

Large-Scale Evaluation of ANTs and FreeSurfer Cortical Thickness Measurements

Nicholas J. Tustison^{a,1}, Philip A. Cook^b, Arno Klein^c, Gang Song^b, Sandhitsu R. Das^b, Jeffrey T. Duda^b, Benjamin M. Kandel^b, Niels van Strien^c, James R. Stone^a, James C. Gee^b, Brian B. Avants^b

^aDepartment of Radiology and Medical Imaging, University of Virginia, Charlottesville, VA

^bPenn Image Computing and Science Laboratory, University of Pennsylvania, Philadelphia, PA

^cSage Bionetworks, Seattle, WA

^dDepartment of Circulation and Medical Imaging, Norwegian University of Science and Technology, Trondheim, Norway

Abstract

Many studies of the human brain have explored the relationship between cortical thickness and cognition, phenotype, or disease. Due to the subjectivity and time requirements in manual measurement of cortical thickness, scientists have relied on robust software tools for automation which facilitate the testing and refinement of neuroscientific hypotheses. The most widely used tool for cortical thickness studies is the publicly available, surface-based FreeSurfer package. Critical to the adoption of such tools is a demonstration of their reproducibility, validity, and the documentation of specific implementations that are robust across large, diverse imaging datasets. To this end, we have developed the automated, volume-based Advanced Normalization Tools (ANTs) cortical thickness pipeline comprising well-vetted components such as SyGN (multivariate template construction), SyN (image registration), N4 (bias correction), Atropos (*n*-tissue segmentation), and DiReCT (cortical thickness estimation). In this work, we have conducted the largest evaluation of automated cortical thickness measures in publicly available data, comparing FreeSurfer and ANTs measures computed on 1205 images from four open data sets (IXI, MMRR, NKI, and OASIS), with parcellation based on the recently proposed Desikan-Killiany-Tourville (DKT) cortical labeling protocol. We found good scan-rescan repeatability with both FreeSurfer and ANTs measures. Given that such assessments of precision do not necessarily reflect accuracy or ability to make statistical inferences, we further tested the neurobiological validity of these approaches by evaluating thickness-based prediction of age and gender. ANTs is shown to have a higher predictive performance than FreeSurfer for both of these measures. In promotion of open science, we make all of our scripts, data, and results publicly available which complements the use of open image data sets and the open source availability of the proposed ANTs cortical thickness pipeline.

Keywords: advanced normalization tools, age prediction, MRI, gender prediction, open science, scientific reproducibility

1. Introduction

Magnetic resonance imaging-based structural analysis of the human brain plays a fundamental role in identifying the relationship between cortical morphology, disease, and cognition. Such research has produced a significant amount of literature concerning cortical variability and its developmental correlates including inquiries into normal aging [64] and gender differences [38]. Conditional abnormalities from Alzheimer's disease and frontotemporal dementia [18, 16] to Parkinson's [31] and Huntington's disease [49] have also demonstrated sensitivity to cortical thickness measurements. Additional exploration has spanned such topics as autism [9], athletic ability [67], male-to-female transsexuality [39], obesity [47], and Tetris-playing ability in female adolescents [26]. Although these findings are subject to debate and interpretation [24], the availability of quantitative computational methods for extracting cortical thickness measures has proven invaluable for developing and refining fundamental neuroscience hypotheses.

Computational methods for analyzing the cortex may be broadly characterized as surface mesh-based or volumetric [52, 10]. Representative of the former is the FreeSurfer cortical modeling software package [12, 21, 19, 20, 22] which owes its popularity to public availability, excellent documentation, good performance, and integration with other toolkits, such as the extensive FMRIB software library [59]. Similar to other surface-based cortical thickness estimation approaches (e.g., [14, 42, 41, 32]), the outer cortical and gray/white matter surfaces from individual subject MR data are modeled with polygonal meshes which are then used to determine local cortical thickness values based on a specified correspondence between the surface models.

Image volumetric (or meshless) techniques vary both in their algorithms as well as in the underlying definitions of cortical thickness. An early, foundational technique is the method of [29] in which the inner and outer surface geometry is used to determine the solution to Laplace's equation where thickness is measured by integrating along the tangents of the resulting field lines spanning the boundary surfaces. Subsequent contributions improved upon the original formulation. For example, in [70], a Eulerian partial differential equation approach was proposed to

¹Corresponding author: PO Box 801339, Charlottesville, VA 22908; T: 434-924-7730; email address: ntustison@virginia.edu.

Partial support provided by the Defense Health Program through U.S. Army Medical Research Acquisition Activity, Grant Number W81XWH-09-2-0055.

facilitate the computation of correspondence paths. Extending the surface-based work of [41], the hybrid approach of [32] uses the discrete Laplacian field to deform the white matter surface mesh towards the outer cortical surface. Other volumetric algorithms employ coupled level sets [71], model-free intelligent search strategies either normal to the gray-white matter interface [52], or using a min-max rule [63]. Most relevant to this work is the DiReCT (Diffeomorphic Registration-based Cortical Thickness) algorithm proposed in [13] where generated diffeomorphic mappings between the gray/white matter and exterior cortical surfaces are used to propagate thickness values through the cortical gray matter.

The general lack of availability of published algorithms [35] (not to mention critical preprocessing components) is a strong deterrent to the use or evaluation of these algorithms by external researchers. For example, one recent evaluation study [10] compared FreeSurfer (a surface-based method) with two volumetric methods, viz., [29, 13]. Whereas the entire FreeSurfer processing pipeline has been made publicly available, refined by the original authors and other contributors, and described in great detail (specifically in terms of suggested parameters), both volumetric methods were implemented and run by the authors of the evaluation (not by the algorithm developers) using unspecified parameters with relatively small, private data sets, making the comparisons less than ideal (see [62] for further discussion concerning the issue of instrumentation bias and scientific reproducibility in the use and evaluation of software). Further complicating such comparisons is the potential for bias, such as interpolation artifacts when converting surface to volume data or vice versa [33].

