# DARTS: DenseUnet-based Automatic Rapid Tool for brain Segmentation

**Aakash Kaku** [†][*]
ark5765@nyu.edu

**Chaitra V. Hegde** [†][*]
cvh255@nyu.edu

**Jeffrey Huang** [‡]
Jeffrey.Huang@nyulangone.org

**Sohae Chung** [‡]
Sohae.Chung@nyulangone.org

**Xiuyuan Wang** [‡]
Xiuyuan.Wang@nyulangone.org

**Matthew Young** [‡]
Matthew.Young3@nyulangone.org

**Alireza Radmanesh** [‡]
Alireza.Radmanesh@nyulangone.org

**Yvonne W. Lui** [‡][§]
Yvonne.Lui@nyulangone.org

**Narges Razavian** [†][‡][§]
Narges.Razavian@nyulangone.org

## Abstract

Quantitative, volumetric analysis of Magnetic Resonance Imaging (MRI) is a fundamental way researchers study the brain in a host of neurological conditions including normal maturation and aging. Despite the availability of open-source brain segmentation software, widespread clinical adoption of volumetric analysis has been hindered due to processing times and reliance on manual corrections. Here, we extend the use of deep learning models from proof-of-concept, as previously reported, to present a comprehensive segmentation of cortical and deep gray matter brain structures matching the standard regions of aseg+aparc included in the commonly used open-source tool, Freesurfer. The work presented here provides a real-life, rapid deep learning-based brain segmentation tool to enable clinical translation as well as research application of quantitative brain segmentation. The advantages of the presented tool include short ($\sim$ 1 minute) processing time and improved segmentation quality. This is the first study to perform quick and accurate segmentation of 102 brain regions based on the surface-based protocol (DMK protocol), widely used by experts in the field. This is also the first work to include an expert reader study to assess the quality of the segmentation obtained using a deep-learning-based model. We show the superior performance of our deep-learning-based models over the traditional segmentation tool, Freesurfer. We refer to the proposed deep learning-based tool as DARTS (DenseUnet-based Automatic Rapid Tool for brain Segmentation). Our tool and trained models are available at https://github.com/NYUMedML/DARTS

## 1 Introduction

Quantitative regional brain volumetrics have been used to study nearly every neurological, developmental and behavioral condition known, from normal aging[40, 9, 42], to schizophrenia[48] to dementia[4, 10, 39, 13, 14], to multiple sclerosis[12, 1] and hydrocephalus [11], just to name a few. Semi-automated segmentation tools have been widely applied for this task, providing measures of whole-brain and regional brain volumes [17],[25]. Segmentation is useful for multiple other research tasks such as coregistration with other imaging modalities (e.g., functional MRI, diffusion MRI, Positron emission tomography) and anatomical localization. Influential research initiatives including the Human Connectome Project [49], Alzheimer's Dementia Neuroimaging Initiative (ADNI) [41], National Alzheimer's

---

[*]Equal contribution

[†]Center for Data Science, New York University, New York, NY 10011

[‡]Department of Radiology, New York University School of Medicine, New York, NY 10016

[§]Co-corresponding authors

Coordinating Center (NACC) [2], and the UKBioBank [46] rely on traditional brain segmentation tools to provide quantitative measures to researchers.

Traditional semi-automated methods are based on Markov random fields and apply boundary determination methods to separate cortical gray matter from subcortical white matter. Despite such tools 1) being available through both open source [17] and commercial visualization products for decades, and 2) having clear potential utility, this technology has failed to translate well to routine clinical care, in part due to the need for manual corrections and off-line processing that can take hours even with modern computing capabilities. In the clinical setting, these aspects place significant practical barriers to successful implementation.

Recent innovations using deep learning for solving problems in computer vision have resulted in a revolution in medical imaging [32]. In particular, there have been novel developments using deep learning for medical imaging segmentation tasks [44, 37, 31, 7].

Previous efforts applying deep-learning-based models to a brain segmentation task [6, 45, 50, 3, 5, 15] provide proof of concept that segmentation for coarse regions of interest (ROIs) ($\sim$ up to 35 regions) is promising. The major practical limitation of these prior works is incomplete segmentation of the brain into finer anatomic regions which are typically available through traditional tools like Freesurfer. There are substantial challenges in terms of how to approach the segmentation of these finer anatomic regions, relating to the small size of these regions containing relatively few voxels and the resulting class imbalance.

Here, we extend the use of deep-learning-based models to perform segmentation of a complete set of cortical and subcortical gray matter structures and ventricular ROIs, matching the regions included in the commonly used, standard tool, Freesurfer (aseg+aparc segmentation libraries), to provide a real-life rapid brain segmentation tool. We employ a weighted loss function, weighing each ROI in inverse proportionality to its average size to address the extreme class imbalance. Additionally, we use dense convolutions in the U-net architecture and show that such architecture (called DenseUNet) provides us substantial gains over the baseline U-net model in terms of Dice Score improvement.

We assess both the quality of segmentation obtained using our deep-learning-based model and time required for segmentation compared against standard Freesurfer segmentation using both quantitative indices as well as expert evaluation. To our knowledge, this is the first report with accompanying source code of a practical tool that can be used both in a research setting to augment standard prior methods as well as in clinical settings to provide fast and accurate quantitative brain measures.

## 2 Related Work

### 2.1 Current Tools for Brain Segmentation

Many different brain segmentation tools such as Freesurfer [17], STAPLE [52] and PICSL [51] are currently used by neuroimaging researchers and radiologists. All of these tools are based on atlas registration via nonrigid registration methods, which are computationally expensive during inference. Of the tools mentioned above, Freesurfer is one of the most commonly used tools. Freesurfer is based on topological surface mappings to detect gray/white matter boundaries followed by nonlinear atlas registration and nonlinear spherical surface registration for each sub-cortical segment. Each step involves an iterative algorithm, and the surface registration is based on inference on Markov Random Fields (MRF) initially trained over manually labeled datasets[17, 43]. Despite the surface registration method being spatially non-stationary, due to 1) non-convex nature of the model, 2) subject- and pathology-specific histological factors that impacts intensity normalization, and 3) iterative process for finding optimal segmentation, Freesurfer creates different results under different initialization settings, even for the same scan. It is known that Freesurfer outputs different results if the previous scans of the patient are taken into account [43]. First released in 2006, Freesurfer has been used innumerable times by researchers, saving the need to perform complete manual segmentation of brain MRIs, which was the prior standard; however, the methodology employed by this tool and others like it suffer from some inherent limitations. Specifically, each transformation of Freesurfer on a single brain volume is computationally intensive and the time required to segment a single 3D MRI volume can be on the order of hours. Additionally, the quality of segmentation of such MRF-based models is also lower than the deep-learning-based models which are demonstrated in the reader study performed in this report. Similar limitations plague all of the other traditional tools.

### 2.2 Deep-learning-based Brain Segmentation Tools/Models

With the advent of deep learning methods for computer vision tasks like classification, object detection, and semantic segmentation, some of the inherent limitations of traditional image processing methods were resolved [30]. Consequently, these techniques were employed in several application domains including segmentation of brain cortical

structures. Researchers have approached the task of segmentation of brain both by using 2D slices [53, 36] and 3D volumes as inputs [50, 23, 5, 15]. Despite 3D models naturally utilizing 3D structural information inherent in brain anatomy, it has been shown that such models do not necessarily yield superior results [18, 36]. Additionally, they tend to be computationally more expensive and therefore slower during inference. 3D-based whole volume methods also require a pre-defined number of slices through the brain volume as input and, in practice, the number of slices varies between protocols, making such models potentially less generalizable. Researchers including [23], and [50] have attempted to address the computational cost by training on volumetric patches; however, inference time remains relatively long (DeepNAT requires ∼1-2 hours and SLANT takes ∼15 mins using multiple GPUs). [53] and [36] have performed segmentation using patches from a 2D slice through the brain, offering 3 and 10 segments respectively. But compared with the over 100 segments available via Freesurfer, these few segments limit the tools' practical utility. In order to take advantage of 3D information while keeping the computational cost low, in QuickNat [45], 2D convolutional neural networks in multiple planes have been trained and combined, but this also requires a complete 3D volume with voxel resolution being $1mm^3$. To perform such a preprocessing, Freesurfer is needed. Additionally, QuickNat only provides coarse segmentation for ∼ 30 ROIs making it less usable for clinical purposes.

These prior works like [23] and [50] clearly show the promise that deep learning models can be used for segmenting the brain into anatomic regions; however, in some prior models, the potential benefit in computation derived from using a deep-learning-based approach is negated by the need for slow pre-processing steps (e.g. registration, intensity normalization, conforming in [23] ) or post-processing steps (e.g., Conditional Random Field in [50]) that are required for these tools to operate. These steps increase the computational cost of the complete pipeline and render them slower.

In summary, our goal is to provide a tool with high accuracy, short inference time and sufficient brain segments to be useful in current research practice and clinical applications.

Contributions of this work are as follows:

- To the best of our knowledge, this is the first presentation of a truly practical, deep-learning-based brain segmentation tool that can provide accurate segmentation of over 100 brain structures, matching regions found in aseg+aparc segmentation libraries from one of the leading industry-standard, registration-based tools, Freesurfer.
- Here, we leverage the benefits of using a surface-based approach (specifically DMK protocol) for brain segmentation.
- We impose no additional, registration-based pre-processing or post-processing steps and achieve inference times of ∼1 minute using just a single GPU machine.
- We show an excellent generalization of our model to different MRI scanners and MRI volumetric acquisition protocols.
- In addition to quantitative assessments against Freesurfer segmented data, we also evaluate our model against manually segmented data and perform an expert reader study to assess the quality of our segmentation tool.

## 3 Methods

The study is conducted in compliance with the local Institutional Review Board at NYU Langone Health.

