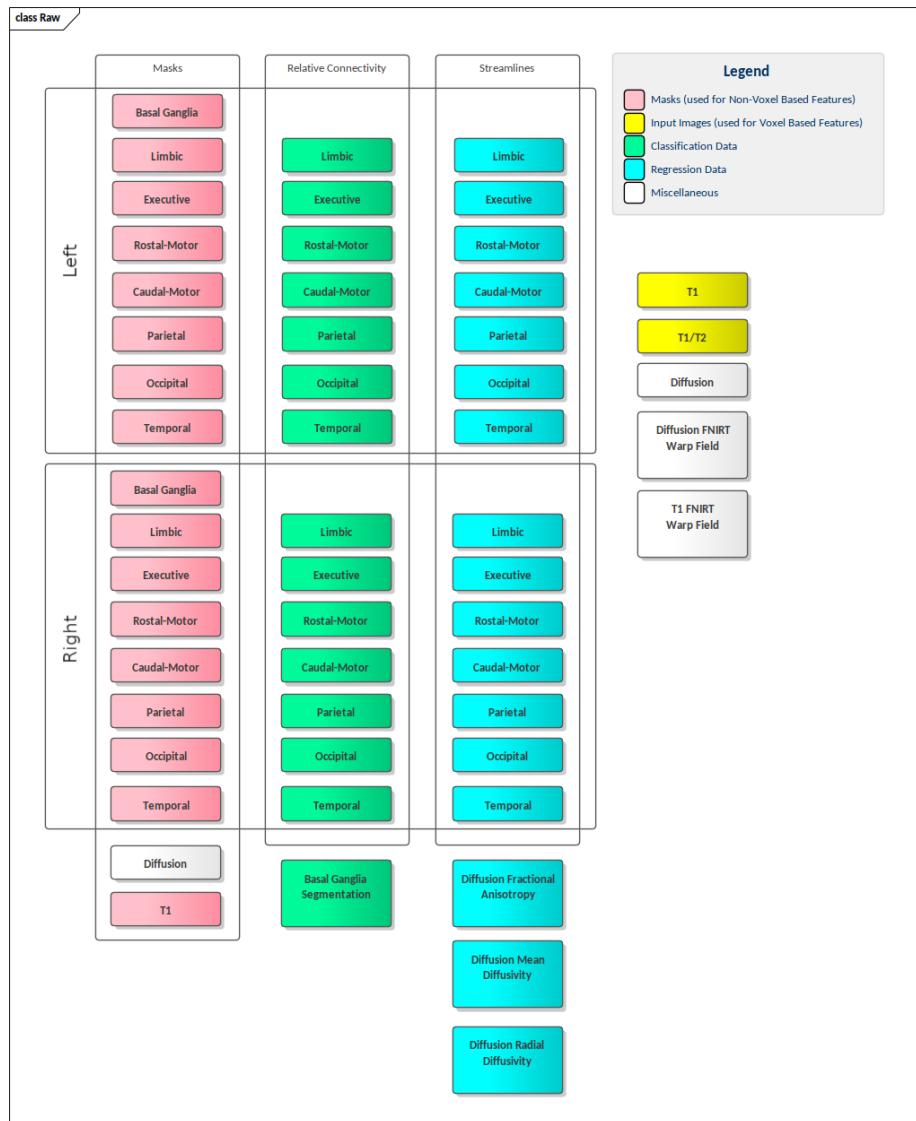


Figure A.1: Files: Raw

The following class diagram documentations are detailed abstractions of the actual source code, as unified modeling language (UML) is not completely Python compliant. Before moving on, the following simple data types are used in the UML diagrams:

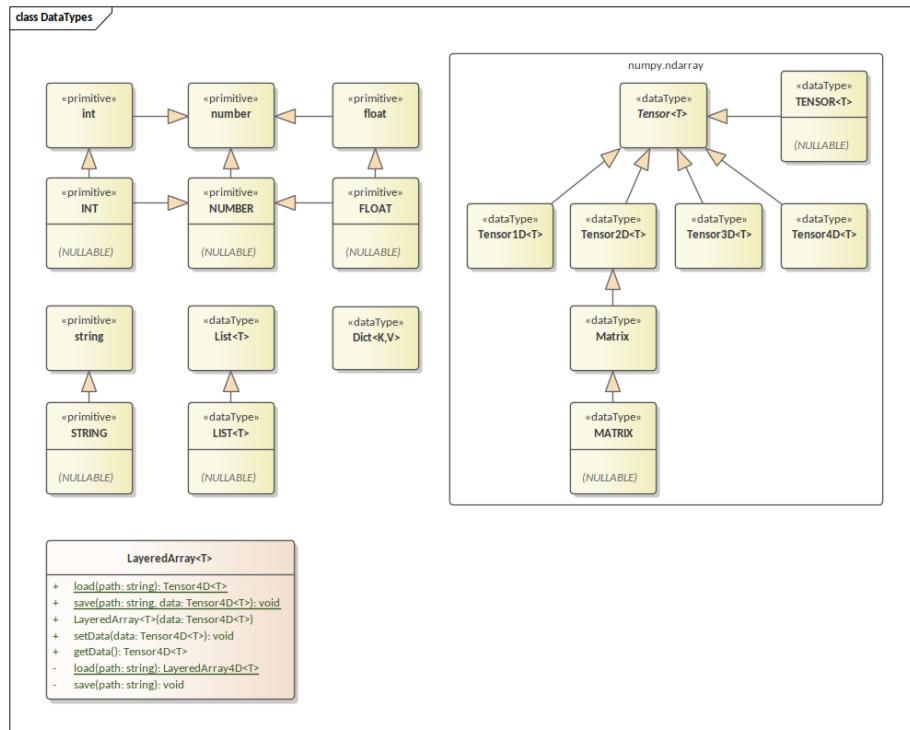


Figure A.2: Class Diagram: Data Types

There is one odd one out in this diagram, that being the `LayeredArray`, which is more than a simple data type or primitive. This class is a simple space efficient data structure of a 4D tensor.

Storing some of the data provided can be highly inefficient by storing the raw tensor. Ultimately the design aims to collect the scattered data from multiple files, into clean logical groups, for example all the cortical target masks into a single 4D tensor where the 4th dimension is reserved for the multiple targets. But this poses a greatly inefficient way of storing the data, as each target region only takes up a small portion of the entire space. The idea behind the `LayeredArray` class, it stores the 4D tensor as a list of 3D tensors, each cropped down to their effective space where there are non-zero voxels. This also requires an additional list of 3D vectors to store each layer's origin, and where to paste it in the original space.

This solution offers a very efficient way of storing data for this use case, as the raw storage solution for the cortical targets were around 50MBs, and this cuts it down to 3MBs. This is even more drastic for the relative connectivity, which is more than a simple boolean mask and proportionally takes up even less space of the entire brain; in which's case it was cut down from 110MBs to 0.5MBs, reducing the disk requirement by several magnitudes.

The original NIfTI format also does a very good job at storing data efficiently, but this solution provides control over the way of storing the data on a much lower level. This is beneficial as the data are stored in numpy format, making it easier to ignore certain data type safety checks, data type conversions, and leaving behind the NIfTI format's additional complexity of the orientation, transformation, and many more nuances that are part of the NIfTI header.

This has the undoubted drawback of not being able to use FMRIB software library (FSL) tools natively on our datatypes, such as fsleyes for simply viewing a record. But thanks to the

opensource nature of the FSL suite, with a few additional lines of code, support can be added for our datatypes (included in Appendix E).

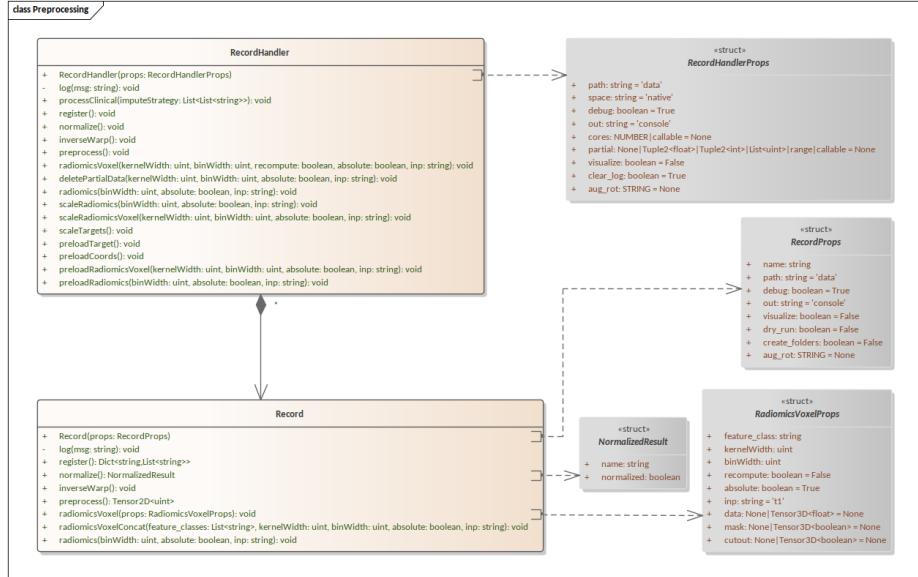


Figure A.3: Class Diagram: Preprocessing

The class diagram above contains the two classes responsible for the preprocessing of the data. The RecordHandler being the controller itself that handles the high level operations on a collection of records. And the records handling the low level computations on the data itself.

A.2 Common Functions

There are set of static common functions, which give the low level backbone of the entire project. They are grouped into two categories, util and visual. Where the prior one contains everything from simple data type castings, external FSL library system calls, to computing radiomics, and more. And the latter one is a collection of functions for visualizing data.

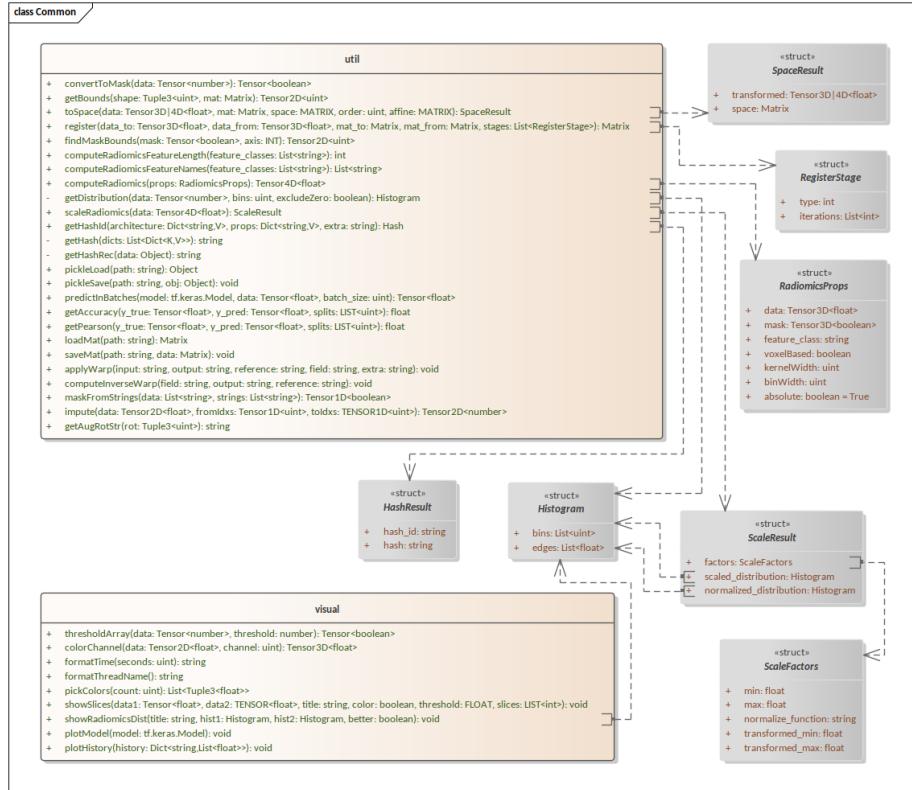


Figure A.4: Class Diagram: Common

A.3 Preprocessed Data

Preprocessing the data are composed from the next few high level operations done by the Record-Handler:

1. Register some records (not needed for most of them)
2. Compute all normalized records
3. Compute the inverse FNIRT warp field
4. Convert native records into our numpy format (apply the affine transformations, merge different sources, etc.)
5. Compute scaling factors across all native records
6. Preload native records
7. Convert normalized records into our numpy format
8. Compute scaling factors across all normalized records
9. Preload normalized records
10. Construct normalized coordinate maps
11. Warp normalized coordinate maps into native space

12. Scale and preload coordinate maps

13. Impute clinical data

After preprocessing, the next set of logical groupings and files are left, split into native and normalized records:



Figure A.5: Files: Preprocessed

A.4 Data Generator

As this project focuses on only the voxels inside the Basal Ganglia, means the model on a datapoint level is never going to operate outside of that volume.