We provide below a brief description of our proposed pipeline, which produces a volumetric cortical thickness map from an individual subject’s T1-weighted MRI. Additionally, we note that it is freely available as part of the Advanced Normalization Tools (ANTs) software package. This includes all the necessary preprocessing steps consisting of well-vetted, previously published algorithms for bias correction [61], brain extraction [2], n -tissue segmentation [4], template construction [5], and image normalization [3]. More importantly, we provide explicit coordination among these components within a set of well-documented shell scripts which are also available in the ANTs repository where parameters have been tuned by ANTs developers (N.T. and B.A.).

Here we demonstrate the use of the described framework in processing 1205 publicly available, T1-weighted brain MR images drawn from four well-known data sets. For comparative evaluation we also process the same data using the standard FreeSurfer cortical thickness processing protocol. Similar to previous work [e.g., 10], we are able to report repeatability assessments for both frameworks using subsets of the data with repeated acquisitions. However, repeatability (or, more generally, *precision*) is not conceptually equivalent to *accuracy* and, thus, does not provide a complete perspective for determination of measurement quality. Although FreeSurfer validation has included histological [50] and image-drawn [36] comparisons, such manual assessments were extremely limited in terms of number of subjects and the number of cortical regions. In ad-

dition, there was no mention in these studies of the number of human observers making these measurements nor discussion of quality assurance. Alternatively, without ground truth, other forms of evidence can be adduced [e.g., 6] in making comparative inferences. In this work we use demographic-based assessments (based on well-studied relationships between cortical thickness and age/gender) to show that ANTs outperforms FreeSurfer-based thickness estimation for these data in terms of prediction.

2. Methods and Materials

2.1. Public data resources

A comparative evaluation between FreeSurfer and ANTs was run on four publicly available data sets: IXI, MMRR, NKI, and OASIS. In addition to these data, we used a subset of the MindBoggle-101 data labeled using the Desikan-Killiany-Tourville (DKT) protocol [34] to define the regions of interest (ROI) in the analysis. This latter data set was not included in the thickness analysis. All five data sets are described below.

2.1.1. Public data for thickness estimation evaluation

Diverse and publicly available data sets collected from multiple sites with a mixture of 3T and 1.5T T1-weighted brain images were analyzed using the ANTs and FreeSurfer pipelines. Subjects in this data set span the age range from 4 to 96 years old. This strategy tested robustness to variation in head position, brain shape, defacing, image contrast, inhomogeneity, imaging artifacts, field strength, and the broad variation in extracerebral tissue. Failure can occur in initial brain extraction, segmentation, registration, or bias correction, any of which can lead to an inaccurate cortical thickness measurement. In total, we processed 1,205 T1-weighted images from four different public data sets to obtain cortical thickness values for both cortical thickness analysis softwares. Below we describe the four data sets:

IXI. Initially, we processed 581 T1-weighted images from the IXI data set, but only 563 subjects (313 females, 250 males) were included in the post-processing analysis due to missing demographic information, which would have prevented an accurate estimate of the age at the time of image acquisition. These data were imaged at three sites with several modalities acquired (T1-weighted, T2-weighted, proton density, magnetic resonance angiography, and diffusion tensor imaging). The database also includes demographic information such as date of birth, date of scan, weight, height, ethnicity, occupation category, educational level, and marital status.

MMRR. The Multi-Modal MRI Reproducibility Resource (MMRR) data set, was originally described in [37] consisting of 21 subjects (10 females, 11 males) and features a rich set of modalities, as well as repeated scans.

Table 1: The 31 cortical labels (per hemisphere) of the DKT atlas.

1) caudal anterior cingulate	17) pars orbitalis
2) caudal middle frontal	18) pars triangularis
3) cuneus	19) pericalcarine
4) entorhinal	20) postcentral
5) fusiform	21) posterior cingulate
6) inferior parietal	22) precentral
7) inferior temporal	23) precuneus
8) isthmus cingulate	24) rosterior anterior cingulate
9) lateral occipital	25) rostral middle frontal
10) lateral orbitofrontal	26) superior frontal
11) lingual	27) superior parietal
12) medial orbitofrontal	28) superior temporal
13) middle temporal	29) supramarginal
14) parahippocampal	30) transverse temporal
15) paracentral	31) insula
16) pars opercularis	

NKI. In support of open science, the 1000 Functional Connectomes Project was initiated on December 11, 2009 by various members of the MRI community seeking to form collaborative partnerships among imaging institutions for sharing well-documented multimodal image sets accompanied by phenotypic data. One such contribution is the Nathan Klein Institute (NKI)/Rockland sample, consisting of 186 T1-weighted images (87 females, 99 males).

OASIS. The initial Open Access Series of Imaging Studies (OASIS) data set consisted of 433 T1-weighted images. We processed all of these data, but 100 were excluded from our analysis due to probable Alzheimer’s disease ($CDR > 0$) and an additional 20 repeat scans were excluded, resulting in 313 individual subject scans included in the normal group statistical analysis (118 males, 195 females). Ages were between 18 and 96 and all subjects are right-handed.

2.1.2. MindBoggle-101 data for ROI definitions

In [34] the authors proposed the DKT cortical labeling protocol—a modification of the popular Desikan-Killiany protocol [15] to improve cortical labeling consistency and to improve FreeSurfer’s cortical classification of 31 cortical regions per hemisphere, listed in Table 1. Forty manually labeled brains were used to construct the DKT40 Gaussian classifier atlas, which is now bundled with current versions of FreeSurfer and used to automate anatomical labeling of MRI data. Since the regional thickness values generated by FreeSurfer follow this protocol, these anatomical labels provide a common standard for comparison between ANTs and FreeSurfer.

The work of [34] also resulted in a publicly available set of manually edited labels following the DKT protocol in 101 T1-weighted brain images from different sources, including a subset of 20 images from the OASIS data set (specifically, the test-retest data). These 20 images are used in the MALF step that defines the volumetric cortical regions for each subject.