### 3.1 Data

The training data comes from the Human Connectome Project (HCP) [19]. Specifically, we used 3D image volumes from a Magnetization Prepared Rapid Acquisition Gradient Echo (MPRAGE) volumetric sequence obtained as part of the HCP protocol. These images were acquired at multiple sites at 3 Tesla (Connectome Skyra scanners; Siemens, Erlangan) with the following parameters: FOV = 224mm x 224mm, resolution = 0.7mm isotropic, TR/TE = 2400/2.14 ms, bandwidth = 210 Hz/pixel. Each MRI volume has a dimension of $256 \times 256 \times 256$. For training the model, a 2-D coronal slice of the MRI volume was used as the input. The coronal plane was selected based on superior performance in our preliminary experiments compared to axial and sagittal. The box-plots for the dice scores of these experiments can be seen in the appendix section (Figures 20, 21, 18, and 19). Freesurfer segmentation results processed for the HCP ([19]) were used as auxiliary ground truth for initial training of the models. These segmentation results had undergone manual inspection for quality control [33]. Further, to fine-tune the model, we used two additional data sources: manually segmented data by [27] and Freesurfer segmented, manually corrected NYU data.

Initially, we focused on 112 regions per individual 3D brain volume based on Freesurfer (aseg+aparc) labels. Figure 8 and figure 9 show the voxel distribution of the 112 segments in our study, revealing a class imbalance challenge: only a

| Dataset Name | Num. of Subjects | Scanners | Acquisition Details |
|---|---|---|---|
| Oasis-TRT [34] | 20 | MP-RAGE 1.5-T Vision scanner (Siemens, Erlangen, Germany) | Resolution = 1mm × 1mm × 1mm, TR/TE = 9.7/4 ms |
| Multi-modal MRI Reproducibility [29] | 21 | MP-RAGE 3T MR scanner (Achieva, Philips Healthcare, Best, The Netherlands) | Resolution = 1mm × 1mm × 1.2mm, TR/TE/TI =6.7/3.1/842ms, FOV = 240 × 204 ×256mm |
| Multi-modal MRI Reproducibility 3T/7T [29] | 2 | MP-RAGE 3T/7T MR scanner (Achieva, Philips Healthcare, Best, The Netherlands) | Resolution = 1mm × 1mm × 1.2mm, TR/TE/TI =6.7/3.1/842ms, FOV = 240 × 204 ×256mm |
| Nathan Kline Institute/Rockland Sample | 22 | MP-RAGE MR Siemens Magnetom | Resolution = 1mm × 1mm × 1mm, TR/TE =1900/2.52, FOV = 250 × 250 ×250mm, Bandwidht = 170 Hz/Px |
| Nathan Kline Institute/Test-Retest | 20 | MP-RAGE MR Siemens Magnetom | Resolution = 1mm × 1mm × 1mm, TR/TE =1900/2.52, FOV = 250 × 250 ×250mm, Bandwidht = 170 Hz/Px |
| Human Language Networks [38] | 12 | MP-RAGE 3T MRI scanner (Philips Medical Systems, Best, Netherlands) | FOV = 240 mm, TE = 35 ms, TR = 2 sec, 4.5 mm thickness |
| Colin Holmes Template [21] | 1 | MP-RAGE Phillips 1.5 T MR | Resolution = 1mm × 1mm × 1mm, TR/TE = 18/10 ms, FOV = 256 mm (SI) × 204 mm (AP) |
| Twins-2 [27] | 2 | MP-RAGE | |
| Afterthought-1 [27] | 1 | MP-RAGE | |

Table 1: MRI acquisition details of all the data in Mindboggle-101 dataset

few regions are large (>60000 voxels) whereas most of the regions are significantly smaller (<20000 voxels). The class imbalance challenge is addressed in section 3.6. The following 10 regions were excluded from the analysis: segments labeled left and right 'unknown', four brain regions not common to a normal brain: White matter and non-white matter hypointensities and left and right frontal and temporal poles, and segments without widely accepted definitions in the neuroradiology community (left and right bankssts) [27].

We randomly divided the cohort into training, validation and (held out) test sets with 60% (667 scans), 20% (222 scans) and 20% (222 scans) ratio. We separated training, validation and held-out test sets according to patients rather than slices to prevent data leakage.

### 3.2 Manually Annotated Dataset

We used Mindboggle-101 [27] manually annotated data to fine-tune and further evaluate our model. The dataset includes manual annotations for all the segments that are present in the Freesurfer aseg+aparc list (except those listed in the previous section). The Mindboggle-101 dataset contains MRI scans of normal individuals from a diverse number and type of MRI scanners ranging in magnetic field strength from 1.5 Tesla to 7 Tesla. Mindboggle-101 contains 101 MRI scans from 101 subjects from multiple open-source studies. The details of the MRI acquisition for each dataset can be found in table 1. These data were also randomly split into training, validation and (held out) test set with the same 60%(60 scans), 20%(21 scans) and 20%(20 scans) ratio, again separated according to patients rather slices to prevent data leakage. The subjects' ages range from 20 to 61 years.

### 3.3 NYU's Manually Corrected Dataset

We also use a small internal NYU dataset consisting of 11 patients to train and assess the generalizability of the segmentation model. The description of the NYU dataset is as follows: MPRAGE (FOV=256 × 256 mm$^2$; resolution=1 × 1 × 1 mm$^3$; matrix=256 × 256; sections, 192; TR=2100 ms;TE=3.19 ms; TI=900 ms; bandwidth=260 Hz/pixel). Imaging was performed on 3T Siemens Skyra and Prisma MRI scanner. Each MRI scan in this dataset was first processed using the standard Freesurfer tool and then underwent manual corrections by an expert neuroimager for ground truth segmentation. Here also, we split the data into the train (6 scans), validation (2 scans) and held-out test (3 scans) sets.
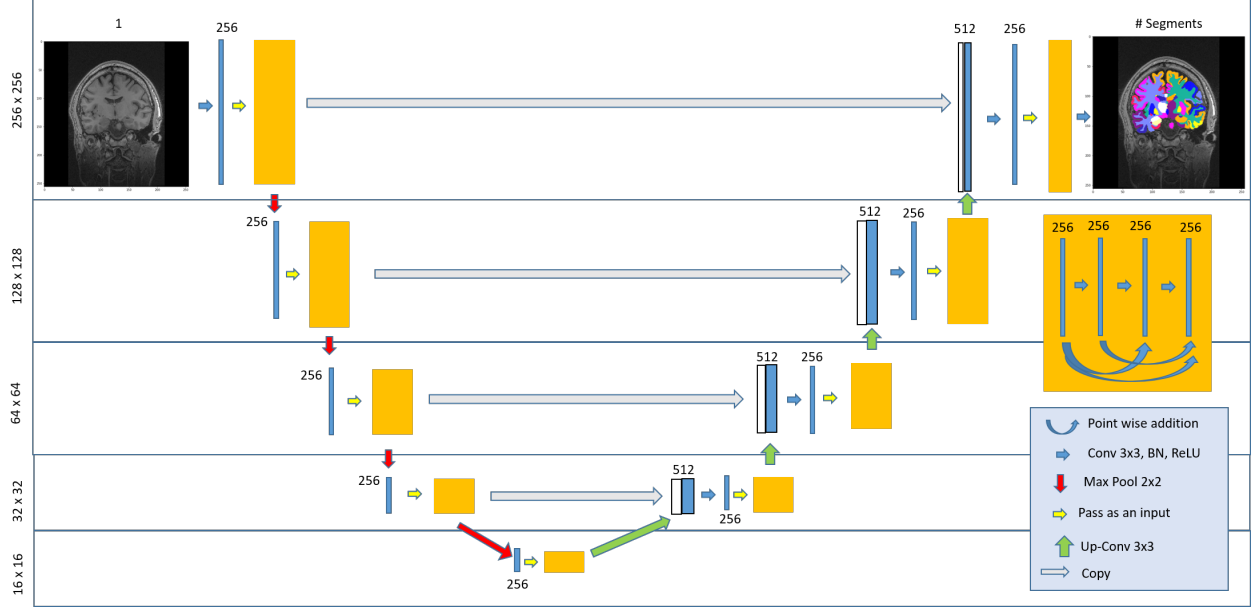
Figure 1: Schematic Diagram of DenseUNet

## 3.4 Data Augmentation

Differences across MRI scanners and acquisition parameters result in differences in image characteristics such as signal-to-noise, contrast, the sharpness of the acquired volume. In addition, there are between-subject differences arising from subject positioning, etc. To improve the generalizability of our models to a broad array of scanner parameters, one or more of the following augmentation methods was applied to the training data at a random 50% of the training batches: First, gaussian blurring (sigma parameter uniformly drawn from [0.65 to 1.0]) and gamma adjustment (gamma parameter uniformly drawn from [1.6 to 2.4]). Second, input MRI and corresponding target labels were rotated by angle theta where theta is a random number uniformly drawn from [-10 to +10] degrees. Third, input and the corresponding labels were shifted up and sideways by dx and dy pixels, where dx and dy were randomly and uniformly drawn between +25 and -25. The input was then normalized between 0 and 1 using min-max normalization. Data augmentation was used only during training.

## 3.5 Deep Learning Model

The model architecture used for performing the segmentation task is an encoder-decoder style fully convolutional neural network. The architecture is partly inspired by U-Net architecture [44] and partly by [24] and [45]. We term the new architecture as DenseUNet. The DenseUNet has a U-Net like architecture where it has four dense blocks in the encoding path and four dense blocks in the decoding path. The encoding pathway is connected to the decoding pathway using another dense block called the connecting dense block. Each of the encoding dense blocks is connected to the corresponding decoding dense block using skip connections. These skip connections facilitate the gradient flow in the network as well as the transfer of spatial information required for decoding the encoded input. The output from the last decoding dense block is passed through a classification block where the network performs the classification and gives the final output in the form of a separate probability map for each brain structure. The schematic diagram of the entire architecture can be seen in figure 1.

We also implemented a vanilla U-Net architecture as described in the original paper [44] which serves as a baseline model. The schematic diagram for the same can be seen in figure 7.

The architectural choice of DenseUNet was also motivated by our empirical results on U-Net. A Dense block has more model capacity compared to standard convolutional block [22]. When larger training data is available, DenseUNet has larger learning and generalization capabilities.

All the parameters of the U-net and DenseUNet were initialized using Xavier initialization [20].

The components of DenseUNet i.e. the encoding dense block, the connecting dense block, the decoding dense block, and the classification block are explained below.