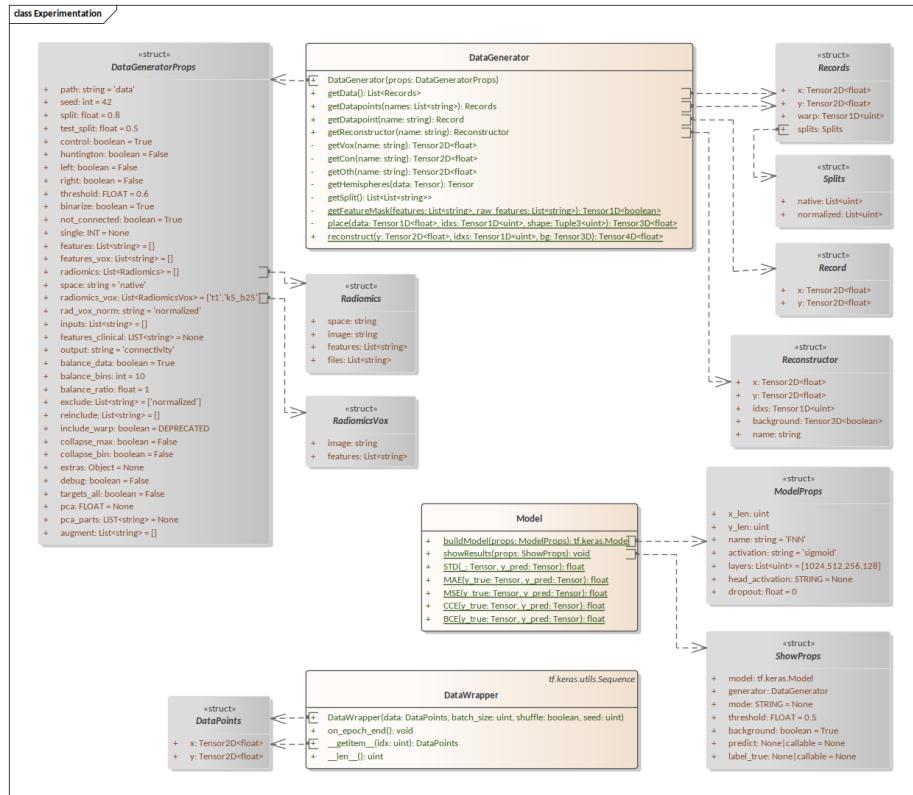


Figure A.6: Class Diagram: Experimentation

The DataGenerator class is responsible for feeding the Model with data, according to the specifications of the experiments. In order to make this more efficient and fast, loading the entire volume of records from disk can be avoided by 'preloading' the ROI from each record. This simply put means that all logical grouping and files specified in Figure A.5 are preloaded into 1D and 2D tensors, where the first dimension is the flattened voxels of the ROI and the second is whatever layers were stored originally in the last dimension of the spatial data (for example the streamlines of the different target regions).

Computing and persisting the preloaded data vastly speeds up the DataGenerator. Although this has the added consequence that the preloaded data yields a pair of tensors instead of a single volume, as there is differentiation between the left and right hemisphere datapoints. And the voxel count in the two hemispheres are not symmetrical, thus the hemisphere differentiation can not be stored as an extra dimension. A design decision was made to store the left and right hemisphere data in separate files.

The properties of the DataGenerator provided in its constructor defines the type and format of the generated data.

Value	Type	Description
path	string	path of the data
seed	int	random seed for the train/val/test splits
split	float	train/all ratio
test_split	float	test/(test+validation) ratio
control	boolean	include control records
huntington	boolean	include huntington records
left	boolean	include left hemisphere datapoints

right	boolean	include right hemisphere datapoints
threshold	FLOAT	if not None it thresholds the labels by setting the values under the provided threshold to zero; if 0 it re-one-hot encodes the labels (sets the maximum label to 1 and the rest to 0)
binarize	boolean	if thresholded and True, it also sets the labels above the threshold to 1 (and the rest to 0)
not_connected	boolean	if thresholded and True, it appends an additional 'not connected' label, which complements the sum of the labels per voxel to 1
single	INT	if not None it only returns the label with the provided index
features	List<string>	used non-voxel based radiomics features (emptylist means all)
features_vox	List<string>	used voxel based radiomics features (emptylist means all)
radiomics: List<Radiomics>	space: string	space of non-voxel based radiomic features (native/normalized)
	image: string	input image of non-voxel based radiomic features (t1/t1t2)
	features: List<string>	list of binning settings of the features (b25, b50, b10r... etc.)
	files: List<string>	list of region mask settings of the features (targets/roi/brain)
space	string	space of voxel based radiomic features (native/normalized)
radiomics_vox: List< RadiomicsVox >	image: string	input image of voxel based radiomic features (t1/t1t2)
	features: List<string>	list of binning and kernel settings of the features (k5_b25, k9_b50, k7_b10r... etc.)
rad_vox_norm	string	use normalization, or only min-max scaling on the voxel based radiomic features (norm/scale)
inps	List<string>	additional voxel inputs (t1/t1t2/diffusion/diffusion_fa/diffusion_md/diffusion_rd)
features_clin	List<string>	additional clinical data inputs (empty array means all)
outp	string	output (connectivity/streamline/basal_seg/diffusion_fa/diffusion_md/diffusion_rd)
balance_data	boolean	enables data balancing
balance_bins	int	number of bins used for continuous data when balancing
balance_ratio	float	ratio of the resampling of the difference between each bin and the max bin when balancing (where 0 is unbalanced and 1 is perfectly balanced)
exclude	List<string>	can manually add missing groups of records to exclude (t1t2/normalized/basal_seg/diffusion_fa)

reinclude	List<string>	can manually re-include (append) missing groups of records to the train split (t1t2/normalized/basal_seg/diffusion_fa)
debug	boolean	only returns 1/1/1 records for train/val/test when True
augment	List<string>	list of record suffixes, used to include augmented records in the training split; for example the suffix '_5_0_0' is used for the rotation augmentation of 5, 0 and 0 degrees on the X, Y and Z axis (naturally these records needs to be computed before using here)

Table A.1: Data Generator Properties

Software Implementation

Detailing the entire implementation of the project is beyond the scope of this report. But some solutions contain non self explanatory ideas and are worthy of detailed explanations. For more details please refer to the source code.

B.1 Interpolation

One crucial detail of applying transformations to the records, is the type of interpolation used. For some images like T1 the standard trilinear interpolation makes sense, as the values inbetween voxels should be continuous. However for some other images like binary masks, the interpolation should be nearest neighbour, as the edge of the mask cant be interpreted as a fraction. And most importantly the relative connectivity and streamline images should also be interpolated with nearest neighbour, preserving their characteristics like their summed total being 1 per voxel for the relative connectivity. In theory the streamline image could be computed with trilinear interpolation, after which the relative connectivity could be computed from the interpolated streamlines. But this would degrade the streamline images due to their unbalanced nature of having a few voxels with high number of streamlines relative to the entire voxel space, effectively eroding the boundaries of the high intensity volumes.

In practice this meant using 0th and 1st order spline interpolations, which are numerically the same as nearest neighbour and trilinear interpolations.

B.2 Multithreading

As most of the preprocessing could not be easily offloaded to the GPU (without re-implementing major libraries such as 'pyradiomics' with GPU support), it was crucial to implement multithreading to save time.

The most straight forward way of doing so was to split the load on a record level, and process multiple records parallel. Testing revealed an interesting property of the feature selection, which is the RAM IO operation bottleneck of the process. Extracting voxel based features happen in a 'voxel batch', which seemed to only increase the RAM usage of the process beyond a certain point, but not the speed. After days of tweaking and running tests, it was concluded that the GLCM feature class is RAM IO operation heavy. Meaning that increasing voxel batch size or the number of threads (computing different records in parallel) will not increase the performance past a certain point, as the threads are waiting on RAM IO operations regardless.

The optimal settings with the used Intel Core i5-12600K CPU and 128GBs of 3200MT/s RAM, is to use 7 threads, out of which 5 are dedicated for computing the GLCM feature class and the rest 2 are for computing other feature classes, with a voxel batch size of 1000.

Trying to increase performance this point by increasing the voxel batch size, or the number of threads only resulted in the same, or worse overall performance due to the computational overhead of splitting and merging the work between more threads, and using a lot more RAM.

B.3 Warp Pre-Computing

As mentioned in subsection 2.2.2, computing the evaluation metrics both in native and normalized space can be done by reconstructing the spatial records from the predicted datapoints and warping them to native/normalized space, and then re-extracting the datapoints from the warped records. However this would be computationally very expensive to do so every time when calculating the evaluation metrics.

The solution is to pre-compute index arrays, by assigning unique indexes to the voxels of the records, warping the records, and mapping the ROI's flattened voxels' unique indexes between the (warped and non-warped) record pairs into an index array which can be used to quickly and efficiently convert the voxels of the ROI between native and normalized spaces.

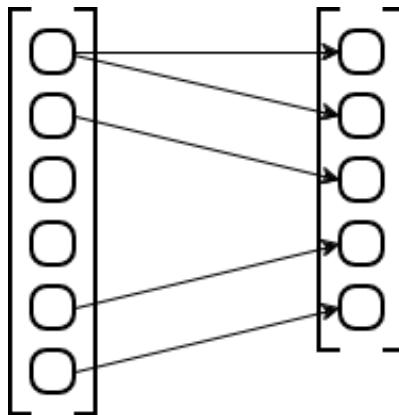


Figure B.1: Index Array Mapping

B.4 Data Normalization

As noted in subsection 2.1.6, determining which voxel-based radiomic features require log scaling for distribution normalization is handled programmatically. This process involves loading each feature across all records and analyzing their distribution both before and after log scaling.

Initially, the feature is min-max scaled, and the distribution of non-zero voxels is calculated using 100 bins (bin size 0.01). Excluding zero values helped to mask the background, which can distort the distribution. However, this exclusion should have been done with the provided brain masks, which would have been more appropriate since zeros could naturally occur in the features. This oversight, while suboptimal, does not significantly affect the results, as non-background zero counts are negligible. The only consequence is a potential false negative artifacts (failing to apply normalization where needed), which does not degrade the raw data by inappropriately normalizing features.

Next, the logarithm is applied to the feature prior to min-max scaling, and the distribution is recalculated using the same 100 bin approach. The criteria for selecting normalized features are based on observing an increased standard deviation and a reduced count of voxels in the largest bin after normalization. This ensures that features with 'flattened' and 'stretched' distributions are selected for normalization.