2.2. ANTs volume-based cortical thickness estimation pipeline

The ANTs cortical thickness estimation workflow is illustrated in Figure 1. The steps are as follows:

1. Initial N4 bias correction on input anatomical MRI
2. Brain extraction using a hybrid segmentation/template-based strategy
3. Alternation between prior-based segmentation and “pure tissue” posterior probability weighted bias correction using Atropos and N4
4. DiReCT-based cortical thickness estimation
5. Optional normalization to specified template and multi-atlas cortical parcellation

Each component, including both software and data, is briefly detailed below with the relevant references for additional information. Although other preprocessing components are possible (e.g., noise reduction as in [57]), the major steps constituting the ANTs pipeline are limited to those enumerated above.

The coordination of all the algorithmic components is encapsulated in the shell script `antsCorticalThickness.sh` with subcomponents delegated to `antsBrainExtraction.sh` and `antsAtroposN4.sh`. A representative script command is reproduced in Listing 1 for a single IXI subject to demonstrate the simplicity and mature status of what we propose in this work and a comparison with the analogous FreeSurfer command. Option descriptions are provided by invoking the help option, i.e., “`antsCorticalThickness.sh -h`”.

```
# Processing calls for subject IXI002-Guys-0828-T1
# ANTs
antsCorticalThickness.sh \
-a IXI/T1/IXI002-Guys-0828-T1.nii.gz \
-e IXI/template/T_template0.nii.gz \
-m IXI/template/T_template0ProbabilityMask.nii.gz \
-f IXI/template/T_template0ExtractionMask.nii.gz \
-p IXI/template/Priors/priors%d.nii.gz \
-o IXI/ANTsResults/IXI002-Guys-0828-
# FreeSurfer
recon-all \
-i IXI/T1/IXI002-Guys-0828-T1.nii.gz \
-s IXI002-Guys-0828 \
-sd IXI/FreeSurferResults/ \
-all
```

Listing 1: Analogous ANTs and FreeSurfer command line calls for a single IXI subject in the evaluation study.

2.2.1. Anatomical template construction

Certain preprocessing steps, such as brain extraction and segmentation, rely on templates and corresponding spatial priors. In addition, normalizing images to a standard coordinate system reduces intersubject variability in population studies. Various approaches exist for determining an optimal template, such as the selection of a preexisting template based on a single individual (e.g., the Talairach atlas [60]) or a publicly available average of multiple individuals (e.g., the MNI [11] or ICBM [43] templates), or an average template constructed from the individuals under study. The work of [5] explicitly models the geometric component of the normalized space during optimization to produce such mean templates. Coupling the intrinsic

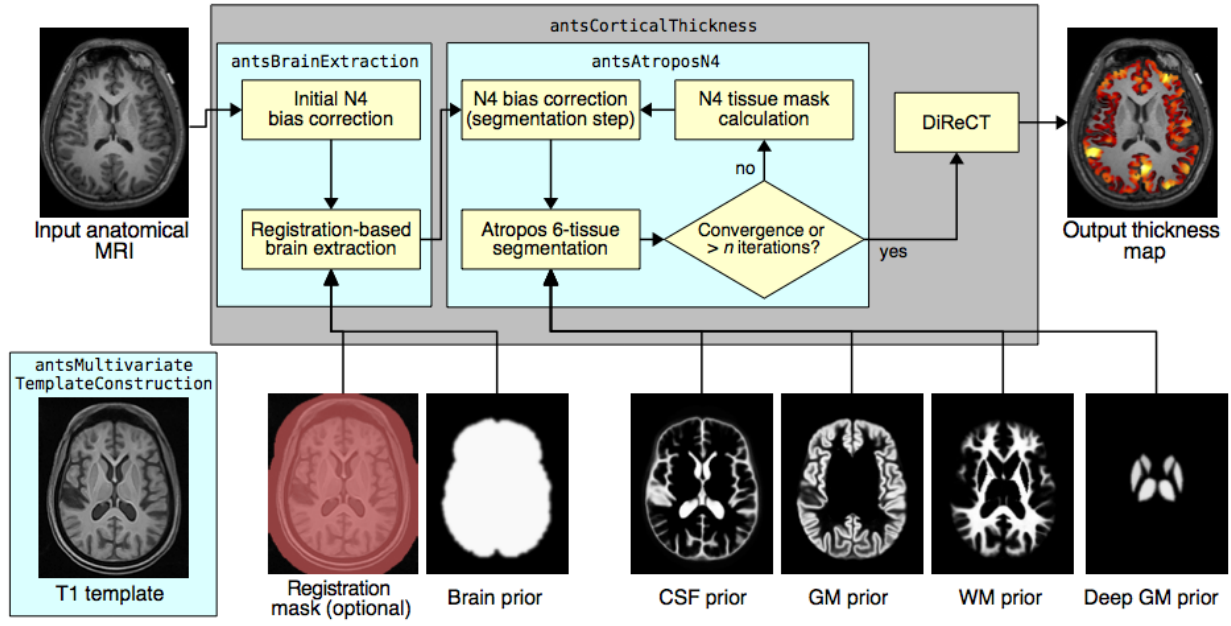


Figure 1: Illustration of the main components of the ANTs processing workflow containing all elements for determining cortical thickness. We also included the domain of operations for the selected scripts. Not shown are the probability maps for the brain stem and cerebellum priors. All template-based prior probability maps are generated prior to pipeline processing of each individual subject.

symmetry of SyN pairwise registration [3] and an optimized shape-based sharpening/averaging of the template appearance, Symmetric Group Normalization (SyGN) is a powerful framework for producing optimal population-specific templates. The five templates used in this evaluation study are represented in Figure 2.

2.2.2. N4 bias field correction

Critical to quantitative processing of MRI is the minimization of field inhomogeneity effects which produce artificial low frequency intensity variation across the image. Large-scale studies, such as ADNI, employ perhaps the most widely used bias correction algorithm, N3 [56], as part of their standard protocol [7].

In [61] we introduced an improvement of N3, denoted as “N4”, which demonstrates a significant increase in performance and convergence behavior on a variety of data. This improvement is a result of an enhanced fitting routine (which includes multi-resolution capabilities) and a modified optimization formulation. For our workflow, the additional possibility of specifying a weighted mask in N4 permits the use of a “pure tissue” probability map (described below) calculated during the segmentation pipeline for further improvement of bias field estimation.