5

### 3.5.1 Encoding Dense Block

The encoding dense block has four convolution layers. The output of all the previous convolution layers is added to the subsequent convolution layers. These additive connections are termed dense connections, that facilitate the gradient flow and allow the network to learn a better representation of the image [22]. The other common connection type is a concatenated connection. QuickNAT [45] uses concatenate connections to build dense blocks. Though concatenated connections can model additive connections, the model complexity in terms of the number of parameters and number of mathematical operations increases significantly leading to out-of-memory issues while training the model for a large number of segments. Therefore, to avoid out-of-memory issues and to achieve low training and inference times, we choose to use additive dense connections as opposed to concatenating dense connections. The output obtained by adding all the previous convolution layers' output is followed by batch normalization and Rectifier Linear Unit (ReLU) non-linearity. Each convolution layer has 256 output channels with the filter size being $3 \times 3$ for all the channels. The output of the encoding dense block is followed by a max-pooling layer with a kernel size of $2 \times 2$ and a stride of 2. The down-sampled output is fed to the next encoding or connecting dense block.

### 3.5.2 Connecting Dense Block

The connecting dense block is similar to the encoding dense block and has four convolution layers with dense connections. The only difference is the output of the dense block is not followed by a downsampling layer like a max-pooling layer.

### 3.5.3 Decoding Dense Block

The decoding dense block is preceded by an upsampling block. The output from the previous decoding or connecting dense block is upsampled using transposed convolution with a filter size of $4 \times 4$ and stride of 2 and padding of 1. The upsampled output is concatenated with the output from the corresponding encoding dense block. The concatenated output serves as an input to a convolution layer which is followed by batch normalization and ReLU activation. The convolution layer has 256 output channels with a filter size of $3 \times 3$.

### 3.5.4 Classification Block

The classification block is a single convolution layer with the number of output channels equal to the number of brain structures to segment (in our case 112 [1]) with a filter size of $1 \times 1$. The output of the convolution layer is passed through a softmax layer to obtain the probability maps for each of the brain structures we are trying to segment.

### 3.6 Loss Function

We model the segmentation task as a multi-class classification problem. Here, since we have 112 [1] tissues of interest, this is a 113-class classification problem, where the last class is "background".

Since the dataset is an imbalanced one, we use a weighted cross-entropy loss and weighted dice loss for our task. The weighted cross entropy loss is defined as:

$$\text{Weighted-CEL} = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{S} w_j (y_{ij} log(p_{ij})) \tag{1}$$

Where $w_j$ = Weight of the jth segment and S is the total number of segmentation classes. Here, $w_j = \frac{\text{median freq}}{\text{freq(j)}}$, where, freq(j) is the number of pixels of class $j$ divided by the total number of pixels of images where $j$ is present, and median freq is the median of these frequencies [16]. N = number of pixels in a 2D MRI image (a slice of the MRI volume), $p_{ij}$ = probability of pixel i to be belonging to segment j, $y_{ij}$ = label of pixel i to be belonging to segment j = 1 or 0.)

Weighted Dice Loss: The weighted dice score is defined as:

$$\text{Weighted-Dice Loss} = 1 - 2 \frac{\sum_{j=1}^{S} w_j \sum_{i=1}^{N} y_{ij} p_{ij}}{\sum_{j=1}^{S} w_j \sum_{i=1}^{N} y_{ij} + p_{ij}} \tag{2}$$

Weights are calculated using the same approach as mentioned for weighted-CEL.

---

[1]Our architecture trains for full aparc+aseg segments i.e. 112 segments, and we omit the 10 excluded segments mentioned in section 3.1 after training.

We experiment with local(using only a specific batch) and global (using the entire train set) estimation of $w_j$. In the local case, $w_j$s were adapted for each batch, and hence the loss function for each batch was slightly different. However, we found that a global $w_j$ gave us the best out of sample performance.

We combine the above two losses in a novel way that we term as loss switching. Loss switching is explained in the next section.

### 3.7 Loss Switching

In segmentation tasks, the dice score is often reported as the performance metric. A loss function that directly correlates with the dice score is the weighted dice loss [47]. Based on our empirical observation, the network trained with only weighted dice loss was unable to escape local optimum and did not converge. Also, empirically it was seen that the stability of the model, in terms of convergence, decreased as the number of classes and class imbalance increased. We found that weighted cross-entropy loss, on the other hand, did not get stuck in any local optima and learned reasonably good segmentations. As the model's performance with regard to dice score flattened out, we switched from weighted cross entropy to weighted dice loss, after which the model's performance further increased by 3-4 % in terms of average dice score. This loss switching mechanism, therefore, is found to be useful to further improve the performance of the model.

### 3.8 Evaluation Metric

Dice score (DSC) is employed here as a primary measure of quality of the segmented images. This is a standard measure against which segmentations are judged and provide direct information on similarity against the ground truth labels.

$$\text{DSC} = \frac{2||PT||_2^2}{||P||_2^2 + ||T||_2^2} \tag{3}$$

where P = Predicted Binary Mask, T = True Binary Mask, PT = element-wise product of P and T, $||X||_2$ is the L-2 norm. Dice score can equivalently be interpreted as the ratio of the cardinality of $(T \cap P)$ with the cardinality of $((T \cup P) + (T \cap P))$.

From the above definition, it can be seen that DSC penalizes both over prediction and under prediction and, hence, is well-suited for segmentation tasks such as the one proposed here and is particularly useful in medical imaging.

### 3.9 Training Procedure

What follows is a description of the training procedure we employed:

- For training each model, we used Adam Optimizer with reducing the learning rate by a factor of 2 after every 10-15 epochs.
- The model was initially trained on the training set of scans from the HCP dataset with the auxiliary labels until convergence. The trained model was then finetuned using the training set of the manually annotated dataset (Mindboggle-101) and the training set of the in-house NYU dataset.
- We trained the model with the HCP dataset using the loss switching procedure described in section 3.7 whereas, for finetuning, the loss function is simply the weighted dice loss as described in section 3.6.
- All the models are trained using early stopping based on the best dice score on the validation set.

## 4 Methodology for Reader Study

### 4.1 Reader Study: Description and Setup

We perform an expert reader evaluation to measure and compare the deep learning models' performance with the Freesurfer model. We use HCP held-out test set scans for reader study. On these scans, Freesurfer results have undergone a manual quality control[33]. We also compare the non-finetuned and fine-tuned model with the Freesurfer model with manual QC. Seven regions of interest (ROIs) were selected: L/R Putamen (axial view), L/R Pallidum (axial view), L/R Caudate (axial view), L/R Thalamus (axial view), L/R Lateral Ventricles (axial view), L/R Insula (axial view) and L/R Cingulate Gyrus (sagittal view). The basal ganglia and thalamus were selected due to their highly interconnected nature with the remainder of the brain, their involvement in many neurological pathologies, and their ease of distinction from surrounding structures. The insular and cingulate gyri were selected to assess the quality of
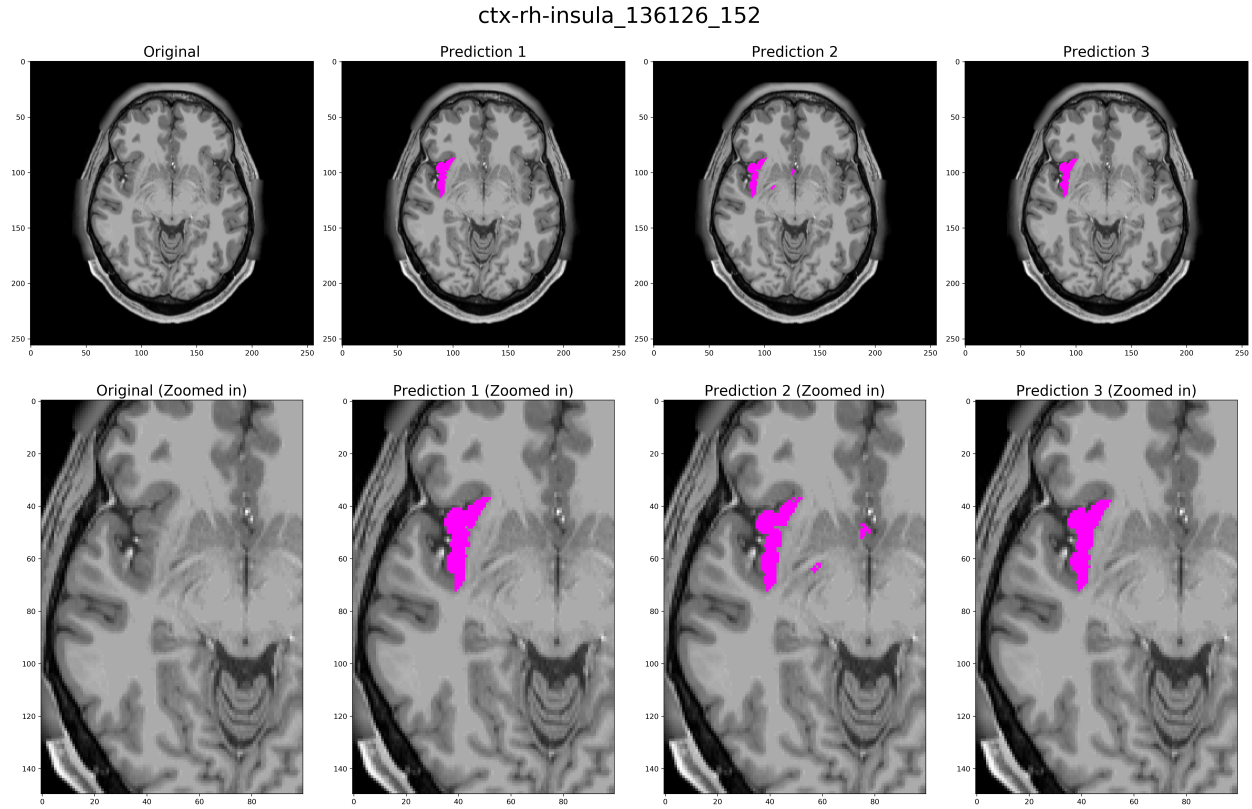
Figure 2: A sample segmentation of Right Insula used in the expert reader evaluation. Here, predictions 1, 2 and 3 are from the Finetuned model, Freesurfer, and non-Finetuned model respectively; though in the reader study, numbering and order of presentation of the predictions are randomized and a total of 280 of examples are presented. Each reader is asked to rate each example for the quality of segmentation on a 5-point Likert-type scale.

cortical segmentation in structures best visualized in different planes and also due to the relatively frequent involvement of the insular gyrus in the middle cerebral artery infarctions. The lateral ventricles were selected to assess for quality of segmentation of cerebrospinal fluid structures, which would help identify pathologies affecting cerebrospinal fluid volumes, including intracranial hypotension, hydrocephalus, and cerebral atrophy.