Additional Figures

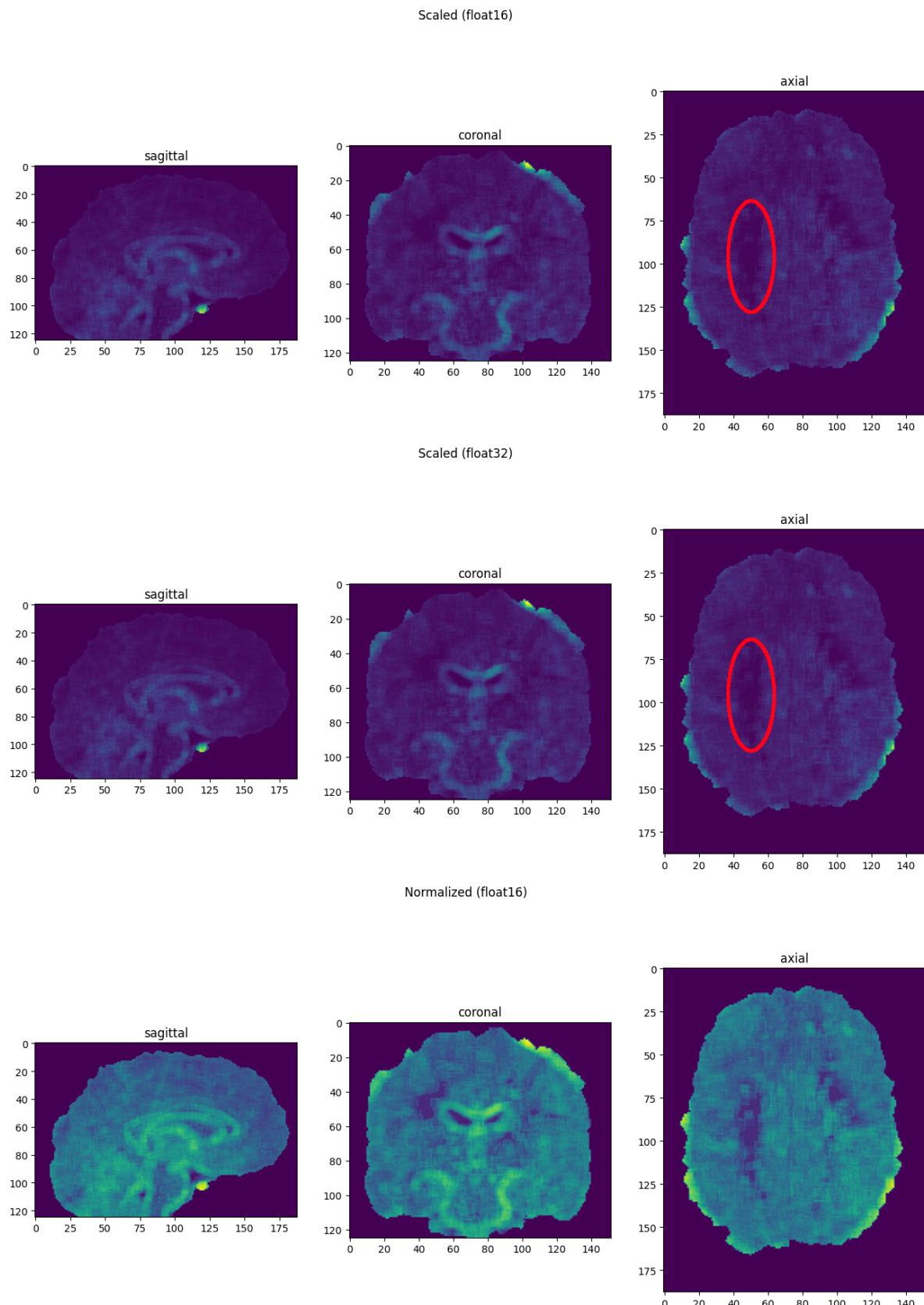


Figure C.1: Slice: GLDM Small Dependence High Gray Level Emphasis

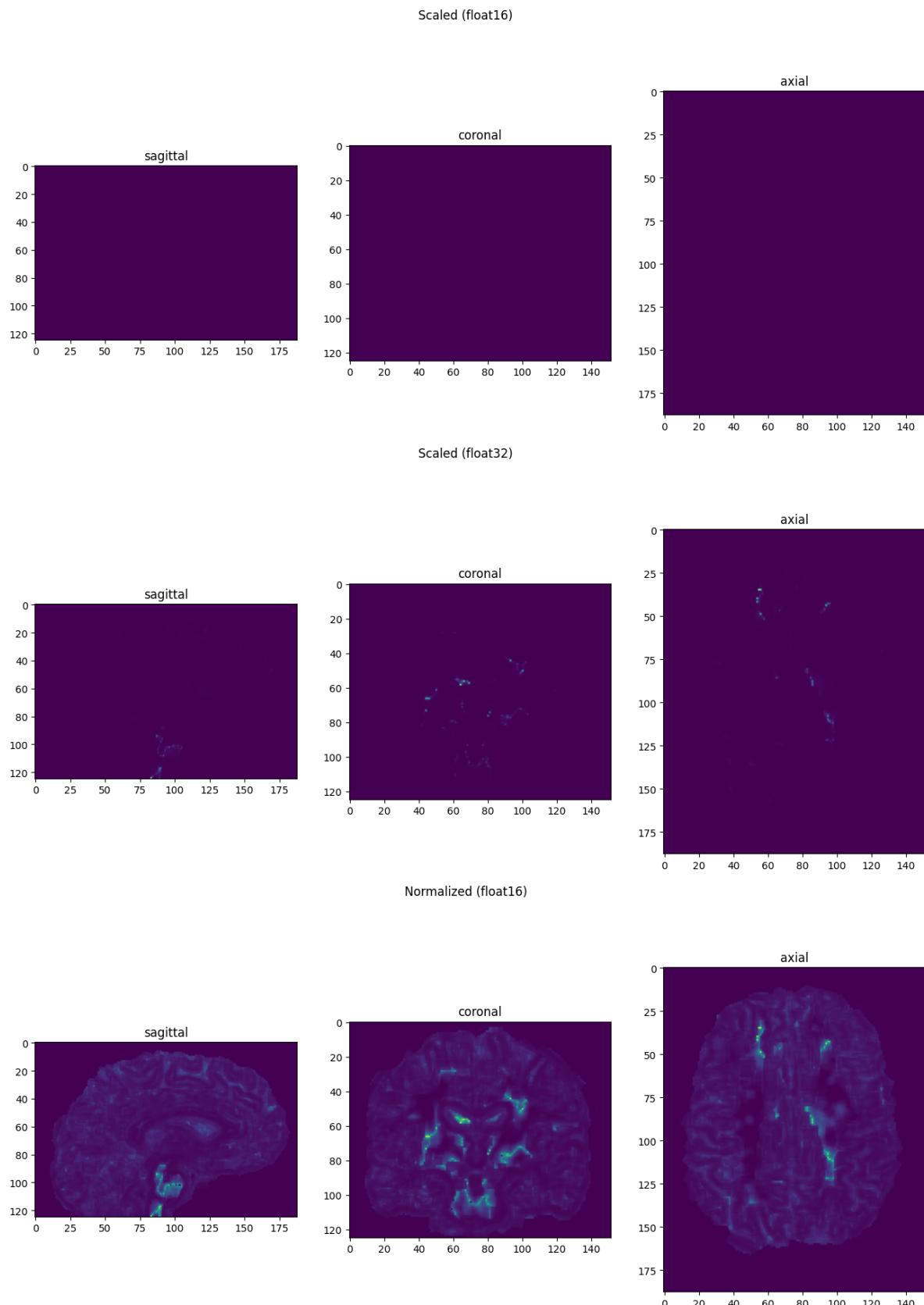


Figure C.2: Slice: NGTDM Busyness

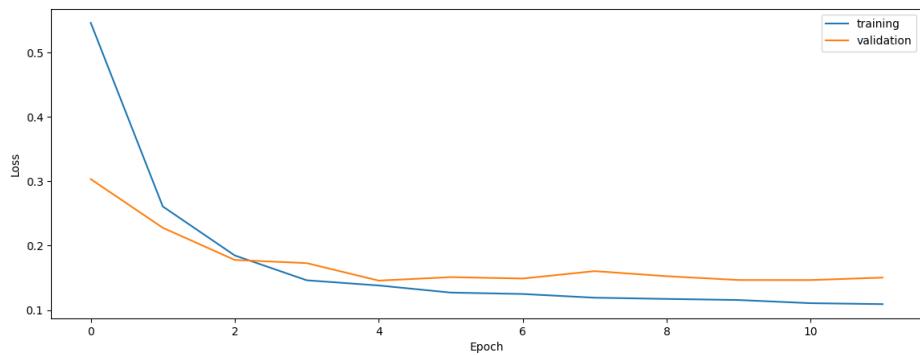


Figure C.3: Training Curve: Subcortical

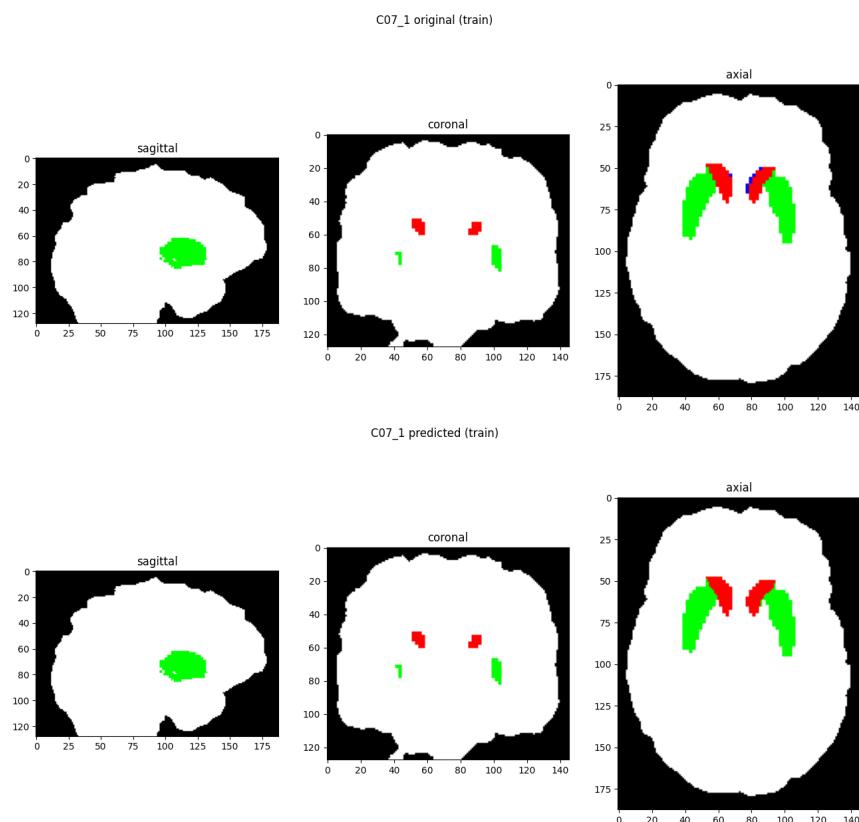


Figure C.4: Train Predictions: Subcortical

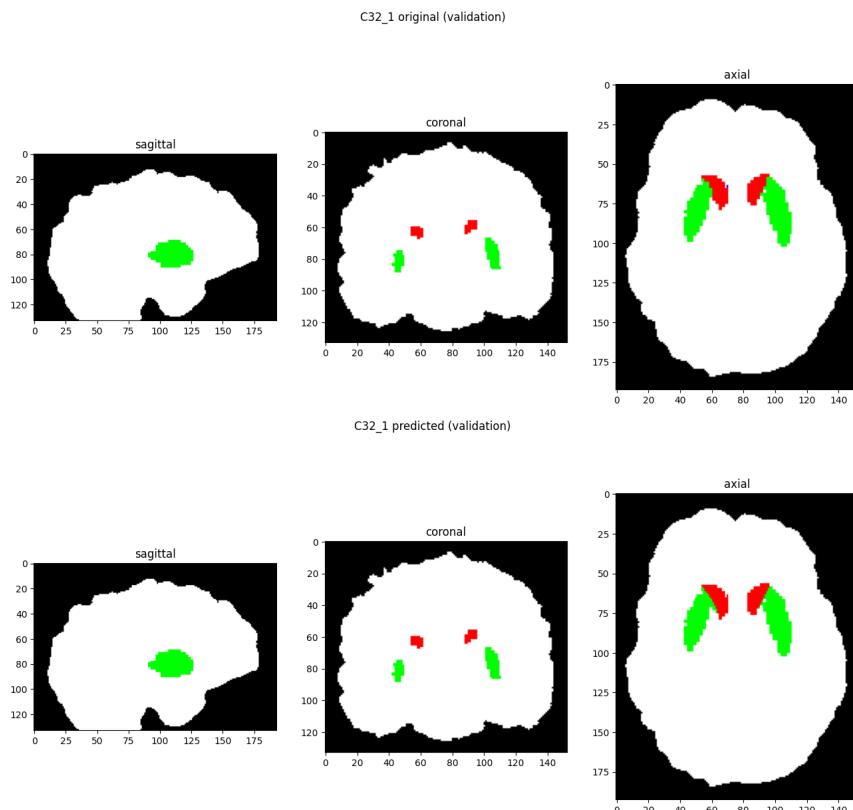


Figure C.5: Validation Predictions: Subcortical

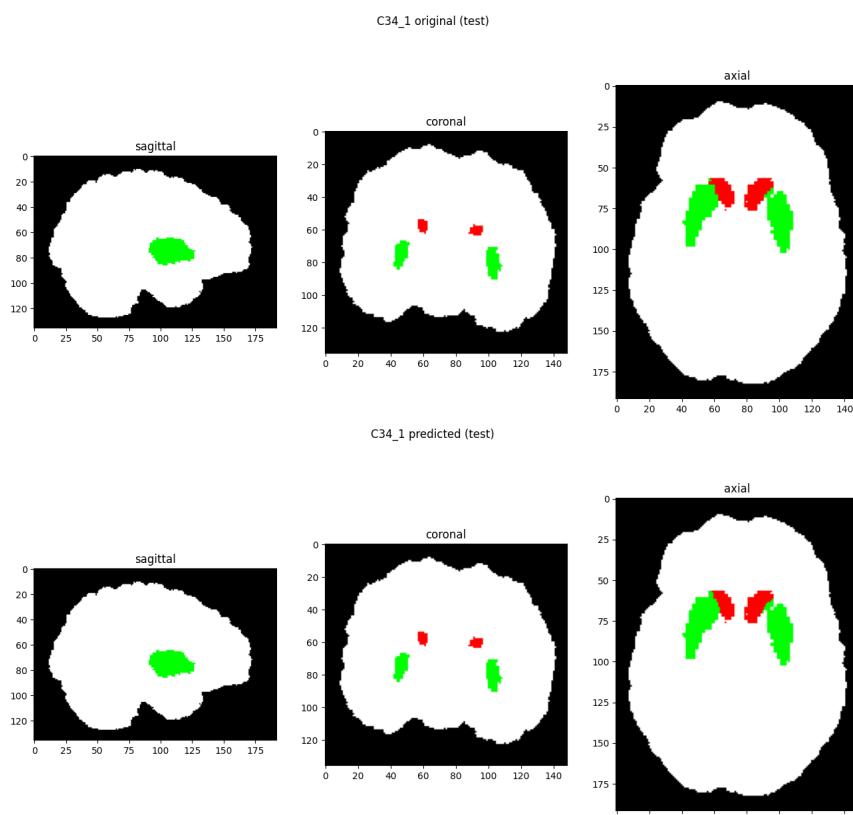


Figure C.6: Test Predictions: Subcortical

Additional Tables

First Order	GLCM	GLSZM
Energy	Autocorrelation	SmallAreaEmphasis
TotalEnergy	JointAverage	LargeAreaEmphasis
Entropy	ClusterProminence	GrayLevelNonUniformity
Minimum	ClusterShade	GrayLevelNonUniformityNormalized
10Percentile	ClusterTendency	SizeZoneNonUniformity
90Percentile	Contrast	SizeZoneNonUniformityNormalized
Maximum	Correlation	ZonePercentage
Mean	DifferenceAverage	GrayLevelVariance
Median	DifferenceEntropy	ZoneVariance
InterquartileRange	DifferenceVariance	ZoneEntropy
Range	JointEnergy	LowGrayLevelZoneEmphasis
MeanAbsoluteDeviation	JointEntropy	HighGrayLevelZoneEmphasis
RobustMeanAbsoluteDeviation	Imc1	SmallAreaLowGrayLevelEmphasis
RootMeanSquared	Imc2	SmallAreaHighGrayLevelEmphasis
Skewness	Idm	LargeAreaLowGrayLevelEmphasis
Kurtosis	MCC	LargeAreaHighGrayLevelEmphasis
Variance	Idmn	
Uniformity	Id	
	Idn	
	InverseVariance	
	MaximumProbability	
	SumEntropy	
	SumSquares	

GLRLM	NGTDM	GLDM
ShortRunEmphasis	Coarseness	SmallDependenceEmphasis
LongRunEmphasis	Contrast	LargeDependenceEmphasis
GrayLevelNonUniformity	Busyness	GrayLevelNonUniformity
GrayLevelNonUniformityNormalized	Complexity	DependenceNonUniformity
RunLengthNonUniformity	Strength	DependenceNonUniformityNormalized
RunLengthNonUniformityNormalized		GrayLevelVariance
RunPercentage		DependenceVariance
GrayLevelVariance		DependenceEntropy
RunVariance		LowGrayLevelEmphasis
RunEntropy		HighGrayLevelEmphasis
LowGrayLevelRunEmphasis		SmallDependenceLowGrayLevelEmphasis
HighGrayLevelRunEmphasis		SmallDependenceHighGrayLevelEmphasis
ShortRunLowGrayLevelEmphasis		LargeDependenceLowGrayLevelEmphasis
ShortRunHighGrayLevelEmphasis		LargeDependenceHighGrayLevelEmphasis
LongRunLowGrayLevelEmphasis		
LongRunHighGrayLevelEmphasis		

Table D.1: Voxel Based Radiomic Features

3D Shape
MeshVolume
VoxelVolume
SurfaceArea
SurfaceVolumeRatio
Sphericity
Maximum3DDiameter
Maximum2DDiameterSlice
Maximum2DDiameterColumn
Maximum2DDiameterRow
MajorAxisLength
MinorAxisLength
LeastAxisLength
Elongation
Flatness

Table D.2: Shape Based Radiomic Features

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
1.	single_vox	58.3	64.9	59.7	60.0	64.3	61.1	60.9	65.0	62.5	92
2.	single_vox-single_targets	66.1	67.1	61.4	65.2	69.0	64.6	65.7	69.1	66.1	1576
3.	single_vox-single_roi	62.1	65.4	58.2	62.5	67.2	61.5	63.0	67.3	63.5	304
4.	single_vox-single_brain	66.1	66.9	66.3	65.9	69.4	66.5	66.4	69.3	67.7	198
5.	single_vox-many_brain	64.8	66.7	63.7	64.4	68.4	64.1	64.8	68.4	65.7	474
6.	many_vox-single_brain	78.6	77.7	77.4	78.5	79.4	77.5	78.7	79.0	78.4	934
7.	many_vox-CLR	80.0	79.3	79.6	80.1	81.1	79.5	80.4	80.9	80.0	828
8.	many_vox-CL	79.6	79.4	76.6	79.6	80.8	76.9	80.7	81.1	78.6	828
9.	many_vox-CR	81.9	82.2	79.8	82.3	83.5	79.9	81.9	83.2	79.3	828