N4 is used in two places during the individual subject processing (cf Figure 1). It is used to generate an initial bias-corrected image for use in brain extraction. The input mask is created by adaptively thresholding the background from the foreground using Otsu’s algorithm [45]. Following brain extraction, six-tissue (cerebrospinal fluid, cortical gray matter,

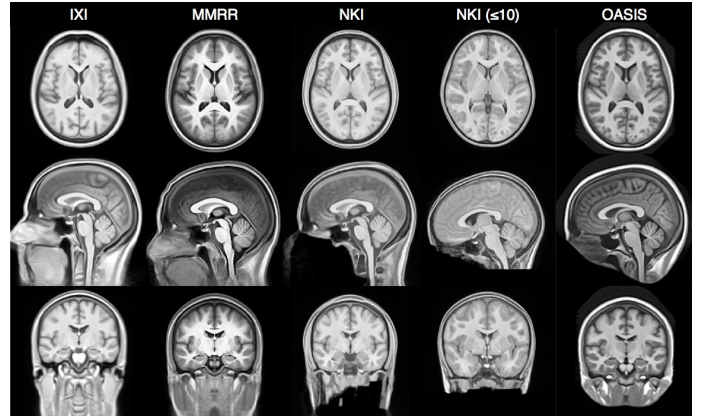


Figure 2: Population-specific templates for each of the four public data sets used for cortical thickness estimation. The benefit of using such population-specific templates is obvious when one sees the variability in acquisition and data preparation (e.g., defacing protocols).

white matter, deep gray matter, brain stem, and cerebellum) segmentation involves iterating between bias field correction using the current pure tissue probability map as a weight mask and then using that bias-corrected image as input for the Atropos segmentation step (described below).

2.2.3. Brain extraction

Brain extraction using ANTs combines template building, high-performance brain image registration, and Atropos segmentation with topological refinements. An optimal template [5], i.e., a mean shape and intensity image representation of a particular cohort, is first constructed using structural MRI data. Template construction iterates between estimating the optimal template and registering each subject to the optimal template. In this work, we perform the additional step of building separate templates for each cohort and propagating the probabilistic mask to each cohort template using registration of the T1-weighted templates (cf Figure 2). A probabilistic brain extraction mask for the new template can then be generated by warping an existing mask to the template or by averaging the warped, whole brain labels of subjects registered to the new template, if such labels are available. Further refinements include thresholding the warped brain probability map at 0.5 and dilating the resulting mask with a radius of two voxels. Atropos is used to generate an initial three-tissue segmentation estimate within the mask region. Each of the three tissue labels undergoes separate morphological operations including hole-filling and erosion. These results are then combined to create the brain extraction mask which is further refined by additional dilation, erosion, and hole-filling operations.

In previous work [2] we compared an earlier version of our extraction method with publicly available brain extraction algorithms, including AFNI’s 3dIntracranial [66], FSL’s BET2 [58], FreeSurfer’s mri_watershed [53], and BrainSuite [17]. Our hybrid registration/segmentation approach performed with an accuracy comparable to FreeSurfer and a parameter-tuned version of BrainSuite. Figure 3 presents a visual comparison of results derived with the current ANTs brain extraction method and results obtained using FreeSurfer.

2.2.4. Atropos six-tissue segmentation

In [4] we presented an open source n -tissue segmentation software tool (which we denote as “Atropos”) that attempts to distill over 20 years of active research in this area, in particular some of the most seminal work (e.g., [72, 1]). Specification of prior probabilities includes spatially varying Markov Random Field modeling, prior label maps, and prior probability maps typically derived from our template building process. Additional capabilities include handling of multivariate data, partial volume modeling [54], a memory-minimization mode, label propagation, a plug-and-play architecture for incorporation of novel likelihood models which includes both parametric and non-parametric models for both scalar and tensorial images, and alternative posterior formulations for different segmentation tasks.

Due to the important interplay between segmentation and bias correction, we perform multiple $N4 \Rightarrow$ Atropos itera-

tions. A pure tissue probability weight mask generated from the posterior probabilities is derived from the segmentation step. Given N labels and the corresponding N posterior probability maps $\{P_1, \dots, P_N\}$ produced during segmentation, the $N4$ weight mask is created at each $N4 \Rightarrow$ Atropos iteration from

$$P_{\text{pure tissue}}(\mathbf{x}) = \sum_{i=1}^N P_i(\mathbf{x}) \prod_{j=1, j \neq i}^N (1 - P_j(\mathbf{x})). \quad (1)$$

One of the key insights of the original N3 development is the observation that inhomogeneities cause the intensity values of pure tissue peaks to spread in the intensity histogram as though convolved with a Gaussian. A core contribution of N3 is the proposed corrective step of deconvolving the intensity histogram to accentuate the tissue peaks, coupled with a spatial smoothing constraint. The pure tissue probability mask is used in N4 to weight more heavily the influence of voxels corresponding to pure tissue types (as determined by the segmentation) during the deconvolution process while minimizing the contribution of regions such as the gray/white matter interface where peak membership is ambiguous.

Atropos enables prior knowledge to guide the segmentation process where template-based priors are integrated into the optimization with a user-controlled weight. Modulating the likelihood and prior contributions to the posterior probability is essential for producing adequate segmentations. Atropos weights the likelihood and priors according to $P(x|y) \propto P(y|x)^{1-\alpha} P(x)^\alpha$ where α is a user-selected parameter which weights the trade-off between the likelihood and priors terms. A weighting of $\alpha = 0.25$ is the default value based on our extensive experimentation with different parameter weights.

Since cortical thickness estimation only requires the cortical gray and white matter, the deep gray and white matter (both labels and posterior maps) are combined to form a single “white matter” set. This white matter set and the cortical gray matter are the only results from the segmentation step that are used by the DiReCT algorithm (described below).