Three expert readers performed visual inspection and assessment of the segmentation results. There were two attending neuroradiologists with 3 and 5 years of experience and one second-year radiology resident. Each reader was asked to rate 40 different predictions for each ROI (20 in each brain hemisphere) such as shown in figure 2. Readers were blinded to the algorithm used to predict the segmentation and examples were presented in a randomized order. Each prediction presented consisted of a single slice containing a minimum of 40 pixels within the ROI, ensuring that enough of the structure being assessed was present on the given image. A sample slice rated by the readers is shown in figure 2.

Each reader rated each example on a Likert-type scale from 1 to 5 with the following definitions:

1. Rating of 1 (Poor): Segmentation has major error(s) in either the major boundary or labeling of areas outside the area of interest. Such segmentation would not yield acceptable quantitative measures.

2. Rating of 2 (Fair): Segmentation has >2 areas of error along the major boundary or includes small but nontrivial areas outside the area of interest that would require manual correction before yielding reasonable quantitative measures to be used in research or clinical scenarios.

3. Rating of 3 (Good): Segmentation has 1 or 2 area(s) of error along the major boundary that would require manual correction before yielding reasonable quantitative measures to be used in research or clinical scenarios. A good segmentation could have minor/few voxels separate from the volume of interest.

4. Rating of 4 (Very Good): Segmentation has minor, small areas of error along the major boundary that would still yield reasonable quantitative anatomic measures without manual corrections, appropriate for research or clinical use.

5. Rating of 5 (Excellent): Segmentation is essentially ideal, has no (or only tiny inconsequential/questionable areas) where manual corrections might be used and would yield highly accurate quantitative anatomic measures, etc. Should have no erroneous areas outside the segment of interest.

## 4.2 Reader Study: Analysis

Using the ratings obtained from three readers the following analyses are performed:

1. Inter-Reader Reliability (IRR): An IRR analysis is performed using a two-way mixed, consistency, average measures ICC (Inter Class Correlation) [35] to assess the degree of agreement between readers. High ICC indicates that the measurement error introduced due to independent readers is low and hence, the subsequent analysis' statistical power is not substantially reduced.

2. Comparison of different models: Based on the readers' ratings, we investigate if there are statistically significant differences between the three methods using paired T-test and Wilcoxon signed-rank test at 95% significance level.

# 5 Results

Here, we report our quantitative evaluation results on the held-out test sets from manually annotated and corrected Mindboggle-101 and NYU dataset. We also report results of a qualitative evaluation via a reader study with expert neuroimaging radiologists on held out HCP scans and their corresponding Freesurfer labels and model's prediction.

## 5.1 Quantitative Evaluation: Performance on the Manually annotated Dataset - Mindboggle-101

Table 2 includes the performance of the Finetuned model and the non-Finetuned model on the manually annotated test set from Mindboggle-101 data.

| Model Name | non - Finetuned | Finetuned |
|---|---|---|
| UNet (Baseline) | 0.7329±0.014 | 0.80±0.013 |
| DenseUNet | 0.7431±0.015 | **0.819±0.011** |

Table 2: Mean Dice Score on 102 segments on Mindboggle-101 dataset. Here, non-Finetuned = Model trained using only HCP dataset, Finetuned = Model initially trained on HCP dataset and subsequently finetuned using Mindboggle-101 and NYU Dataset

Detailed dice scores for the 102 segments are included in figures 3 and 4. More comprehensive comparison of fine-tuned vs non-fine-tuned results are also included in the Appendix (Figure 12, figure 13).

## 5.2 Quantitative Evaluation: Performance on an External NYU Dataset

Performance of the Finetuned model and non-Finetuned model on the external manually corrected NYU Test dataset is presented in the table 3.

| Model Name | non - Finetuned | Finetuned |
|---|---|---|
| UNet (Baseline) | 0.787±0.01 | 0.785±0.014 |
| DenseUNet | 0.795±0.013 | **0.800±0.012** |

Table 3: Mean Dice Score on 102 segments on NYU-Dataset dataset. Here, non-Finetuned = Model trained using only HCP dataset, Finetuned = Model initially trained on HCP dataset and subsequently finetuned using Mindboggle-101 and NYU Dataset

For detailed results of each segment (with a box plot), please refer to the figure 10 and figure 11 in the Appendix.

## 5.3 Quantitative Evaluation: Time

Time to perform segmentation, and comparison with other models can be seen in the table 4. Our model gives the most informative segmentation (102 segments) in the least amount of time ($\sim$ 1 min).
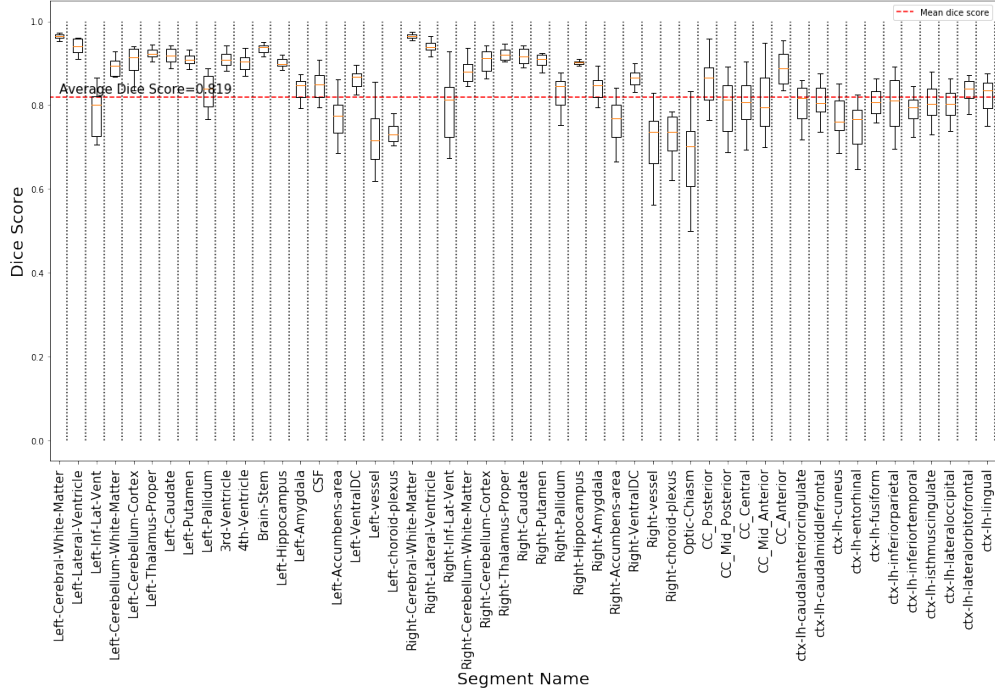
Figure 3: Box plot showing the dice scores for the first 51 regions of interest of Mindboggle dataset

| Model Name | Num. of Segments | Time (for one brain Scan) |
|---|---|---|
| DenseUNet (our Model) | 102 | 65.6 secs (± 0.9353 secs) (Single GPU Machine) |
| Unet (our Baseline Model) | 102 | 54.2 secs (± 0.7908 secs) |
| FreeSurfer [45] | ∼190 | ∼4 hrs |
| PISCL [45] | ∼190 | ∼30 hrs |
| DeepNAT [45] | 27 | ∼1 hr (on a Multi-GPU Machine) |
| QuickNAT [45] | 27 | ∼20 secs (on a Multi-GPU Machine) |
| SLANT [23] | 133 | ∼15 mins (on a Multi-GPU machine) |

Table 4: Time to perform segmentation for a single MRI scan

## 5.4 Qualitative Evaluation: Reader study

### 5.4.1 Inter Reader Reliability

As mentioned in the section 4, IRR is assessed using a two-way mixed, consistency, average-measures ICC (Inter Class Correlation). The result can be seen in table 5.

| Model | Inter Class Correlation (Lower Bound - Upper Bound at alpha = 95%) |
|---|---|
| Freesurfer | 0.66 (0.58 - 0.72) |
| non-Finetuned model | 0.82 (0.78 - 0.85) |
| Finetuned model | 0.80 (0.76 - 0.84) |

Table 5: Inter Class Correlation between expert readers in assessment of quality of segmentations shows excellent agreement between readers for our model.

Figure 4: Box plot showing the dice scores for the next 51 regions of interest of Mindboggle dataset

Since the resulting IRR is in good and excellent ranges [8] for all the ratings, the statistical power and significance of the paired T-test and Wilcoxon signed-rank test to compare different model's performance would be reliable.

### 5.4.2 Comparison of Different Models with Freesufer

The table 6 showcases the average ratings obtained by each model (including Freesurfer). It also showcases whether the difference between any two model is statistically significant or not (at 95% significance level) using paired T-test and Wilcoxon signed rank test.

## 6 Discussion

### 6.1 First fast and accurate segmentation for 102 regions of interest, consistent with aseg+aparc

This work presents a deep-learning based model that performs extensive segmentation of the brain into 102 ROIs, matching the regions provided by the aseg+aparc segmentation libraries of the tool Freesurfer with inference time of ∼1 minute for a single 3D brain volume. Freesurfer is a widely used tool by neuroimaging researchers as well as the basis for many industry products, and thus the regions labelled under aseg+aparc are useful for current research and clinical applications.

In our reader study, quality of segmentations of Freesurfer as well as our developed models are evaluated by neuroradiologists. We find that the deep learning model trained on Freesurfer labels leads to significantly improved segmentation quality in 13 out of 14 regions assessed. We note that compared with MRF-based segmentation which forms the basis for Freesurfer, our deep learning model improves the quality of segmentation by producing smoother boundaries that follow the anatomic border more closely. Finetuning on manually annotated Mindboggle data further improves quality of segmentations for Insula and Pallidum ROIs. Interestingly, these are the areas where Freesurfer has the most reported boundary errors.

Although the Mindboggle dataset is useful for improving performance for most of the ROIs, there are regions such as the Cingulate gyrus that are better segmented by the non-finetuned model.