10.	many_vox-HLR	78.2	79.0	74.6	77.9	78.6	77.1	78.2	78.7	78.6	828
11.	many_vox-HL	76.7	73.9	73.2	76.1	75.9	74.7	78.2	78.2	78.7	828
12.	many_vox-HR	79.3	80.9	73.7	79.4	80.3	76.6	78.2	78.9	76.6	828
13.	many_vox-HLR-clinical_CAP	80.7	76.4	81.0	80.3	77.1	81.1	80.7	78.3	79.4	829
14.	many_vox-HLR-clinical_UHDRS	76.8	73.3	76.6	76.5	74.3	77.5	76.9	75.6	75.7	832
15.	many_vox-HLR-clinical_all	75.0	66.7	72.9	74.1	72.1	75.3	74.5	73.4	73.4	919
16.	many_vox-CHLR	79.3	80.5	78.5	79.3	81.2	80.0	79.5	81.2	81.0	828
17.	many_vox-CHL	77.5	77.4	74.3	77.5	79.1	76.6	79.0	80.4	79.3	828
18.	many_vox-CHR	81.2	82.7	79.9	81.3	83.3	81.1	80.6	82.6	80.7	828
19.	many_vox-CLR_coords	79.5	79.3	79.0	79.7	81.0	79.0	80.0	80.8	79.6	831
20.	many_vox-CLR_scale	79.6	79.0	78.9	79.8	80.6	78.8	80.0	80.3	79.3	828
21.	many_vox-CLR_b50	80.1	79.5	77.8	80.2	80.2	77.9	80.5	80.1	78.4	828
22.	many_vox-CLR-balance_025	51.6	68.5	73.9	80.3	81.0	79.1	80.5	80.7	79.6	828
23.	many_vox-CLR-balance_050	45.8	66.6	72.7	79.7	80.9	79.1	79.9	80.7	79.6	828
24.	many_vox-CLR-balance_075	41.8	63.5	70.6	80.6	81.0	79.0	80.8	80.8	79.6	828
25.	many_vox-CLR-balance_100	41.1	63.1	67.9	79.1	80.1	78.0	79.3	79.8	78.8	828
27.	many_vox-CLR-reinclude	80.5	79.1	79.3	80.7	81.1	79.3	81.1	80.9	79.7	828

Table D.3: Hyperparameter Tuning: Diffusion Fractional Anisotropy - Native T1

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
1.	single_vox	63.1	70.1	62.2	66.2	69.8	68.2	66.6	69.9	69.2	92
2.	single_vox-single_targets	70.6	71.2	64.3	69.4	71.3	66.5	69.5	71.1	67.6	1576
3.	single_vox-single_roi	67.1	69.0	65.5	66.2	69.8	66.7	66.5	69.8	67.6	304
4.	single_vox-single_brain	67.1	70.6	64.6	67.0	70.8	67.0	67.3	70.7	68.0	198
5.	single_vox-many_brain	67.6	70.3	66.5	67.3	70.8	67.5	67.5	70.8	68.4	474
6.	many_vox-single_brain	81.0	80.7	78.8	80.9	80.5	79.4	81.1	80.5	79.9	934
7.	many_vox-CLR	80.1	79.7	76.3	80.5	79.8	77.8	80.8	79.9	78.7	828
8.	many_vox-CL	77.2	78.1	72.5	77.9	78.5	77.2	79.1	79.5	79.7	828
9.	many_vox-CR	79.9	80.0	77.2	80.6	79.3	78.3	80.1	79.0	78.3	828
10.	many_vox-HLR	79.9	74.3	78.6	79.5	76.2	78.5	79.7	76.1	80.5	828
11.	many_vox-HL	72.1	58.9	71.4	72.6	69.4	71.5	75.3	72.2	76.3	828
12.	many_vox-HR	75.6	75.7	74.9	75.7	76.2	75.0	74.2	74.7	75.0	828
13.	many_vox-HLR-clinical_CAP	75.1	75.2	72.5	75.4	76.0	73.4	75.8	76.8	72.6	829
14.	many_vox-HLR-clinical_UHDRS	80.2	77.4	78.2	79.8	77.7	78.2	80.2	78.5	77.2	832
15.	many_vox-HLR-clinical_all	76.2	70.5	71.5	75.3	73.5	72.9	75.7	74.7	71.8	919
16.	many_vox-CHLR	75.8	74.4	72.6	76.4	75.8	76.4	76.6	75.8	77.7	828
17.	many_vox-CHL	77.4	72.1	72.6	77.6	74.2	77.3	79.2	75.8	80.3	828
18.	many_vox-CHR	75.5	75.6	73.1	76.2	76.0	75.3	75.1	74.8	75.4	828
19.	many_vox-CLR_coords	80.2	80.2	77.4	80.6	80.2	78.4	80.9	80.4	79.3	831
20.	many_vox-CLR_scale	80.8	80.3	75.9	81.2	80.1	77.5	81.5	80.2	78.4	828
21.	many_vox-CLR_b50	79.8	80.6	74.9	80.3	80.2	77.2	80.6	80.1	78.3	828
22.	many_vox-CLR-balance_025	57.0	68.0	74.6	80.8	80.0	78.2	81.1	80.2	79.2	828
23.	many_vox-CLR-balance_050	51.8	64.9	72.6	81.0	80.2	78.3	81.2	80.3	79.2	828
24.	many_vox-CLR-balance_075	47.3	62.6	71.0	80.9	79.8	77.4	81.2	80.0	78.4	828
25.	many_vox-CLR-balance_100	43.8	62.5	70.4	79.6	79.3	77.8	79.9	79.5	78.8	828
28.	many_vox-single_brain-coords	81.4	81.0	79.0	81.3	80.8	79.7	81.6	80.8	80.3	937