To generate the priors for each T1 template, we used the multi-atlas label fusion (MALF) algorithm of [65] in conjunction with a labeled subset of the OASIS data set.² First, we normalized the labeled OASIS subset to the template. We then performed MALF on the template using the normalized labeled data as input. This resulted in a labeled, parcellated template consisting of 100+ labels defining the different brain regions. We then condensed this template-specific labeling to the six needed for our analysis, viz., cerebrospinal fluid (CSF), gray matter (GM), white matter (WM), deep gray matter, brain stem, and cerebellum. For example, all cortical regions were assigned a single label representing the gray matter.

These binary masks were then smoothed using Gaussian convolution with a one voxel-width kernel. Since the labelings

² These data were originally acquired by the first and last authors to aid in the MICCAI 2012 Grand Challenge and Workshop on Multi-Atlas Labeling. The data was released under the Creative Commons Attribution-NonCommercial license. Labelings were provided by Neuromorphometrics, Inc. (<http://Neuromorphometrics.com/>) under academic subscription.

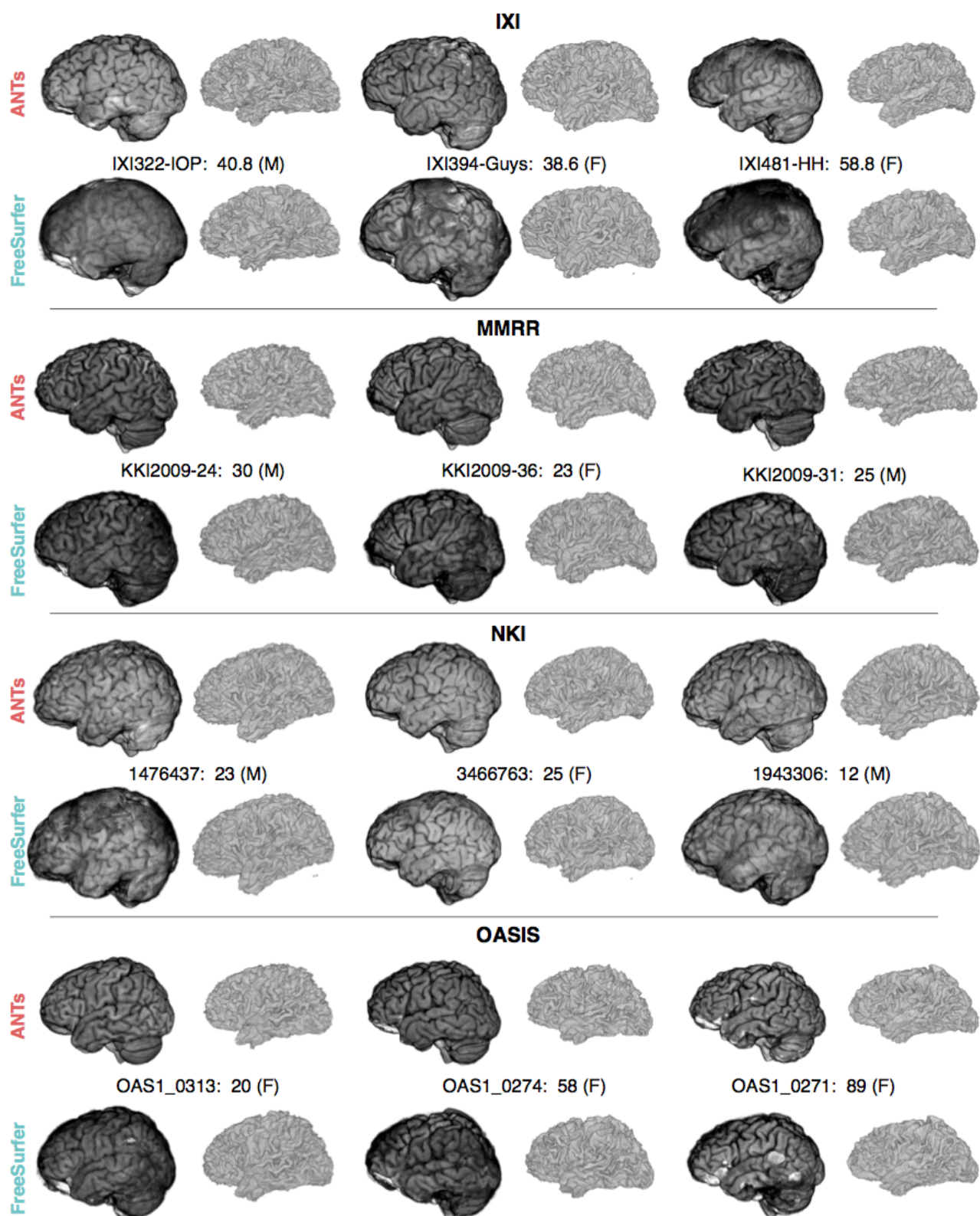


Figure 3: Representative sample of volume brain renderings from the four different cohorts (IXI = rows 1 and 2, MMRR = rows 3 and 4, NKI = rows 5 and 6, OASIS = rows 7 and 8), illustrating the qualitative difference between ANTs and FreeSurfer results, which are arranged top-and-bottom for each subject. Each brain was rigidly registered to the OASIS template for rendering purposes. With each subject we provide subject ID, age, and gender.