Based on our findings, optimal finetuning of the model would be to only finetune for regions that are not already well segmented, using the manually annotated dataset. This is part of our future explorations.

| Region of Interest | Mean Rating ± One standard deviation | | | Statistically Significant Difference (s) (as per paired T-test and Wilcoxon test) |
|---|---|---|---|---|
| | **FS** | **NFT** | **FT** | |
| Insula | L-3.13±0.88<br>R-2.85±1.00 | L-3.90±1.01<br>R-3.48±1.26 | **L-4.23±0.83**<br>**R-4.21±0.83** | L/R - FS and NFT, NFT and FT, FS and FT |
| Caudate | L- 4.26±0.79<br>R- 3.97±0.75 | **L- 4.46±0.70**<br>R- 4.17±0.76 | L- 4.45±0.75<br>**R- 4.26±0.67** | L/R - FS and NFT, FS and FT |
| Cingulate-Gyrus | L- 2.59±0.76<br>R- 2.72±0.74 | **L- 2.91±0.97**<br>**R- 3.0±0.88** | L- 2.53±1.05<br>R- 2.49±0.89 | L - FS and NFT, NFT and FT<br>R - FS and FT, FS and NFT, NFT and FT |
| Lateral-Ventricles | L- 4.14±0.83<br>R- 4.16±0.73 | **L- 4.46±0.66**<br>R- 4.39±0.73 | L- 4.44±0.73<br>**R- 4.41±0.72** | L - FS and FT, FS and NFT<br>R - FS and FT, FS and NFT |
| Pallidum | L- 3.20±0.80<br>R- 3.07±0.65 | L- 3.30±0.79<br>R- 3.72±0.78 | **L- 3.87±0.85**<br>**R- 3.93±0.78** | L - NFT and FT, FS and FT<br>R - FS and NFT, NFT and FT, FS and FT |
| Putamen | L- 3.22±2.14<br>R- 3.10±1.04 | L- 3.16±1.11<br>**R- 3.44±1.13** | **L- 3.23±1.16**<br>R- 3.19±1.13 | L - No difference is statistically significant (as per paired T-test)<br>L - FS and NFT, FS and FT (as per Wilcoxon test)<br>R - FS and NFT, NFT and FT |
| Thalamus | L- 3.40±0.75<br>R- 3.28±0.83 | **L- 4.0±0.88**<br>R- 3.96±0.82 | L- 3.95±0.94<br>**R- 4.02±0.87** | L - FS and NFT, FS and FT<br>R - FS and NFT, FS and FT |
| All Regions | 3.41 ± 1.12 | 3.78±1.04 | **3.86±1.07** | FS and NFT, NFT and FT, FS and FT |

Table 6: Reader study results comparing Freesurfer(FS), Non-fine-tuned(NFT) and Fine-tuned(FT) models' segmentation on a total of 20 evaluations per ROI, per Left(L) and Right(R) hemisphere, per model. Statistical tests are performed at 95% significance level.

## 6.2 Model Generalizability and Usability

The model's generalizability was tested using a held out dataset from the Mindboggle dataset as well as a held out dataset from NYU. Our model achieves a mean dice score of 0.819 and 0.80 on the two datasets, respectively, showcasing good generalization capabilities on the unseen data.

Additionally, if we see the detailed box plot shown in figure 3 and 4, we notice that the dice scores for few regions such as Optic-Chiasm and Right-vessel are relatively lower than the dice scores for other regions. We investigated such regions and found that the regions with low dice scores are the regions with low mean voxel count (see figure 22 in appendix). This is likely due to the fact that small errors in segmentation represents as a higher percentage of a small sized segment. Hence, overall dice scores for smaller regions tend to be lower than the dice scores for larger regions.

The use of 2D images as inputs also adds to the practical usability of this model. Often in real-life clinical and research practices, incomplete 3D image volumes may be encountered. In such cases, performing segmentation using our tool remains straightforward, distinct from most currently available semi-automated tools and 3D deep-learning-based models that require entire 3D brain volumes with pre-specified number of slices as inputs.

Another advantage of the proposed the model is relatively short inference time. As seen in table 4, the model provides fairly comprehensive segmentation in $\sim$ 1 minute, comparing favorably against other available tools. This, combined with lack of dependency on pre-processing, makes our tool feasible on-demand in a clinical setting. Moreover, this work adheres to Freesurfer's region-naming convention, already in wide use and familiar to neuroimagers.

### 6.2.1 Errors of Freesurfer Segmentation: Putamen, Insula and Pallidum

The work of [26] reveals errors in Freesurfer segmentation, thus necessitating manual corrections for quality control. Our reader study results also demonstrates the low quality nature of Freesurfer segmentations and confirms the findings of [26]. As per our readers, the highest difference in the quality of segmentation between the Freesurfer segmentation and Finetuned model segmentation is seen for those regions which showed major boundary errors in the Freesurfer segmentation as demonstrated in figure 5. Those regions are Insula and Pallidum.

In the Freesurfer segmentation, we see boundary errors, inclusion of discontiguous areas, as well as stair-step artifacts along the boundary that render a noisy and non-natural-appearing result. One potential benefit of a deep-learning-based brain segmentation tool over the traditional MRF-based tool is that by training over multiple examples, the model learns that jagged or stair-step boundaries are not consistent, and can not be explained by naturally visible MRI boundaries. The model therefore simply fails to learn the arbitrary jaggedness. For additional examples of Freesurfer vs DenseUNet putamen and pallidum segmentations please refer to figure 16 and figure 17 in the appendix.
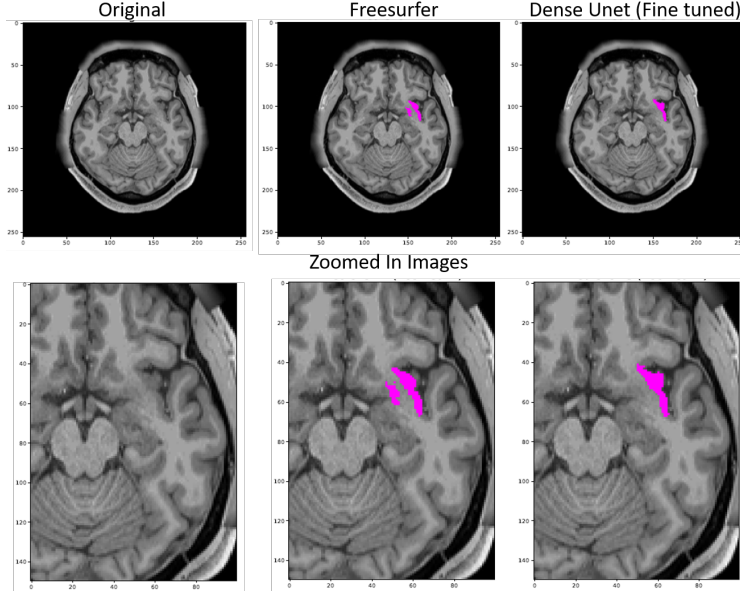
Figure 5: MR image in the axial plane through the level of the insula (left), Freesurfer (FS) (center) and our proposed model, DenseUNet (right) prediction of the insula. It is evident that the FS prediction has errors in both over and underestimation along the boundary, includes discontiguous voxels and is also somewhat non-natural looking with stair-step artifact and a noisy appearance. DenseUNet insula segmentation, on the other hand, obeys well the segment boundaries and lacks the stair-step artifact, rendering a smooth contoured, more natural appearing segmentation.

### 6.3 Limited high quality manual data for quantitative evaluation and finetuning

As we see in section 6.2.1 and from the results in table 6, the Freesurfer generated segmentation labels have low quality and hence, the model trained only on such labels would also be susceptible to generate low quality segmentation. A manually annotated high quality segmentation, therefore, becomes an important resource for training/finetuning the model and quantatively evaluating its performance.

There are two major open source datasets which contain manually annotated complete brain segmentations. Most prior work in the area of developing deep-learning-based brain segmentation models use the 2012 MICCAI Multi-Atlas labelling challenge dataset [28]. The ground truth labels of this dataset when visualized, show corruption in sagittal and axial views. This is due to the manual segmentation being performed only in coronal view without correction in other planes. These artifacts, visualized in the figure 6, introduce potential unwanted biases in models that train on these data.

Because of these problems with the manual segmentation ground truth observed in the 2012 MICCAI Multi-Atlas labelling challenge dataset, we use the Mindboggle 101 dataset which contains 102, manually segmented regions in each brain MRI. Mindboggle 101 has the benefit of containing image data from a diverse set of scanners and sites (discussed in the dataset section) which gives our model the opportunity to learn the invariances across the image protocols and scanners, making a more robust and generalizable model.



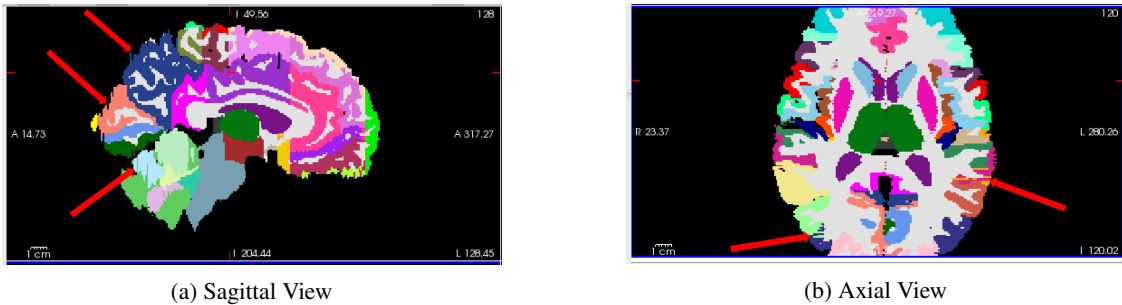(a) Sagittal View



(b) Axial View

Figure 6: Image showcasing artifacts induced by performing manual correction in one view (here coronal) and extending to other views (here sagittal and axial)

### 6.4 Limitations of this study

Limitations of this study include:

1. At the current time, the model offers segmentation of gray matter sub-structures (aparc) and does not include white matter sub-structure segmentation (wmparc), though this is an ongoing part of our future work.

2. Some regions such as white matter hyperintensities, which are not commonly represented in the HCP dataset, were excluded.

3. In the neuroimaging community, there are different approaches to segmentation. This work adheres closely to the DKT protocol followed by Mindboggle and Freesurfer [27]; other segmentation protocols would require re-training of the model.

4. More extensive reader study involving the full 102 regions of interest was not present in this work and is deferred to future studies.