Table D.4: Hyperparameter Tuning: Diffusion Fractional Anisotropy - Native T1/T2

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
1.	single_vox	52.5	58.9	53.4	54.2	60.3	54.9	54.4	59.3	54.2	92
2.	single_vox-single_targets	60.1	56.7	58.8	60.3	64.0	61.4	60.5	63.8	60.8	1380
3.	single_vox-single_roi	65.2	67.1	63.3	64.0	68.9	65.6	64.3	68.6	65.2	276
4.	single_vox-single_brain	65.3	66.1	65.8	64.5	68.9	66.0	64.8	68.6	65.9	184
5.	single_vox-many_brain	60.6	62.7	61.8	61.5	66.1	62.5	61.8	65.8	62.0	460
6.	many_vox-single_brain	81.9	82.3	80.8	81.4	83.2	81.0	81.6	83.4	81.5	920
7.	many_vox-CLR	82.3	82.1	80.8	82.0	83.4	81.3	82.3	83.7	81.7	828
8.	many_vox-CL	81.8	81.3	79.8	82.5	83.9	82.9	81.8	83.5	81.6	828
9.	many_vox-CR	83.6	83.2	82.8	82.8	83.5	81.5	83.9	84.2	82.9	828
10.	many_vox-HLR	76.4	72.9	64.0	76.4	75.8	70.0	75.4	74.3	68.5	828
11.	many_vox-HL	77.0	71.8	66.3	79.2	77.9	73.2	76.7	74.4	69.3	828
12.	many_vox-HR	75.0	73.5	60.9	73.0	73.5	66.6	74.1	74.5	67.9	828
13.	many_vox-HLR-clinical_CAP	77.2	77.5	66.1	77.6	78.0	64.1	76.5	76.3	64.0	829
14.	many_vox-HLR-clinical_UHDRS	75.7	75.6	64.6	76.2	76.6	62.8	75.1	74.8	62.6	832
15.	many_vox-HLR-clinical_all	71.8	69.7	65.3	72.4	72.3	65.7	71.0	70.7	64.8	919
16.	many_vox-CHLR	80.8	80.1	75.7	80.6	81.5	79.1	80.3	81.1	78.7	828
17.	many_vox-CHL	79.9	77.7	73.9	81.2	81.9	79.6	79.6	80.0	77.5	828
18.	many_vox-CHR	81.5	82.0	77.5	80.5	81.5	79.0	81.3	82.2	79.9	828
19.	many_vox-CLR_coords	81.7	81.6	81.0	81.5	83.3	81.4	81.7	83.5	81.7	831
20.	many_vox-CLR_scale	84.4	82.9	82.0	84.1	84.4	82.4	84.3	84.8	82.6	828
22.	many_vox-CLR-balance_025	58.2	80.7	85.6	82.7	83.5	81.7	82.9	83.8	82.0	828
23.	many_vox-CLR-balance_050	55.3	79.8	86.5	84.0	84.2	82.4	84.3	84.6	82.7	828
24.	many_vox-CLR-balance_075	51.2	79.2	86.0	83.4	84.0	81.9	83.7	84.4	82.1	828
25.	many_vox-CLR-balance_100	46.2	78.2	84.8	81.9	83.1	81.2	82.1	83.4	81.5	828
27.	many_vox-CLR-reinclude	84.6	83.3	82.3	84.4	84.6	82.8	84.6	84.9	82.9	828
28.	many_vox-scale-bal50-reinclude	96.3	83.8	81.7	76.8	77.1	77.3	76.8	76.6	76.9	828
29.	many_vox-bal50-reinclude	96.2	83.1	83.9	76.9	77.1	77.7	76.9	76.6	77.3	828

Table D.5: Hyperparameter Tuning: Diffusion Fractional Anisotropy - Normalized T1

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
1.	single_vox	58.7	66.4	59.8	61.4	66.7	65.0	61.8	66.2	65.1	92
2.	single_vox-single_targets	59.8	58.0	58.8	58.8	63.9	61.5	59.0	63.8	61.4	1380
3.	single_vox-single_roi	64.3	63.0	66.7	63.0	65.8	66.1	63.3	66.0	66.2	276
4.	single_vox-single_brain	66.6	68.8	67.1	65.8	69.1	67.4	66.1	68.9	67.8	184
5.	single_vox-many_brain	67.3	68.8	66.9	66.3	69.0	67.2	66.5	68.9	67.6	460
6.	many_vox-single_brain	83.4	82.3	82.2	83.1	82.3	82.7	83.2	82.8	82.7	920
7.	many_vox-CLR	80.7	80.9	79.2	80.9	80.3	81.1	80.9	80.7	81.0	828
8.	many_vox-CL	79.1	79.8	73.8	80.4	80.3	81.0	79.5	79.5	79.1	828
9.	many_vox-CR	81.9	81.8	81.9	81.4	80.6	80.8	82.0	81.3	81.8	828
10.	many_vox-HLR	73.9	62.3	66.0	74.5	69.4	68.5	73.6	66.6	67.2	828
11.	many_vox-HL	78.1	58.6	66.3	79.9	70.1	72.1	77.7	65.1	67.3	828
12.	many_vox-HR	74.0	66.5	68.0	72.5	69.0	66.0	73.7	68.2	68.9	828
13.	many_vox-HLR-clinical_CAP	78.0	72.9	69.8	78.6	74.2	70.2	77.7	73.8	68.5	829
14.	many_vox-HLR-clinical_UHDRS	80.3	77.2	71.2	80.6	77.6	71.8	79.7	77.1	70.0	832
15.	many_vox-HLR-clinical_all	69.6	69.0	61.2	70.1	71.2	62.2	68.7	70.0	59.9	919
16.	many_vox-CHLR	79.9	73.1	75.5	79.9	75.7	78.0	79.5	74.6	77.3	828
17.	many_vox-CHL	78.9	70.7	69.7	80.4	75.9	78.1	78.9	73.0	75.2	828
18.	many_vox-CHR	72.1	73.1	69.5	71.4	72.8	69.5	72.7	73.4	72.0	828
19.	many_vox-CLR_coords	83.2	82.6	80.8	83.3	81.6	82.1	83.3	82.2	82.0	831
20.	many_vox-CLR_scale	79.9	80.3	79.1	79.9	79.6	80.8	80.0	80.0	80.8	828
22.	many_vox-CLR-balance_025	68.0	80.2	87.3	81.2	80.3	81.3	81.2	80.9	81.3	828
23.	many_vox-CLR-balance_050	61.6	78.2	87.4	82.0	80.8	81.7	82.1	81.3	81.6	828
24.	many_vox-CLR-balance_075	63.0	78.0	86.5	79.7	79.5	80.5	79.8	79.9	80.6	828
25.	many_vox-CLR-balance_100	59.3	77.4	88.9	83.4	81.4	82.0	83.4	82.0	81.8	828
28.	many_vox-brain-coords-bal50	96.5	82.4	83.6	77.6	79.9	77.1	77.3	80.1	76.7	923
29.	many_vox-brain-coords	83.0	82.5	81.8	82.7	82.4	82.4	82.7	82.9	82.3	923

Table D.6: Hyperparameter Tuning: Diffusion Fractional Anisotropy - Normalized T1/T2

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
27.	normalized-t1-many_vox-CLR-reinclude	84.6	83.3	82.3	84.4	84.6	82.8	84.6	84.9	82.9	828
1.	layers_2048	82.0	82.2	81.1	81.9	83.1	81.3	82.1	83.4	81.7	828
2.	layers_1024	84.1	83.1	82.1	84.0	84.4	82.5	84.1	84.7	82.8	828
3.	layers_512	83.2	83.1	81.5	83.1	84.3	82.2	83.2	84.6	82.5	828
4.	activation_sigmoid	83.4	82.9	81.9	83.3	84.2	82.3	83.5	84.5	82.6	828
5.	activation_relu	85.7	83.1	82.3	85.5	84.0	82.3	85.7	84.3	82.8	828
6.	activation_silu	84.7	83.1	82.7	84.6	84.3	82.4	84.7	84.6	83.0	828
7.	activation_elu	7.2	4.9	7.4	14.3	14.8	14.9	12.5	11.1	13.8	828
8.	batch_10e5	80.2	81.2	80.1	80.3	82.3	80.2	80.4	82.6	80.7	828
9.	batch_10e4	83.1	82.7	81.7	83.0	83.8	82.1	83.2	84.2	82.4	828
10.	batch_10e3	85.5	83.8	82.5	85.3	84.9	82.8	85.5	85.2	83.1	828
11.	batch_10e2	89.3	83.4	83.6	89.1	84.2	84.2	89.3	84.4	84.3	828
12.	dropout_000	89.0	84.8	83.0	88.9	86.0	83.6	88.9	86.0	84.0	828
13.	dropout_025	89.4	84.2	83.9	89.2	85.0	84.0	89.3	85.4	84.4	828
14.	dropout_050	84.6	83.5	82.2	84.4	84.5	82.6	84.6	84.9	82.9	828
15.	dropout_075	89.5	84.2	83.0	89.3	85.7	83.3	89.3	85.9	83.7	828
16.	dropout_090	84.5	83.8	82.4	84.3	84.8	82.7	84.5	85.1	82.9	828
17.	patience_4	85.3	83.4	82.6	85.2	84.9	82.7	85.3	85.1	83.1	828
18.	patience_7	86.3	83.1	82.8	86.2	84.9	83.0	86.3	85.1	83.4	828
19.	patience_10	87.1	84.4	83.3	87.1	85.1	83.8	87.2	85.3	84.2	828
20.	rate_10e-4	82.9	82.7	80.8	82.8	84.0	81.6	82.9	84.3	81.9	828
21.	rate_10e-3	85.2	83.9	82.1	85.1	85.0	82.7	85.3	85.2	83.0	828
22.	rate_10e-2	nan	nan	nan	nan	nan	nan	nan	nan	nan	828