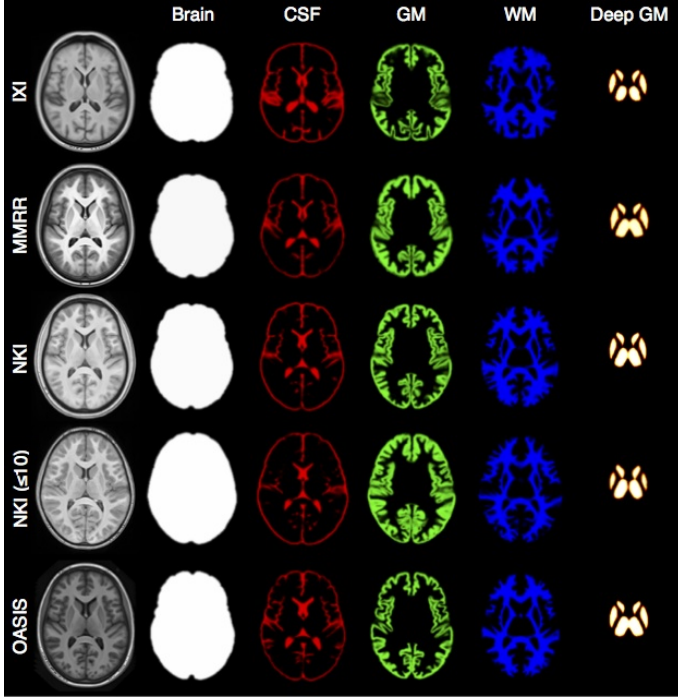


Figure 4: Axial slices from each of the five T1 templates including the corresponding probability masks used for brain extraction and brain tissue segmentation. Not shown are the prior probability maps for brain stem and cerebellum regions.

did not describe the extracerebral CSF, we augmented the CSF prior image with the CSF posterior output from running each template through the segmentation component of the above-described pipeline. This new CSF prior was then subtracted from each of the other five prior probability images and limited to the probability range of $[0, 1]$. The prior probabilities for the five templates used in this evaluation are given in Figure 4.

2.2.5. DiReCT cortical thickness estimation

DiReCT was introduced in [13] and was made available in ANTs as the program *KellySlater*. Since then several improvements have been made and incorporated into the program *KellyKapowski*.³ The more recent implementation has made numerous advances including multi-threading, written in rigorous ITK coding style, and has been made publicly available through ANTs, complete with a unique command line interface design developed specifically for ANTs tools.

2.3. Processing miscellany

Given the documented variability in FreeSurfer results with version and operating system [25] (as we would ex-

pect with our own ANTs pipeline), all data were processed using the same ANTs and FreeSurfer versions on the same hardware platform. Processing was performed using the Linux (CentOS release 6.4) cluster at the University of Virginia (<http://www.uvacse.virginia.edu>) using single-threading with a maximal requested memory footprint of 8 gb for ANTs and 4 gb for FreeSurfer. The development version of ANTs was used for processing (git commit tag: 69d3a5a6c7125ccf07a9e9cf6ef29f0b91e9514f, date Dec. 11, 2013). FreeSurfer version 5.3 x86_64 for CentOS was downloaded on 5 December, 2013 (“freesurfer-Linux-centos6_x86_64-stable-pub-v5.3.0”, release date: 15 May, 2013). The brain extraction and segmentation results from both pipelines were visually inspected to screen for major problems. No manual changes were made for any component of either pipeline and no change was made to the settings of either processing pipeline.

3. Evaluation

Traditional assessment approaches, such as manual labeling, are inadequate for evaluating large-scale performance. We therefore sought to minimize failure rate, quantify the repeatability of cortical thickness measures, and determine whether the ANTs pipeline reveals biologically plausible relationships between the cortex, gender,⁴ and age and how its performance compares to the current de facto standard of FreeSurfer-derived thickness estimation. Collectively, these surrogate measurements allow us to establish data-derived relative performance standards. Additionally, for completeness, we include timing results as that factors into usability.

3.1. Repeatability

Repeat scans of 40 subjects (20 MMRR subjects and 20 OASIS subjects) were used to determine the repeatability of regional cortical thickness measurements, T . Similar to the reproducibility assessment given in [30], we demonstrate this in terms of the percent variability error:

$$\varepsilon = \frac{|T_{scan} - T_{rescan}|}{0.5 \times (T_{scan} + T_{rescan})}. \quad (2)$$

Comparison of the ANTs and FreeSurfer percent variability errors for the 62 DKT regions for both the OASIS and MMRR scan-rescan data sets are given in Figure 5. Mean values are given in Table 2. Although the variance is slightly greater for the set of ANTs measurements, statistical testing per cortical region (two-tailed paired t-test, corrected using false discovery rate) did not indicate non-zero mean differences for either approach for any region.

³ Traditional academic discourse encountered in the published literature rarely contextualizes peculiarities such as algorithmic nomenclature. We briefly mention that this was the source of a rare disagreement between the first and last authors based, as many disagreements are, on a simple misunderstanding and not an affronting existential statement concerning a certain favorite sitcom of the first author’s youth.

⁴ We recognize the distinction often made between “sex” and “gender” (cf <http://www.who.int/gender/whatisgender/en/>). As the demographic information collected during the course of the imaging studies is presumably self-reported, we assume that most self-identify in terms of gender and, therefore, use the term “gender” in data descriptions.