## 7 Future Work and Conclusion

We have a number of active directions currently under investigation: Firstly,as discussed in section 6.4, we have not validated white matter regional segmentation, which is underway. Secondly, lately, there is much interest to learn from noisy labels. Since the first step of this model involves learning from auxiliary and somewhat noisy Freesurfer labels, employing additional methods to more efficiently learn from noisy labels could be explored. Lastly, we are exploring ways of selectively fine tuning the model only for those regions with high quality manual segmentation available.

In conclusion, we present a deep-learning based model that performs an extensive, anatomical brain segmentation yielding 102 brain regions that match a commonly used tool, Freesurfer, in one minute (for a single GPU machine) for a single brain volume. This fast and accurate open-source segmentation tool can finally make on-demand clinical utilization of brain segmentation feasible, enabling translation of wealth of neurological research into clinic. Our proposed model does not need a complete 3D Brain MR for performing the segmentation. Even a single slice of Brain T1 MRI could be segmented using our proposed model, and the model operates without any additional pre- or post-processing steps. We demonstrate successful generalization of our model to a variety of scanner devices and resolutions. Finally, we also performed first reader study to evaluate the segmentation quality of Freesurfer and proposed model and show that the proposed model's segmentations are of superior quality compared to Freesurfer.

## 8 Open Source Tool

Training and inference code, a jupyter notebook for full instructions, and pre-trained models including Coronal, Axial and Sagittal finetuned and non-finetuned models can be accessed from `https://github.com/NYUMedML/DARTS`.

## References

[1] Christina J Azevedo, Steven Y Cen, Amir Jaberzadeh, Ling Zheng, Stephen L Hauser, and Daniel Pelletier. Contribution of normal aging to brain atrophy in ms. *Neurology-Neuroimmunology Neuroinflammation*, 6(6):e616, 2019.

[2] Duane L Beekly, Erin M Ramos, William W Lee, Woodrow D Deitrich, Mary E Jacka, Joylee Wu, Janene L Hubbard, Thomas D Koepsell, John C Morris, Walter A Kukull, and et al. The national alzheimer's coordinating center (nacc) database: the uniform data set, 2007.

[3] José Bernal, Kaisar Kushibar, Mariano Cabezas, Sergi Valverde, Arnau Oliver, and Xavier Lladó. Quantitative analysis of patch-based fully convolutional neural networks for tissue segmentation on brain magnetic resonance imaging. *CoRR*, abs/1801.06457, 2018.

[4] Esther E Bron, Marion Smits, Wiesje M Van Der Flier, Hugo Vrenken, Frederik Barkhof, Philip Scheltens, Janne M Papma, Rebecca ME Steketee, Carolina Méndez Orellana, Rozanna Meijboom, et al. Standardized evaluation of algorithms for computer-aided diagnosis of dementia based on structural mri: the caddementia challenge. *NeuroImage*, 111:562–579, 2015.

[5] Toan Duc Bui, Jitae Shin, and Taesup Moon. 3d densely convolutional networks for volumetric segmentation. *CoRR*, abs/1709.03199, 2017.

[6] Hao Chen, Qi Dou, Lequan Yu, and Pheng-Ann Heng. Voxresnet: Deep voxelwise residual networks for volumetric brain segmentation. *CoRR*, abs/1608.05895, 2016.

[7] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *CoRR*, abs/1706.05587, 2017.

[8] Domenic V Cicchetti. Guidelines, criteria, and rules of thumb for evaluating normed and standardized assessment instruments in psychology. *Psychological assessment*, 6(4):284, 1994.

[9] James H Cole, Robert Leech, David J Sharp, and Alzheimer's Disease Neuroimaging Initiative. Prediction of brain age suggests accelerated atrophy after traumatic brain injury. *Annals of neurology*, 77(4):571–581, 2015.

[10] Pierrick Coupé, José Vicente Manjón, Enrique Lanuza, and Gwenaelle Catheline. Lifespan changes of the human brain in alzheimer's disease. *Scientific reports*, 9(1):3998, 2019.

[11] Benito Pereira Damasceno. Neuroimaging in normal pressure hydrocephalus. *Dementia & neuropsychologia*, 9(4):350–355, 2015.

[12] Nicola De Stefano, Laura Airas, Nikolaos Grigoriadis, Heinrich P Mattle, Jonathan O'Riordan, Celia Oreja-Guevara, Finn Sellebjerg, Bruno Stankoff, Agata Walczak, Heinz Wiendl, et al. Clinical relevance of brain volume measures in multiple sclerosis. *CNS drugs*, 28(2):147–156, 2014.

[13] Davangere P Devanand, Xinhua Liu, Matthias H Tabert, Gnanavalli Pradhaban, Katrina Cuasay, Karen Bell, Mony J de Leon, Richard L Doty, Yaakov Stern, and Gregory H Pelton. Combining early markers strongly predicts conversion from mild cognitive impairment to alzheimer's disease. *Biological psychiatry*, 64(10):871–879, 2008.

[14] DP Devanand, G Pradhaban, X Liu, A Khandji, S De Santi, S Segal, H Rusinek, GH Pelton, LS Honig, R Mayeux, et al. Hippocampal and entorhinal atrophy in mild cognitive impairment: prediction of alzheimer disease. *Neurology*, 68(11):828–836, 2007.

[15] Jose Dolz, Christian Desrosiers, and Ismail Ben Ayed. 3d fully convolutional networks for subcortical segmentation in MRI: A large-scale study. *CoRR*, abs/1612.03925, 2016.

[16] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE international conference on computer vision*, pages 2650–2658, 2015.

[17] Bruce Fischl. Freesurfer. *Neuroimage*, 62(2):774–781, 2012.

[18] Mohsen Ghafoorian, Nico Karssemeijer, Tom Heskes, Inge van Uden, Clara I. Sánchez, Geert J. S. Litjens, Frank-Erik de Leeuw, Bram van Ginneken, Elena Marchiori, and Bram Platel. Location sensitive deep convolutional neural networks for segmentation of white matter hyperintensities. *CoRR*, abs/1610.04834, 2016.

[19] Matthew F. Glasser, Stamatios N. Sotiropoulos, J. Anthony Wilson, Timothy S. Coalson, Bruce Fischl, Jesper L. Andersson, Junqian Xu, Saad Jbabdi, Matthew Webster, Jonathan R. Polimeni, David C. Van Essen, and Mark Jenkinson. The minimal preprocessing pipelines for the human connectome project. *NeuroImage*, 80:105–124, oct 2013.

[20] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In Yee Whye Teh and Mike Titterington, editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR.

[21] Colin J. Holmes, Rick Hoge, Roger P. Woods, Alan C. Evans, and Arthur W. Toga. Enhancement of t2 and proton density mr images using registration for signal averaging. *NeuroImage*, 3(3), 1996.

[22] Gao Huang, Zhuang Liu, and Kilian Q. Weinberger. Densely connected convolutional networks. *CoRR*, abs/1608.06993, 2016.

[23] Yuankai Huo, Zhoubing Xu, Yunxi Xiong, Katherine Aboud, Prasanna Parvathaneni, Shunxing Bao, Camilo Bermudez, Susan M. Resnick, Laurie E. Cutting, and Bennett A. Landman. 3d whole brain segmentation using spatially localized atlas network tiles. *CoRR*, abs/1903.12152, 2019.

[24] Simon Jégou, Michal Drozdzal, David Vázquez, Adriana Romero, and Yoshua Bengio. The one hundred layers tiramisu: Fully convolutional densenets for semantic segmentation. *CoRR*, abs/1611.09326, 2016.

[25] Mark Jenkinson, Christian F. Beckmann, Timothy E.J. Behrens, Mark W. Woolrich, and Stephen M. Smith. Fsl. *NeuroImage*, 62(2):782 – 790, 2012. 20 YEARS OF fMRI.

[26] Eduard T. Klapwijk, Ferdi van de Kamp, Mara van der Meulen, Sabine Peters, and Lara M. Wierenga. Qoala-t: A supervised-learning tool for quality control of freesurfer segmented mri data. *NeuroImage*, 189:116 – 129, 2019.

[27] Arno Klein and Jason Tourville. 101 labeled brain images and a consistent human cortical labeling protocol. *Frontiers in Neuroscience*, 6:171, 2012.

[28] B Landman and S Warfield. Miccai 2012 workshop on multi-atlas labeling. In *Medical image computing and computer assisted intervention conference*, 2012.

[29] Bennett A. Landman, Alan J. Huang, Aliya Gifford, Deepti S. Vikram, Issel Anne L. Lim, Jonathan A.d. Farrell, John A. Bogovic, Jun Hua, Min Chen, Samson Jarso, and et al. Multi-parametric neuroimaging reproducibility: A 3-t resource study. *NeuroImage*, 54(4):2854–2866, 2011.

[30] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436, 2015.

[31] Guosheng Lin, Anton Milan, Chunhua Shen, and Ian D. Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. *CoRR*, abs/1611.06612, 2016.

[32] Geert J. S. Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen A. W. M. van der Laak, Bram van Ginneken, and Clara I. Sánchez. A survey on deep learning in medical image analysis. *CoRR*, abs/1702.05747, 2017.

[33] Daniel S Marcus, Michael P Harms, Abraham Z Snyder, Mark Jenkinson, J Anthony Wilson, Matthew F Glasser, Deanna M Barch, Kevin A Archie, Gregory C Burgess, Mohana Ramaratnam, et al. Human connectome project informatics: quality control, database services, and data visualization. *Neuroimage*, 80:202–219, 2013.

[34] Daniel S. Marcus, Tracy H. Wang, Jamie Parker, John G. Csernansky, John C. Morris, and Randy L. Buckner. Open access series of imaging studies (oasis): Cross-sectional mri data in young, middle aged, nondemented, and demented older adults. *Journal of Cognitive Neuroscience*, 19(9):1498–1507, 2007.

[35] Kenneth O McGraw and Seok P Wong. Forming inferences about some intraclass correlation coefficients. *Psychological methods*, 1(1):30, 1996.

[36] Fausto Milletari, Seyed-Ahmad Ahmadi, Christine Kroll, Annika Plate, Verena E. Rozanski, Juliana Maiostre, Johannes Levin, Olaf Dietrich, Birgit Ertl-Wagner, Kai Bötzel, and Nassir Navab. Hough-cnn: Deep learning for segmentation of deep brain regions in MRI and ultrasound. *CoRR*, abs/1601.07014, 2016.