Table D.7: Architecture Tuning: Diffusion Fractional Anisotropy

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
1.	single_vox	92.6	94.3	94.2	93.1	94.6	94.2	93.9	94.9	95.6	92
2.	single_vox-single_targets	93.9	93.6	93.7	93.6	94.4	94.5	94.3	94.6	95.7	1576
3.	single_vox-single_roi	93.7	94.9	92.8	93.6	94.6	94.6	94.3	94.8	95.8	304
4.	single_vox-single_brain	93.4	94.9	94.9	93.6	94.9	94.8	94.3	95.1	96.0	198
5.	single_vox-many_brain	93.2	94.6	94.5	93.4	94.5	94.6	94.1	94.8	95.8	474
6.	many_vox-single_brain	94.7	95.8	94.8	94.7	95.5	95.1	95.4	95.7	96.3	934
7.	many_vox-CLR	94.0	95.3	94.9	94.2	95.2	94.9	94.9	95.4	96.1	828
8.	many_vox-CL	93.8	94.4	94.6	94.2	94.6	94.7	95.7	95.7	96.6	828
9.	many_vox-CR	95.0	96.4	94.7	95.2	96.1	94.9	95.0	95.6	95.0	828
10.	many_vox-HLR	92.4	94.1	93.8	92.4	94.2	93.8	93.3	94.6	94.9	828
11.	many_vox-HL	91.7	94.2	95.8	91.9	94.7	95.5	94.1	96.1	97.3	828
12.	many_vox-HR	91.8	93.7	92.2	91.6	95.1	93.4	91.2	94.3	93.6	828
13.	many_vox-HLR-clinical_CAP	93.3	89.4	95.2	93.2	90.9	95.1	93.9	92.3	95.3	829
14.	many_vox-HLR-clinical_UHDRS	93.6	89.8	94.7	93.4	90.8	94.9	94.1	92.3	95.1	832
15.	many_vox-HLR-clinical_all	93.2	88.6	94.3	92.5	90.4	94.6	93.3	92.0	94.9	919
16.	many_vox-CHLR	93.6	95.1	94.8	93.7	95.4	95.0	94.4	95.7	96.1	828
17.	many_vox-CHL	93.6	94.3	94.4	94.0	95.3	95.4	95.6	96.4	97.0	828
18.	many_vox-CHR	93.7	95.5	93.4	93.9	95.9	94.5	93.5	95.3	94.8	828
19.	many_vox-CLR_coords	94.5	95.7	95.0	94.6	95.5	95.1	95.3	95.7	96.2	831
20.	many_vox-CLR_scale	94.1	95.3	94.9	94.2	95.2	95.0	94.9	95.4	96.1	828
21.	many_vox-CLR_b50	93.8	95.1	94.2	93.9	94.8	94.5	94.8	95.1	95.8	828
22.	many_vox-CLR-balance_025	95.5	95.9	95.8	94.3	95.3	95.0	95.0	95.5	96.1	828
23.	many_vox-CLR-balance_050	95.3	95.8	95.8	94.3	95.2	95.1	95.0	95.4	96.3	828
24.	many_vox-CLR-balance_075	95.2	95.3	95.8	94.0	95.0	94.9	94.8	95.2	96.1	828
25.	many_vox-CLR-balance_100	95.4	95.7	95.8	94.3	95.2	95.0	95.0	95.4	96.2	828
27.	many_vox-CLR-reinclude	94.2	95.2	94.9	94.3	95.0	94.9	95.1	95.3	96.2	828
28.	many_vox-bal25-coords-brain	97.4	96.7	96.7	93.2	94.0	93.6	94.1	94.4	95.1	937
29.	many_vox-single_brain-coords	94.0	95.3	94.9	94.2	95.1	95.0	94.9	95.4	96.2	937

Table D.8: Hyperparameter Tuning: Mean Diffusivity - Native T1

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
1.	single_vox	93.0	93.8	94.1	93.6	94.7	94.3	94.3	95.0	95.6	92
2.	single_vox-single_targets	94.4	93.9	91.7	93.9	94.7	93.8	94.4	94.9	95.0	1576
3.	single_vox-single_roi	94.2	95.2	94.3	93.8	94.8	94.6	94.5	95.1	95.7	304
4.	single_vox-single_brain	93.9	95.5	92.2	93.8	94.9	94.0	94.4	95.1	95.2	198
5.	single_vox-many_brain	94.0	95.4	91.7	93.7	94.8	93.8	94.3	95.0	95.1	474
6.	many_vox-single_brain	94.1	95.2	93.8	94.1	95.0	94.5	94.7	95.3	95.7	934
7.	many_vox-CLR	94.9	95.1	94.4	95.1	95.8	94.8	95.6	96.1	95.6	828
8.	many_vox-CL	94.1	93.2	95.5	94.6	95.3	95.4	95.9	96.4	96.9	828
9.	many_vox-CR	95.0	96.3	91.4	95.3	96.6	92.8	94.9	96.1	92.0	828
10.	many_vox-HLR	92.2	94.6	95.5	92.6	94.7	95.3	93.3	94.9	96.4	828
11.	many_vox-HL	92.9	94.0	96.1	93.5	94.5	95.9	95.3	96.0	97.6	828
12.	many_vox-HR	92.7	94.9	94.6	93.3	95.1	95.5	92.8	94.2	95.4	828
13.	many_vox-HLR-clinical_CAP	92.7	90.5	93.3	92.9	91.9	94.2	93.6	93.0	94.3	829
14.	many_vox-HLR-clinical_UHDRS	93.0	90.6	93.5	92.8	91.8	94.1	93.5	93.0	94.2	832
15.	many_vox-HLR-clinical_all	92.8	88.9	93.9	92.2	90.9	94.3	93.1	92.1	94.5	919
16.	many_vox-CHLR	93.3	94.7	94.9	93.7	95.3	95.0	94.4	95.6	96.1	828
17.	many_vox-CHL	93.0	92.8	95.6	93.7	95.0	95.6	95.3	96.3	97.1	828
18.	many_vox-CHR	93.9	96.2	94.5	94.5	96.4	94.9	94.1	95.8	94.7	828
19.	many_vox-CLR_coords	94.8	95.2	94.3	95.0	95.9	94.6	95.5	96.2	95.6	831
20.	many_vox-CLR_scale	94.4	94.9	94.3	94.7	95.8	94.5	95.3	96.1	95.4	828
21.	many_vox-CLR_b50	94.7	95.1	94.5	94.9	95.8	94.6	95.4	96.1	95.6	828
22.	many_vox-CLR-balance_025	96.4	95.4	95.7	94.7	95.8	94.6	95.2	96.1	95.6	828
23.	many_vox-CLR-balance_050	96.6	95.5	95.1	94.9	95.8	94.7	95.5	96.1	95.6	828
24.	many_vox-CLR-balance_075	96.3	95.2	95.4	94.6	95.7	94.6	95.2	96.0	95.6	828
25.	many_vox-CLR-balance_100	96.3	95.3	95.2	94.7	95.7	94.9	95.3	96.0	95.9	828

Table D.9: Hyperparameter Tuning: Mean Diffusivity - Native T1/T2

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
1.	single_vox	87.6	88.6	88.7	88.7	91.7	91.5	87.8	90.8	90.0	92
2.	single_vox-single_targets	90.1	89.7	91.2	90.1	92.3	92.9	89.3	91.5	91.5	1380
3.	single_vox-single_roi	90.5	91.3	88.7	90.6	92.6	93.0	89.7	91.7	91.5	276
4.	single_vox-single_brain	89.8	90.2	89.9	90.3	92.5	93.1	89.5	91.6	91.7	184
5.	single_vox-many_roi	90.4	90.8	89.4	90.4	92.5	92.9	89.5	91.7	91.4	828
6.	many_vox-single_roi	94.3	93.7	92.0	94.2	95.0	94.8	93.8	94.8	93.9	1012
7.	many_vox-CLR	93.9	93.6	92.5	94.1	95.0	95.1	93.8	94.7	94.3	828
8.	many_vox-CL	94.3	92.3	93.1	95.4	95.2	96.2	94.4	94.1	95.2	828
9.	many_vox-CR	92.2	94.1	91.4	90.4	93.9	91.6	92.2	94.8	92.5	828
10.	many_vox-HLR	92.0	89.2	86.9	92.9	90.9	90.8	91.6	89.7	88.1	828
11.	many_vox-HL	92.1	90.0	88.2	93.8	91.7	92.3	91.9	89.7	89.0	828
12.	many_vox-HR	87.7	85.3	81.0	85.4	84.8	84.4	87.0	86.3	85.9	828
13.	many_vox-HLR-clinical_CAP	90.0	90.2	78.1	91.1	92.0	85.1	89.5	90.8	82.3	829
14.	many_vox-HLR-clinical_UHDRS	82.9	87.5	77.7	86.1	88.2	85.7	83.6	86.7	83.1	832
15.	many_vox-HLR-clinical_all	84.8	83.3	80.7	86.5	88.0	86.3	84.1	86.5	83.8	919
16.	many_vox-CHLR	91.8	91.9	90.3	92.7	92.6	93.2	91.7	91.9	91.6	828
17.	many_vox-CHL	93.1	92.2	91.6	94.7	93.9	94.8	93.3	92.5	92.9	828
18.	many_vox-CHR	91.1	92.1	89.5	89.4	90.8	89.6	91.0	91.9	90.4	828
19.	many_vox-CLR_coords	94.0	93.8	93.2	94.2	95.1	95.3	93.9	94.9	94.5	831
20.	many_vox-CLR_scale	93.7	93.5	92.4	94.0	94.9	94.9	93.6	94.6	94.0	828
22.	many_vox-CLR-balance_025	95.2	95.0	94.6	94.1	95.0	95.1	93.7	94.7	94.3	828
23.	many_vox-CLR-balance_050	95.0	94.5	94.2	94.1	95.0	95.2	93.7	94.7	94.4	828
24.	many_vox-CLR-balance_075	95.0	94.3	94.6	94.2	95.0	95.2	93.8	94.7	94.4	828
25.	many_vox-CLR-balance_100	95.4	94.3	94.5	94.3	95.0	95.1	93.9	94.7	94.3	828
27.	many_vox-CLR-reinclude	93.5	93.3	93.6	93.9	94.8	95.2	93.4	94.6	94.5	828
28.	many_vox-single_roi-coords	94.8	93.8	92.6	94.6	95.2	95.1	94.2	94.9	94.3	1015

Table D.10: Hyperparameter Tuning: Mean Diffusivity - Normalized T1

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
1.	single_vox	86.5	89.4	90.2	88.9	90.8	92.6	88.2	90.7	91.7	92
2.	single_vox-single_targets	90.5	91.1	91.9	90.1	91.8	93.3	89.5	91.7	92.3	1380
3.	single_vox-single_roi	90.5	91.3	90.5	90.4	91.8	93.5	89.8	91.7	92.5	276
4.	single_vox-single_brain	91.1	91.3	92.0	90.9	92.0	93.6	90.3	92.0	92.6	184
5.	single_vox-many_brain	90.7	91.4	92.6	90.5	91.9	93.7	89.9	91.8	92.7	460
6.	many_vox-single_brain	93.5	93.6	94.1	93.7	94.7	94.8	93.3	94.7	94.1	920
7.	many_vox-CLR	92.4	91.2	93.5	93.0	93.8	94.4	92.6	93.9	93.8	828
8.	many_vox-CL	90.1	90.2	92.1	92.4	93.2	94.4	91.0	92.2	93.3	828
9.	many_vox-CR	91.3	91.9	92.5	90.0	93.2	90.6	91.9	94.4	92.1	828
10.	many_vox-HLR	90.0	84.2	89.0	91.4	86.1	91.4	90.1	84.6	89.7	828
11.	many_vox-HL	90.2	84.0	90.5	92.8	88.1	93.2	91.2	85.3	90.9	828
12.	many_vox-HR	86.3	81.8	86.3	85.0	79.6	86.7	86.5	81.8	88.3	828
13.	many_vox-HLR-clinical_CAP	87.1	91.1	84.9	89.7	91.6	88.3	88.0	90.7	86.6	829
14.	many_vox-HLR-clinical_UHDRS	83.0	88.7	76.7	87.3	89.6	85.1	85.3	88.5	82.6	832
15.	many_vox-HLR-clinical_all	88.2	89.2	81.3	89.1	91.5	87.6	87.5	90.6	86.0	919
16.	many_vox-CHLR	91.7	88.8	90.7	92.8	90.6	93.1	92.0	89.7	92.0	828
17.	many_vox-CHL	91.7	88.7	91.6	93.7	91.6	94.5	92.4	89.8	92.8	828
18.	many_vox-CHR	90.6	88.7	89.7	89.1	87.2	89.0	90.7	88.6	90.4	828
19.	many_vox-CLR_coords	93.3	91.1	93.9	93.7	94.2	94.7	93.4	94.2	94.1	831
20.	many_vox-CLR_scale	92.5	91.3	93.4	93.1	93.9	94.4	92.7	94.0	93.8	828
22.	many_vox-CLR-balance_025	93.6	91.7	95.2	93.3	93.9	94.6	92.9	94.0	94.0	828
23.	many_vox-CLR-balance_050	93.4	91.3	94.9	93.3	94.0	94.5	92.9	94.0	93.9	828
24.	many_vox-CLR-balance_075	93.2	90.7	94.9	93.3	93.9	94.5	93.0	94.0	93.9	828
25.	many_vox-CLR-balance_100	93.2	90.4	94.8	93.4	94.0	94.5	93.0	94.0	93.9	828
28.	many_vox-bal100-coords-brain	96.9	95.8	93.0	92.3	93.8	93.3	91.4	93.3	92.2	923
29.	many_vox-single_brain-coords	93.7	93.7	94.2	93.8	94.8	94.9	93.5	94.8	94.3	923