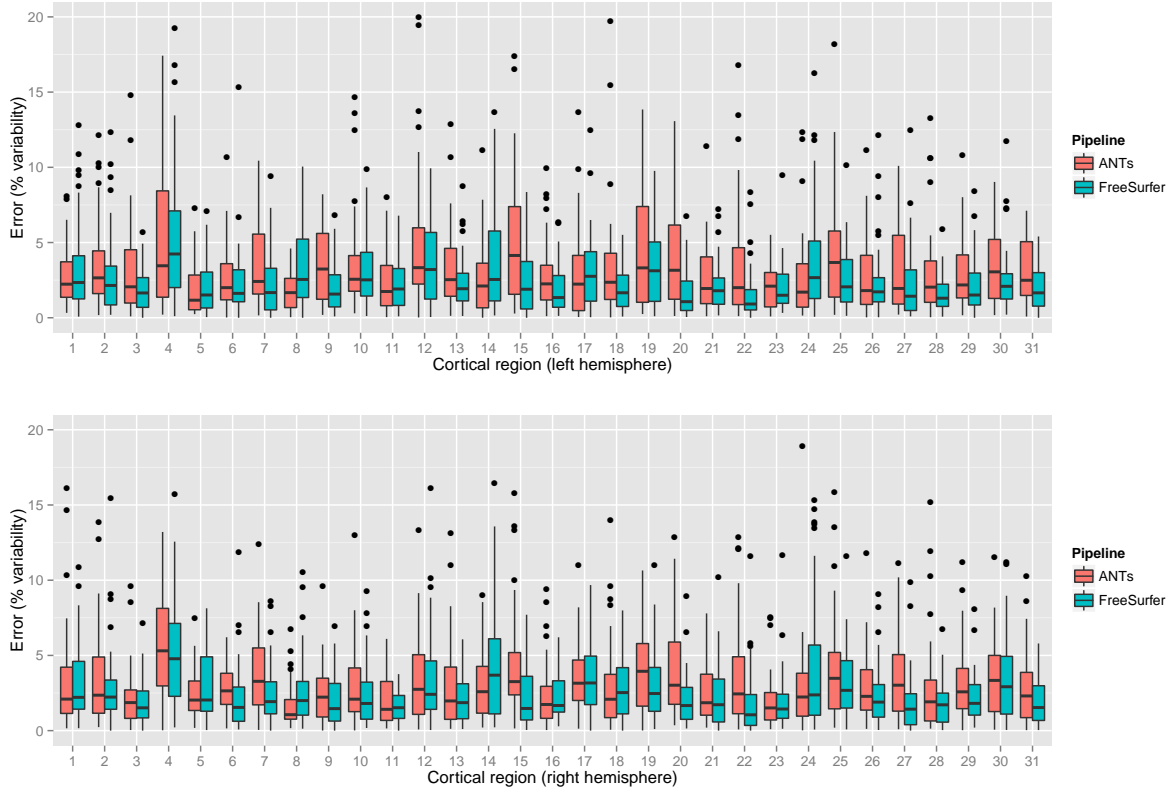


Figure 5: Percent error variability for both ANTs and FreeSurfer pipelines over the left and right hemispheres of both the MMRR and OASIS data subsets within the 62 regions defined by the Desikan-Killiany-Tourville atlas. Both methods demonstrate good repeatability qualities.

Table 2: Mean repeatability error over all regions.

	MMRR	OASIS
ANTs	3.2%	3.3%
FreeSurfer	2.5%	2.8%

We also calculated the intraclass correlation coefficient (“ICC(2,1)” in the notation of [55]) to assess scan/rescan reliability. The ANTs thickness pipeline produced an ICC value of 0.98 and the FreeSurfer thickness pipeline yielded an ICC value of 0.97, indicating good scan/rescan reliability for both ANTs and FreeSurfer.

3.2. Age prediction assessment

Despite good repeatability with both ANTs and FreeSurfer, such measures do not provide an assessment of accuracy or even relative utility. For example, strong priors can yield good repeatability measures but potentially at the expense of data fidelity thus compromising the quality of models (statistical or otherwise) built from such results. Given that ground truth is not available for these data nor for the many studies looking at brain morphology, an indirect method (or set of methods) is required for determining the quality of thickness estimation.

For our first assessment, we modeled age versus regional cortical thickness values to determine which framework produces better predictive thickness estimates. We first subdivided the

thickness data into training and testing subsets with an even split between the two subsets.⁵ We then used the training data to create two models for each pipeline: 1) standard linear regression and 2) random forests (a non-parametric machine learning technique) [8], for estimating age from both ANTs and FreeSurfer thickness values in the testing data.

Table 3: Mean RMSE for age prediction in years.

	Linear Model	Random Forest
ANTs (Combined)	10.7	10.2
FreeSurfer (Combined)	12.3	11.9
ANTs (IXI)	9.3	8.6
FreeSurfer (IXI)	12.3	11.7
ANTs (NKI)	NA [†]	10.9
FreeSurfer (NKI)	NA [†]	13.3
ANTs (OASIS)	15.0	12.4
FreeSurfer (OASIS)	15.0	11.4

[†]Fitting error.

⁵ We tried various training proportions between 10 and 90% (in increments of 10%) to see if that had an effect on relative performance for both age and gender prediction comparisons. Although age predictive capabilities for both pipelines showed improvement (gender prediction was mostly unaffected), the relative outcomes were the same.

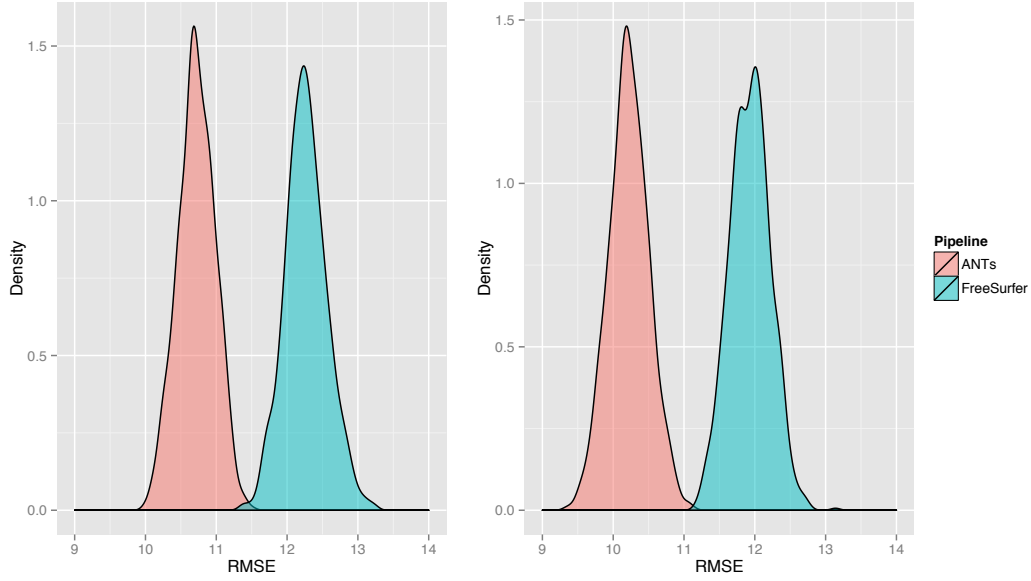


Figure 6: Age prediction RMSE distributions of linear (left) and random forest (right) models for the ANTs- and FreeSurfer-derived thickness values over the combined four cohorts. For both prediction models ANTs RMSE error is lower.