[37] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. *CoRR*, abs/1606.04797, 2016.

[38] Victoria L. Morgan, Arabinda Mishra, Allen T. Newton, John C. Gore, and Zhaohua Ding. Integrating functional and diffusion magnetic resonance imaging for analysis of structure-function relationship in the human language network. *PLoS ONE*, 4(8), 2009.

[39] Alexis Moscoso, Jesús Silva-Rodríguez, Jose Manuel Aldrey, Julia Cortés, Anxo Fernández-Ferreiro, Noemí Gómez-Lado, Álvaro Ruibal, Pablo Aguiar, Alzheimer's Disease Neuroimaging Initiative, et al. Prediction of alzheimer's disease dementia with mri beyond the short-term: Implications for the design of predictive models. *NeuroImage: Clinical*, 23:101837, 2019.

[40] Ruth Peters. Ageing and the brain. *Postgraduate medical journal*, 82(964):84–88, 2006.

[41] R C Petersen, P S Aisen, L A Beckett, M C Donohue, A C Gamst, D J Harvey, C R Jack, W J Jagust, L M Shaw, A W Toga, and et al. Alzheimer's disease neuroimaging initiative (adni): clinical characterization, Jan 2010.

[42] Lorenzo Pini, Michela Pievani, Martina Bocchetta, Daniele Altomare, Paolo Bosco, Enrica Cavedo, Samantha Galluzzi, Moira Marizzoni, and Giovanni B Frisoni. Brain atrophy in alzheimer's disease and aging. *Ageing research reviews*, 30:25–48, 2016.

[43] Martin Reuter, Nicholas J Schmansky, H Diana Rosas, and Bruce Fischl. Within-subject template estimation for unbiased longitudinal image analysis. *Neuroimage*, 61(4):1402–1418, 2012.

[44] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. *CoRR*, abs/1505.04597, 2015.

[45] Abhijit Guha Roy, Sailesh Conjeti, Nassir Navab, and Christian Wachinger. Quicknat: Segmenting MRI neuroanatomy in 20 seconds. *CoRR*, abs/1801.04161, 2018.

[46] Cathie Sudlow, John Gallacher, Naomi Allen, Valerie Beral, Paul Burton, John Danesh, Paul Downey, Paul Elliott, Jane Green, Martin Landray, Bette Liu, Paul Matthews, Giok Ong, Jill Pell, Alan Silman, Alan Young, Tim Sprosen, Tim Peakman, and Rory Collins. Uk biobank: An open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLOS Medicine*, 12(3):1–10, 03 2015.

[47] Carole H. Sudre, Wenqi Li, Tom Vercauteren, Sébastien Ourselin, and M. Jorge Cardoso. Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. *CoRR*, abs/1707.03237, 2017.

[48] Theo GM van Erp, Derrek P Hibar, Jerod M Rasmussen, David C Glahn, Godfrey D Pearlson, Ole A Andreassen, Ingrid Agartz, Lars T Westlye, Unn K Haukvik, Anders M Dale, et al. Subcortical brain volume abnormalities in 2028 individuals with schizophrenia and 2540 healthy controls via the enigma consortium. *Molecular psychiatry*, 21(4):547, 2016.

[49] David C Van Essen, Stephen M Smith, Deanna M Barch, Timothy EJ Behrens, Essa Yacoub, Kamil Ugurbil, Wu-Minn HCP Consortium, et al. The wu-minn human connectome project: an overview. *Neuroimage*, 80:62–79, 2013.

[50] Christian Wachinger, Martin Reuter, and Tassilo Klein. Deepnat: Deep convolutional neural network for segmenting neuroanatomy. *CoRR*, abs/1702.08192, 2017.

[51] Hongzhi Wang and Paul Yushkevich. Multi-atlas segmentation with joint label fusion and corrective learning—an open source implementation. *Frontiers in neuroinformatics*, 7:27, 2013.

[52] Simon K Warfield, Kelly H Zou, and William M Wells. Simultaneous truth and performance level estimation (staple): an algorithm for the validation of image segmentation, Jul 2004.

[53] Wenlu Zhang, Rongjian Li, Houtao Deng, Li Wang, Weili Lin, Shuiwang Ji, and Dinggang Shen. Deep convolutional neural networks for multi-modality isointense infant brain image segmentation, Mar 2015.

# Appendix

**Schematic diagram for U-Net model**

Figure 7 shows the schematic diagram for the vanilla U-Net model.



Figure 7: Schematic Diagram of Vanilla U-Net

**Plots for Voxel Count**

Figure 8 and figure 9 shows the mean voxel count for each class in an MRI.

**Detailed Dice scores for all the ROIs for DenseUNet model**

**NYU Data**

The plots 10 and 11 show the box-plot of dice scores for all the ROIs for NYU dataset.

**Comparison of Finetuned v/s non-Finetuned model**

**DenseUNet**

The plots 12 and 13 show the difference in the model's performance on the Mindboggle Data with and without finetuning.

Figure 8: Plot of mean voxel count for top 53 Region of Interests (ROIs): The count is taken over 20 datasets and the average is reported. Average count stands for the average number of voxels for a particular label in the Freesurfer output



Figure 9: Plot of mean voxel count for bottom 54 Region of Interests (ROIs): The count is taken over 20 datasets and the average is reported. Average count stands for the average number of voxels for a particular label in the Freesurfer output

**Comparison of DenseUNet v/s U-Net model**

The plots 14 and 15 show the difference in the model's performance on the Mindboggle Data for DenseUNet and U-Net.

**Faulty Freesurfer Segmentation**

**Putamen**

Figure 16 shows the difference in the segmentation outputs of Freesurfer and the proposed model.

18

Figure 10: Box plot showing the detailed dice scores for the First 51 ROIs of NYU Dataset

**Pallidum**

Figure 17 shows the difference in the segmentation outputs of Freesurfer and the proposed model.

**Comparison of models trained from different views**

Plots 20, 21, 18, and 19 show the difference in the performance of the model trained using 2D slices from different views (coronal, sagittal and axial). For performing the initial view selection, we train the models only using 30,000 2D slices of HCP data and validated using 5,000 2D slices. The validation set is also from the HCP dataset. Based on the preliminary experiment, it was seen that model trained using the 2D coronal slices performed the best for all the ROIs.

**Correlation of dice score with the size of ROI**

The plot 22 shows the variation of dice score with the size of the ROI.

Figure 11: Box plot showing the detailed dice scores for the last 51 ROIs of NYU Dataset



Figure 12: Box plot showing the comparison of dice scores for the Finetuned and non-Finetuned DenseUNet (First 51 ROIs of Mindboggle dataset)

Figure 13: Box plot showing the comparison of dice scores for the Finetuned and non-Finetuned DenseUNet (Last 51 ROIs of Mindboggle dataset)

Figure 14: Box plot showing the comparison of dice scores for DenseUNet and U-Net model(First 51 ROIs of Mindboggle dataset)

Figure 15: Box plot showing the comparison of dice scores for DenseUNet and U-Net model(Last 51 ROIs of Mindboggle dataset)



Figure 16: The image shows the original slice of an MRI, Freesurfer's (FS's) prediction for Putamen and the proposed model's prediction of the same. It is evident that FS's prediction don't obey the boundaries and are non-natural looking grainy segmentation whereas the proposed model's prediction obey the segment boundaries and are much more natural looking
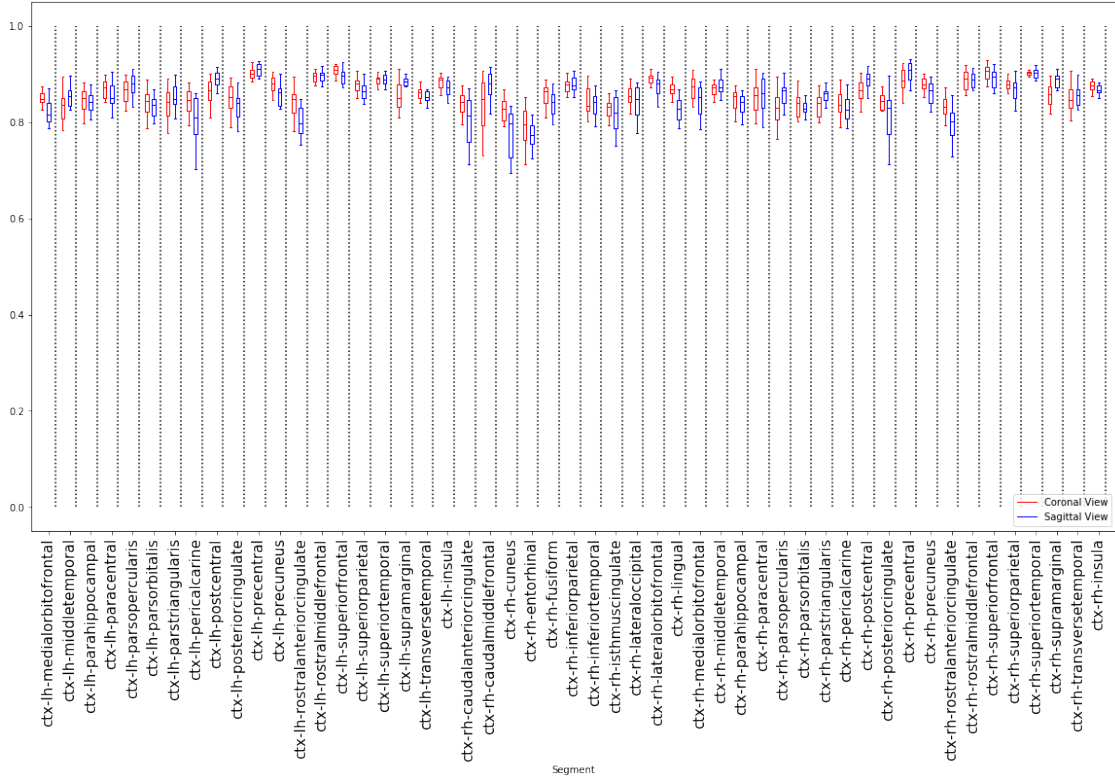
Figure 17: The image shows the original slice of an MRI, Freesurfer's (FS's) prediction for Pallidum and the proposed model's prediction of the same. It is evident that FS's prediction don't obey the boundaries and are non-natural looking grainy segmentation whereas the proposed model's prediction obey the segment boundaries and are much more natural looking