Table D.11: Hyperparameter Tuning: Mean Diffusivity - Normalized T1/T2

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
6.	native-t1-many_vox-single_brain	94.7	95.8	94.8	94.7	95.5	95.1	95.4	95.7	96.3	934
1.	layers_2048	94.6	95.8	95.1	94.6	95.5	95.2	95.3	95.7	96.4	934
2.	layers_1024	94.2	95.5	95.0	94.3	95.2	95.0	95.0	95.4	96.2	934
3.	layers_512	94.3	95.4	94.9	94.4	95.2	95.0	95.1	95.4	96.2	934
4.	activation_sigmoid	94.2	95.4	95.0	94.3	95.2	95.0	95.0	95.4	96.2	934
5.	activation_relu	94.6	95.5	94.9	94.6	95.2	95.0	95.3	95.4	96.2	934
6.	activation_silu	95.3	95.8	94.5	95.1	95.4	94.9	95.7	95.5	96.0	934
8.	batch_10e5	93.0	94.6	94.3	93.3	94.3	94.3	94.1	94.6	95.8	934
9.	batch_10e4	94.0	95.4	94.9	94.1	95.1	94.9	94.8	95.3	96.1	934
10.	batch_10e3	95.1	95.9	94.7	95.0	95.5	95.1	95.6	95.7	96.2	934
11.	batch_10e2	95.7	95.8	94.5	95.5	95.3	94.9	96.0	95.4	95.8	934
12.	dropout_000	94.9	95.8	94.8	94.9	95.5	95.1	95.5	95.6	96.2	934
13.	dropout_025	95.6	95.9	94.5	95.5	95.4	94.9	96.0	95.5	96.0	934
14.	dropout_050	94.9	95.9	94.7	94.9	95.6	95.0	95.5	95.7	96.2	934
15.	dropout_075	94.5	95.7	95.0	94.5	95.4	95.2	95.2	95.6	96.3	934
16.	dropout_090	94.7	95.8	94.7	94.7	95.4	95.0	95.4	95.6	96.2	934
17.	patience_4	94.2	95.5	95.0	94.3	95.2	95.0	95.0	95.4	96.2	934
18.	patience_7	95.2	95.9	94.7	95.1	95.5	95.1	95.7	95.6	96.1	934
19.	patience_10	94.6	95.9	94.9	94.6	95.5	95.1	95.3	95.7	96.2	934
20.	rate_10e-4	94.4	95.6	94.9	94.5	95.3	95.2	95.2	95.5	96.3	934
21.	rate_10e-3	94.9	95.9	94.8	94.8	95.5	95.2	95.5	95.6	96.3	934
22.	rate_10e-2	nan	nan	nan	nan	nan	nan	nan	nan	nan	934

Table D.12: Architecture Tuning: Mean Diffusivity

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
1.	single_vox	63.7	66.6	66.4	63.8	66.3	66.8	63.9	66.3	67.0	92
2.	single_vox-single_targets	65.4	65.6	62.1	65.5	65.3	62.8	65.3	65.0	63.2	1576
3.	single_vox-single_roi	64.9	65.9	61.2	65.0	65.5	61.9	64.7	65.2	62.3	304
4.	single_vox-single_brain	64.2	66.1	65.6	64.4	65.8	66.0	64.0	65.6	66.3	198
5.	single_vox-many_brain	64.2	66.0	65.9	64.3	65.6	66.3	64.2	65.5	66.7	474
6.	many_vox-single_targets	72.6	70.1	67.2	72.7	70.1	67.8	72.2	69.1	68.0	2312
7.	many_vox-single_roi	71.1	71.6	67.6	71.1	71.3	68.0	70.5	70.8	67.7	1040
8.	many_vox-single_brain	72.0	72.2	70.9	72.1	72.0	71.0	71.6	71.2	71.0	934
9.	many_vox-CLR	72.1	72.2	71.4	72.2	71.9	71.5	71.8	71.3	71.5	828
10.	many_vox-CL	71.6	69.3	71.3	71.7	68.7	71.4	71.4	68.1	72.2	828
11.	many_vox-CR	73.9	73.9	69.1	73.9	73.9	69.1	73.6	73.3	68.2	828
12.	many_vox-HLR	70.0	67.6	66.1	69.9	67.9	66.9	69.8	67.1	66.7	828
13.	many_vox-HL	65.8	63.9	66.3	65.6	63.8	67.3	65.2	63.1	68.2	828
14.	many_vox-HR	70.4	67.6	68.3	70.3	68.0	68.5	70.1	66.7	67.2	828
15.	many_vox-HLR-clinical_CAP	70.3	68.8	64.6	70.3	69.0	64.6	70.0	69.6	64.4	829
16.	many_vox-HLR-clinical_UHDRS	69.7	68.1	67.1	69.7	68.5	67.1	69.3	68.8	67.2	832
17.	many_vox-HLR-clinical_all	69.6	65.5	60.4	69.7	66.0	60.4	69.3	66.3	60.2	919
18.	many_vox-CHLR	71.6	70.7	70.4	71.5	70.5	70.2	71.2	70.1	70.2	828
19.	many_vox-CHL	70.9	69.0	70.5	70.8	68.8	70.8	70.6	69.1	71.9	828
20.	many_vox-CHR	70.6	71.7	69.4	70.4	71.7	69.6	70.3	71.1	69.0	828
21.	many_vox-CLR_coords	72.7	72.5	70.8	72.7	72.3	70.8	72.4	71.8	71.0	831
22.	many_vox-CLR_scale	72.0	71.7	69.8	72.0	71.4	69.7	71.5	70.4	69.5	828
23.	many_vox-CLR_b50	70.6	71.0	71.2	70.6	70.8	71.2	70.4	70.6	71.7	828
24.	many_vox-feature_10	71.7	71.7	70.8	71.8	71.5	70.8	71.4	70.9	70.7	738
25.	many_vox-feature_20	71.7	72.1	71.2	71.8	71.9	71.1	71.4	71.1	71.1	648
26.	many_vox-feature_30	71.6	72.1	70.4	71.7	71.9	70.4	71.5	71.4	70.1	558
27.	many_vox-feature_40	71.9	71.9	69.9	72.0	71.7	69.9	71.7	71.2	69.9	468
28.	many_vox-balance_025	63.1	60.8	61.4	65.2	65.6	61.0	65.0	64.9	62.1	828
29.	many_vox-balance_050	67.0	61.5	61.5	60.3	59.8	53.6	60.2	58.8	54.5	828
30.	many_vox-balance_075	69.3	61.7	62.6	59.7	59.4	51.3	59.8	58.2	52.4	828
31.	many_vox-balance_100	64.9	61.3	62.0	52.3	53.7	46.8	52.4	53.2	48.0	828
32.	many_vox-reinclude	72.6	71.9	70.6	72.6	71.6	70.6	72.2	71.0	70.6	828
33.	many_vox-CLR-aug	72.8	72.6	70.3	72.9	72.3	70.3	72.5	71.7	70.6	828

Table D.13: Hyperparameter Tuning: Relative Connectivity - Native T1

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
1.	single_vox	62.7	65.8	66.1	62.8	65.4	66.4	62.7	65.3	66.4	92
2.	single_vox-single_targets	65.0	66.1	61.4	65.2	65.8	61.7	64.9	65.3	61.8	1576
3.	single_vox-single_roi	65.8	66.8	60.3	65.9	66.5	60.4	65.6	66.1	60.2	304
4.	single_vox-single_brain	63.9	66.0	65.6	64.0	65.7	65.7	63.8	65.5	65.5	198
5.	single_vox-many_brain	64.0	65.4	62.7	64.1	65.2	62.7	63.8	64.9	62.3	474
6.	many_vox-single_targets	72.3	71.4	69.7	72.4	71.2	69.9	72.1	70.6	69.8	2312
7.	many_vox-single_roi	73.2	71.6	70.0	73.2	71.5	70.2	72.9	71.0	70.1	1040
8.	many_vox-single_brain	72.2	71.8	71.4	72.3	71.7	71.3	71.9	71.1	71.1	934
9.	many_vox-CLR	70.8	71.8	69.8	70.9	71.5	70.0	70.6	71.4	69.7	828
10.	many_vox-CL	70.4	70.1	68.0	70.5	69.7	68.2	70.1	70.4	68.5	828
11.	many_vox-CR	70.7	72.4	68.6	70.8	72.3	68.7	70.5	71.9	68.5	828
12.	many_vox-HLR	69.0	61.7	68.4	69.0	62.4	68.7	68.7	62.1	68.7	828
13.	many_vox-HL	67.9	64.0	68.3	67.9	64.2	69.6	67.8	64.0	70.6	828
14.	many_vox-HR	68.6	63.0	69.8	68.7	64.0	69.8	68.2	63.5	70.1	828
15.	many_vox-HLR-clinical_CAP	70.3	68.7	65.6	70.3	68.8	65.6	70.0	69.2	65.1	829
16.	many_vox-HLR-clinical_UHDRS	70.6	67.7	66.0	70.6	68.0	66.0	70.3	68.5	65.6	832
17.	many_vox-HLR-clinical_all	69.5	65.0	59.7	69.6	65.3	59.8	69.1	65.6	59.8	919
18.	many_vox-CHLR	70.0	69.2	70.2	69.9	69.1	70.6	69.6	68.6	70.2	828
19.	many_vox-CHL	71.4	68.4	70.1	71.3	68.0	70.9	71.0	67.5	71.0	828
20.	many_vox-CHR	70.5	69.5	71.2	70.4	69.6	71.2	70.2	69.7	71.2	828
21.	many_vox-CLR_coords	73.3	72.6	70.7	73.4	72.5	70.9	73.2	72.7	70.5	831
22.	many_vox-CLR_scale	71.8	71.2	69.6	71.9	71.0	69.8	71.6	71.2	69.2	828
23.	many_vox-CLR_b50	72.3	71.7	70.5	72.4	71.4	70.6	72.1	71.2	70.2	828
24.	many_vox-feature_10	72.2	71.7	69.4	72.3	71.5	69.6	72.0	71.2	69.0	738
25.	many_vox-feature_20	70.9	71.2	68.7	71.0	71.0	69.0	70.7	71.0	68.5	648
26.	many_vox-feature_30	70.9	71.3	69.1	71.0	71.1	69.5	70.7	71.1	69.2	558
27.	many_vox-feature_40	70.3	70.8	68.6	70.4	70.6	68.9	70.1	70.4	68.5	468
28.	many_vox-balance_025	61.4	58.6	57.7	63.3	62.5	58.8	63.2	62.2	58.4	828
29.	many_vox-balance_050	60.2	56.3	56.4	57.4	56.4	52.3	57.3	55.4	52.2	828
30.	many_vox-balance_075	61.8	56.5	58.5	52.7	51.0	48.3	52.8	50.5	48.3	828
31.	many_vox-balance_100	62.2	57.1	58.1	51.0	49.3	46.0	50.9	48.7	46.4	828