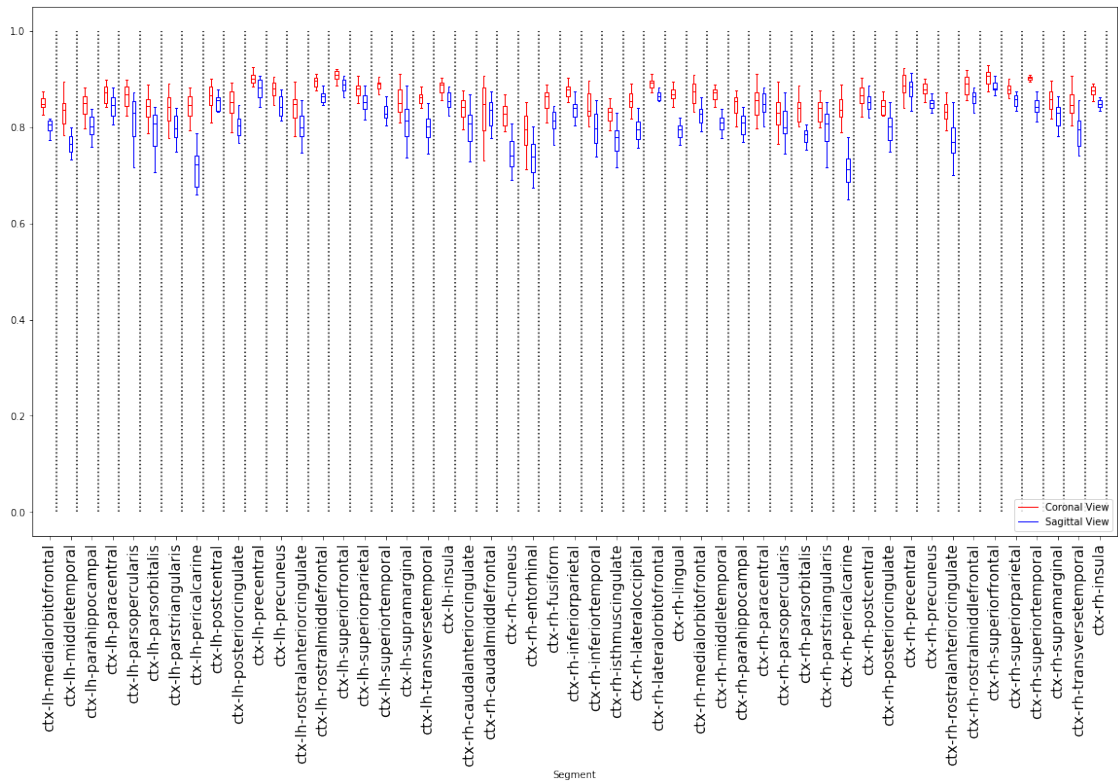Figure 18: Box plot showing the comparison of dice scores for the coronally and sagittally trained model (First 51 ROIs of HCP data)

Figure 19: Box plot showing the comparison of dice scores for the coronally and sagittally trained model (Last 51 ROIs of HCP data)
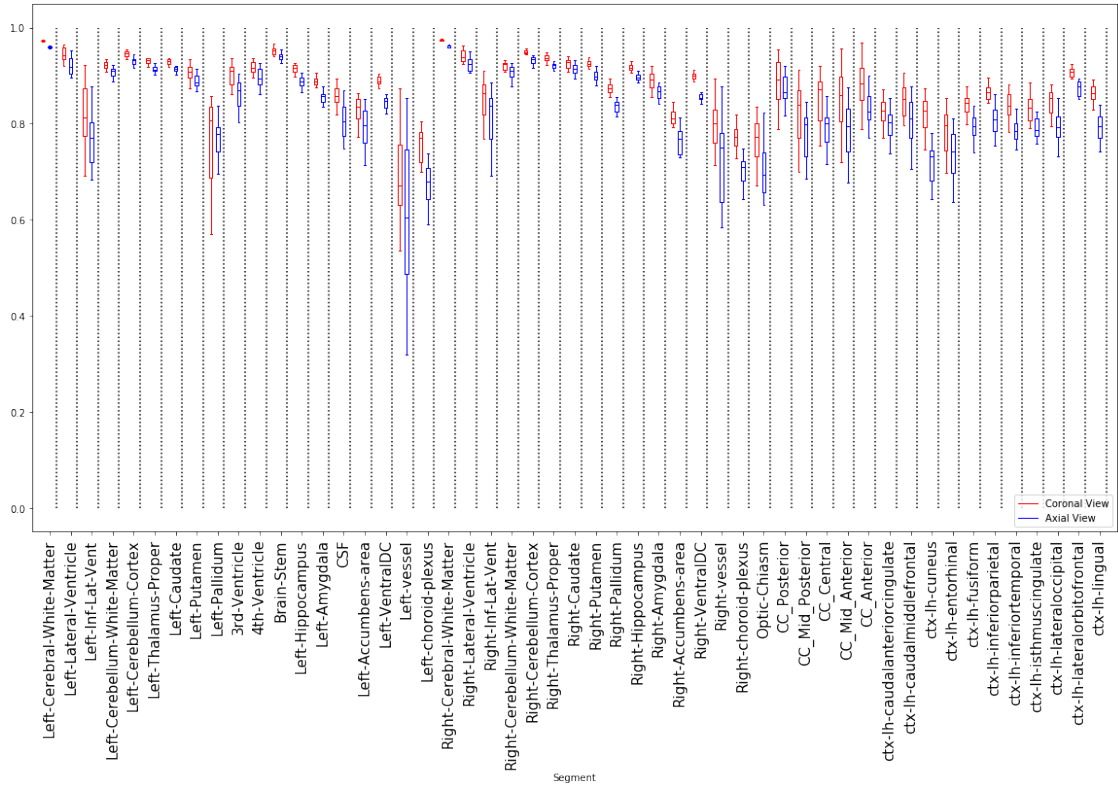
Figure 20: Box plot showing the comparison of dice scores for the coronally and axially trained model (First 51 ROIs of HCP data)
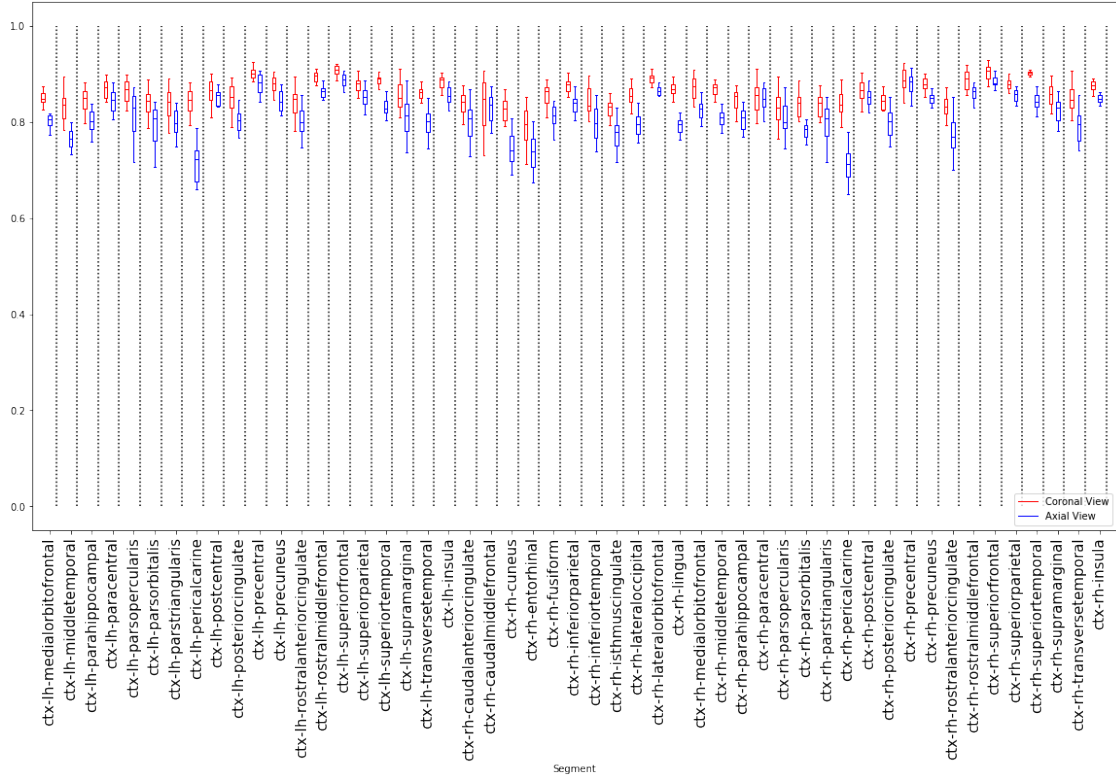
Figure 21: Box plot showing the comparison of dice scores for the coronally and axially trained model (Last 51 ROIs of HCP data)
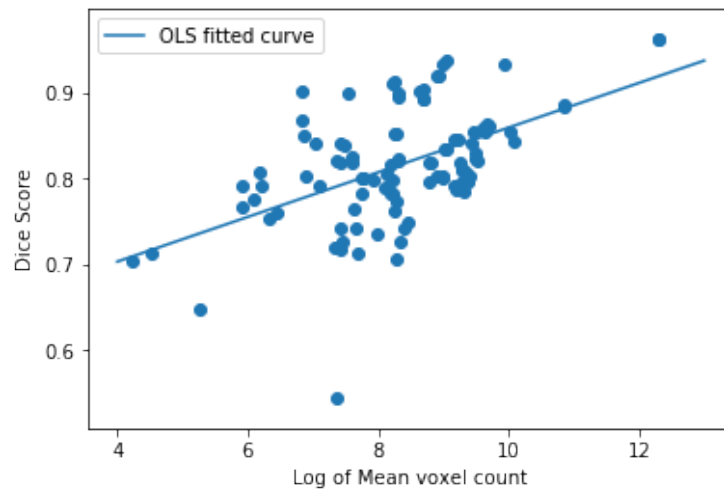


Figure 22: Plot shows the dice score vs log(mean voxel count) for all 102 ROIs (each ROI is a represented as a point on the plot). Pearson's correlation coefficient = 0.521 and is statistically significant with 99.9% significance level.