Table D.14: Hyperparameter Tuning: Relative Connectivity - Native T1/T2

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
1.	single_vox	63.1	65.1	66.5	63.6	66.0	66.7	63.2	65.5	66.2	92
2.	single_vox-single_targets	63.3	64.2	66.5	63.8	65.3	66.9	63.4	64.7	66.3	1380
3.	single_vox-single_roi	63.5	64.4	59.7	64.0	65.6	59.7	63.7	64.9	59.5	276
4.	single_vox-single_brain	63.5	65.1	66.7	64.0	66.1	67.3	63.7	65.5	66.4	184
5.	single_vox-many_brain	63.2	64.9	66.7	63.7	65.9	67.4	63.4	65.3	66.5	460
6.	many_vox-single_targets	73.9	69.5	69.8	73.5	69.6	69.8	74.1	69.8	69.6	2116
7.	many_vox-single_roi	72.7	69.1	62.2	72.2	69.5	61.9	72.9	69.5	61.9	1012
8.	many_vox-single_brain	70.1	70.5	67.4	69.4	71.1	66.9	70.2	70.8	67.1	920
9.	many_vox-CLR	71.9	71.4	73.0	71.4	71.9	73.1	72.0	71.8	72.9	828
10.	many_vox-CL	70.0	66.1	70.6	69.8	67.7	71.7	70.2	66.2	70.6	828
11.	many_vox-CR	73.1	73.5	71.3	73.1	73.3	70.9	73.1	73.7	71.3	828
12.	many_vox-HLR	68.0	67.4	63.9	68.5	69.3	66.1	67.9	67.8	65.0	828
13.	many_vox-HL	68.8	67.0	63.2	69.6	68.9	67.2	68.9	67.6	65.0	828
14.	many_vox-HR	69.2	68.2	64.3	69.9	70.5	66.3	69.1	68.7	65.0	828
15.	many_vox-HLR-clinical_CAP	69.1	67.9	63.9	69.9	68.1	65.6	69.2	68.0	63.1	829
16.	many_vox-HLR-clinical_UHDRS	68.2	66.9	64.6	68.9	67.0	66.5	68.3	66.9	64.0	832
17.	many_vox-HLR-clinical_all	64.9	62.8	60.7	65.7	62.9	63.2	65.1	63.0	60.7	919
18.	many_vox-CHLR	70.5	70.3	72.1	70.5	71.3	72.7	70.4	70.4	72.0	828
19.	many_vox-CHL	70.8	68.5	70.7	71.1	70.2	72.2	70.8	68.8	70.7	828
20.	many_vox-CHR	71.3	70.8	70.9	71.5	71.8	70.8	71.2	70.9	70.9	828
21.	many_vox-CLR_coords	73.0	71.5	73.0	72.6	72.1	73.3	73.2	71.9	73.0	831
22.	many_vox-CLR_scale	71.2	71.0	72.1	71.0	71.4	72.1	71.3	71.2	72.0	828
23.	many_vox-feature_10	70.3	70.8	71.5	70.1	71.6	71.5	70.4	71.1	71.4	738
24.	many_vox-feature_20	71.5	70.5	71.6	71.3	71.3	71.9	71.6	70.8	71.5	648
25.	many_vox-feature_30	69.4	69.8	70.1	69.1	70.2	70.0	69.6	70.1	70.1	558
26.	many_vox-feature_40	69.9	69.9	70.4	69.8	70.6	70.4	70.1	70.2	70.3	468
28.	many_vox-balance_025	63.1	58.3	63.3	64.6	64.0	64.1	65.1	63.5	63.4	828
29.	many_vox-balance_050	66.5	57.1	63.5	60.6	59.4	56.4	60.8	58.5	55.8	828
30.	many_vox-balance_075	66.8	57.1	64.7	55.8	55.1	52.3	56.1	54.2	51.7	828
31.	many_vox-balance_100	67.6	56.9	64.7	54.7	54.5	50.7	54.9	53.1	50.2	828
32.	many_vox-reinclude	71.5	71.0	71.2	71.3	71.8	71.5	71.7	71.4	71.1	828

Table D.15: Hyperparameter Tuning: Relative Connectivity - Normalized T1

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
1.	single_vox	62.4	64.2	66.5	62.8	65.2	66.6	62.6	64.7	66.2	92
2.	single_vox-single_targets	62.4	64.2	66.5	62.8	65.2	66.6	62.6	64.7	66.2	1380
3.	single_vox-single_roi	63.1	63.9	66.3	63.5	65.0	66.6	63.3	64.3	66.1	276
4.	single_vox-single_brain	63.4	63.0	64.0	63.9	63.7	63.8	63.5	63.2	63.9	184
5.	single_vox-many_brain	63.2	62.4	63.5	63.8	63.2	63.4	63.4	62.4	63.4	460
6.	many_vox-single_targets	71.0	68.5	66.1	70.7	68.9	65.3	71.1	68.6	65.9	2116
7.	many_vox-single_roi	70.5	68.9	69.6	70.4	69.4	69.3	70.6	69.0	69.6	1012
8.	many_vox-single_brain	71.8	69.6	68.0	71.6	69.7	67.3	71.9	69.6	68.0	920
9.	many_vox-CLR	71.6	70.3	71.1	71.7	71.5	70.4	71.7	70.5	70.9	828
10.	many_vox-CL	70.6	68.0	69.9	70.7	70.3	70.0	70.7	68.1	69.9	828
11.	many_vox-CR	72.0	72.1	71.3	72.2	72.0	70.6	72.1	72.3	71.2	828
12.	many_vox-HLR	68.4	63.6	70.0	68.9	65.5	71.8	68.3	63.7	70.3	828
13.	many_vox-HL	66.8	61.9	67.8	67.3	64.4	70.1	66.8	62.1	67.9	828
14.	many_vox-HR	68.9	62.6	69.4	69.7	65.0	71.1	68.8	63.2	69.6	828
15.	many_vox-HLR-clinical_CAP	66.4	66.5	59.3	67.2	66.2	61.7	66.5	66.0	59.1	829
16.	many_vox-HLR-clinical_UHDRS	68.5	66.8	62.1	69.3	66.9	65.0	68.5	66.6	62.2	832
17.	many_vox-HLR-clinical_all	65.4	64.6	62.1	66.2	64.8	64.7	65.4	64.5	62.1	919
18.	many_vox-CHLR	70.0	66.9	71.1	70.2	68.1	71.2	70.0	67.0	71.0	828
19.	many_vox-CHL	69.4	66.1	70.7	69.7	68.0	71.3	69.4	66.1	70.6	828
20.	many_vox-CHR	69.4	67.3	70.2	69.9	68.2	70.2	69.3	67.3	70.2	828
21.	many_vox-CLR_coords	73.8	71.1	71.8	73.6	71.9	70.8	74.0	71.3	71.7	831
22.	many_vox-CLR_scale	71.5	69.5	71.9	71.5	70.6	71.3	71.7	69.7	71.8	828
23.	many_vox-feature_10	70.2	69.8	71.5	70.4	71.2	71.1	70.3	70.0	71.4	738
24.	many_vox-feature_20	71.4	70.3	71.0	71.3	71.3	70.2	71.5	70.7	70.8	648
25.	many_vox-feature_30	72.3	69.9	70.4	72.3	71.1	69.6	72.4	70.2	70.1	558
26.	many_vox-feature_40	70.1	69.8	71.5	70.3	71.0	70.8	70.3	70.0	71.3	468
28.	many_vox-balance_025	67.6	55.9	62.5	67.3	61.7	60.4	67.4	61.0	60.4	828
29.	many_vox-balance_050	63.0	54.4	60.9	58.0	54.9	52.8	58.3	53.8	52.7	828
30.	many_vox-balance_075	61.5	54.9	61.5	50.2	47.8	45.0	50.6	46.9	45.4	828
31.	many_vox-balance_100	64.2	54.9	64.0	49.4	46.4	43.8	49.7	45.2	44.2	828

Table D.16: Hyperparameter Tuning: Relative Connectivity - Normalized T1/T2

ID	Experiment	Raw			Native			Normalized			Input Layer
		Train	Val	Test	Train	Val	Test	Train	Val	Test	
21.	normalized-t1-many_vox-CLR_coords	73.0	71.5	73.0	72.6	72.1	73.3	73.2	71.9	73.0	831
1.	layers_2048	72.7	71.6	72.8	72.4	72.4	73.0	72.9	71.9	72.6	831
2.	layers_1024	71.7	71.7	73.4	71.5	72.6	73.4	71.8	72.0	73.2	831
3.	layers_512	71.9	71.7	72.2	71.6	72.7	72.3	72.0	72.1	72.1	831
4.	activation_sigmoid	71.7	71.7	73.4	71.5	72.6	73.4	71.8	72.0	73.2	831
5.	activation_relu	73.4	71.8	72.1	73.2	72.8	72.3	73.5	72.1	72.0	831
6.	activation_silu	72.5	70.0	70.5	72.2	70.5	70.7	72.6	70.3	70.3	831
7.	activation_elu	71.4	70.9	72.9	71.1	71.6	72.7	71.5	71.3	72.8	831
8.	batch_10e5	69.4	71.1	70.6	69.1	71.4	70.4	69.5	71.4	70.5	831
9.	batch_10e4	71.7	71.7	73.4	71.5	72.6	73.4	71.8	72.0	73.2	831
10.	batch_10e3	73.0	71.7	72.7	72.7	72.5	72.6	73.1	72.1	72.6	831
11.	batch_10e2	70.7	70.5	70.8	70.4	70.8	71.0	70.9	70.7	70.7	831
12.	dropout_000	73.0	71.7	72.7	72.7	72.5	72.6	73.1	72.1	72.6	831
13.	dropout_025	73.5	71.1	73.1	73.2	71.7	73.3	73.6	71.5	73.0	831
14.	dropout_050	72.8	70.6	71.8	72.7	71.4	72.0	72.9	70.9	71.8	831
15.	dropout_075	73.2	71.1	72.7	73.0	72.0	72.6	73.4	71.4	72.6	831
16.	dropout_090	71.2	71.5	71.3	70.9	72.1	71.3	71.3	71.8	71.2	831
17.	patience_4	71.9	71.6	72.4	71.9	72.5	72.5	72.0	71.9	72.2	831
18.	patience_7	73.2	71.1	72.7	73.0	72.0	72.6	73.4	71.4	72.6	831
19.	patience_10	72.8	71.1	71.3	72.6	71.7	71.8	72.9	71.3	71.2	831
20.	rate_10e-4	73.4	71.9	73.4	73.3	72.9	73.4	73.5	72.3	73.4	831
21.	rate_10e-3	73.2	71.1	72.7	73.0	72.0	72.6	73.4	71.4	72.6	831
22.	rate_10e-2	62.4	64.2	66.5	62.8	65.2	66.6	62.6	64.7	66.2	831

Table D.17: Architecture Tuning: Relative Connectivity

Source Code

```
source
├── data
│   ├── models
│   ├── native
│   │   ├── preloaded
│   │   ├── preprocessed
│   │   └── raw
│   ├── normalized
│   │   ├── preloaded
│   │   ├── preprocessed
│   │   └── raw
│   ├── preprocessed
│   └── raw
├── logs
├── fsleyes
└── distributed
experiments
├── subcortical
│   ├── diffusion_md
│   │   ├── native-t1
│   │   ├── native-t1t2
│   │   ├── normalized-t1
│   │   ├── normalized-t1t2
│   │   └── architecture
│   ├── diffusion_fa
│   │   ├── native-t1
│   │   ├── native-t1t2
│   │   ├── normalized-t1
│   │   ├── normalized-t1t2
│   │   └── architecture
│   ├── connection
│   │   ├── native-t1
│   │   ├── native-t1t2
│   │   ├── normalized-t1
│   │   ├── normalized-t1t2
│   │   └── architecture
│   ├── streamline
│   └── misc
└── report
    ├── progress
    └── project
```