The formula (in the notation of [68]) for the linear model is

$$AGE \sim VOLUME + GENDER + \sum_{i=1}^{62} T(DKT_i) \quad (3)$$

where $T(DKT_i)$ is the average thickness value in region DKT_i and $VOLUME$ is total intracranial volume. Similarly, the random forest model was specified as a combination of all terms using the `randomForest` package in R with the default settings and 200 trees.

In order to ensure a fair comparison, the procedure described above consisting of training and testing steps was performed for $n = 1000$ permutations to elicit a performance distribution which we measure using the relative mean square error (RMSE):

$$RMSE = \sqrt{\frac{\sum (AGE_{true} - AGE_{predicted})^2}{N}}. \quad (4)$$

Due to the limited range in ages across data sets, we restricted training and testing to the age range [20, 80]. The resulting distributions are illustrated in Figure 6. In addition to a combined assessment, we also perform separate model prediction for each of the three larger data sets (i.e., IXI, NKI, and OASIS).

ANTs-based RMSE values were lower for both models and each of the four different subset comparisons except for the random forest model constructed from the OASIS data set. All mean RMSE values are provided in Table 3.

To further elucidate the regional differences in predictive power specifically in the random forest model, we provide variable importance plots for both pipelines using the mean decrease in accuracy measure in Figure 7. During random forest model construction (specifically the out-of-bag error calculation stage), the decrease in prediction accuracy with the

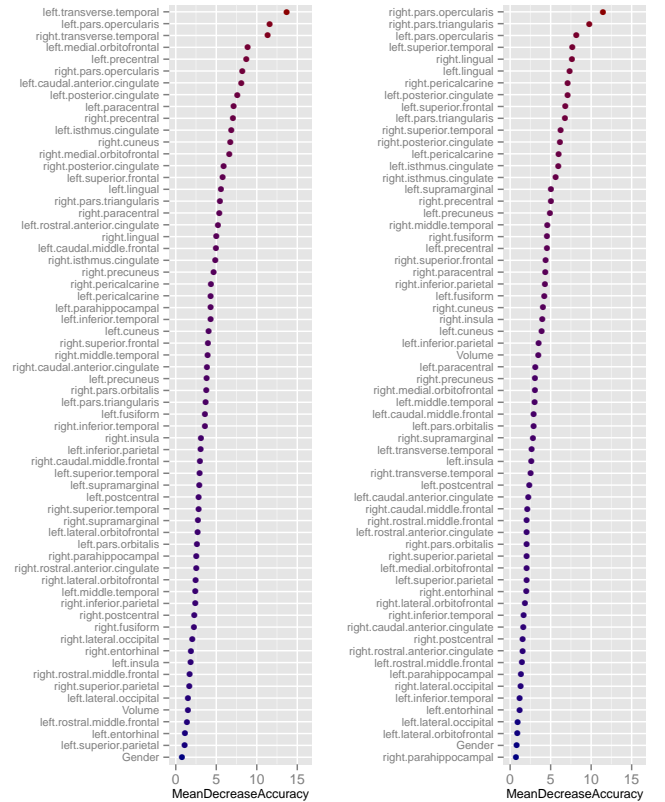


Figure 7: Regional importance random forest plots for (left) ANTs and (right) FreeSurfer using "MeanDecreaseAccuracy" ranking all model variables specified by Equation (3).

omission of a single feature or variable is tracked and averaged. Thus, those features which have the greatest decrease in mean accuracy are considered to be the most discriminative. It should be noted that correlative effects are not considered in the rankings.

3.3. Gender prediction assessment

We also performed a similar prediction assessment using gender as the regressand. The binomial generalized linear model is

$$GENDER \sim VOLUME + AGE + \sum_{i=1}^{62} T(DKT_i) \quad (5)$$

where $T(DKT_i)$ is the average thickness value in region DKT_i and $VOLUME$ is total intracranial volume. We then characterized performance using a ROC curve for both methods (see Figure 9) where we averaged over 1000 permutations. The mean area under the curve (AUC) for both methods was also quantified with values of $ANTS_{AUC} = 0.87$ and $FreeSurfer_{AUC} = 0.83$.

3.4. Computation time

All images underwent the ANTs and FreeSurfer pipeline processing using the computational cluster at the University of Virginia. Processing times varied approximately between 10–20 hours per subject for both pipelines for the entire cortical thickness estimation procedure although ANTs processing, on average, took slightly longer. Averaged over all cohorts, ANTs required 15.7 ± 2.0 hours per subject and FreeSurfer required 14.1 ± 2.9 hours per subject.

The propagation of the DKT labels to each subject using label fusion as described earlier was performed in parallel and took anywhere between 40 and 80 hours per subject for 16 serial image registrations and application of the joint label fusion algorithm [65]. For each subject, 20 atlas registrations are used to generate the labeling for that subject. Therefore to do the MALF labeling for the entire cohort, approximately $1200 \times 20 = 24000$ registrations were performed. The `antsMalFLabeling.sh` script mentioned earlier parallelizes the registration component which decreases the time for parallel computation platforms.

4. Discussion

In the absence of ground truth, we used repeatability and prediction of demographic variables to compare the ANTs and FreeSurfer cortical thickness pipelines. The only major failure was the FreeSurfer brain extraction of a single IXI subject (IXI430-IOP-0990). Also, three NKI subjects were not processed to completion with FreeSurfer (1713515, 18755434, and 2674565) and were not included in the analysis. Although researchers might quibble over processing minutiae such as the inclusion of too much (or not enough) of the meninges, we approached our evaluation using more objective criteria which concern all those engaged in this type of research. We are currently trying to develop methods to facilitate data inspection for quick quality assurance/control.

4.1. Repeatability of thickness measurements

The OASIS data set and the MMRR data set allow us to test whether the same thickness values emerge from T1-weighted MRI collected on the same subject but at different times of the day or over a time separation within a few weeks. Although the ANTs cortical thickness pipeline produced similar repeatability assessments as FreeSurfer in these data, there are many additional issues to explore with the ANTs-based framework. Pre-analysis confounds such as short-term alterations in cortical morphology due to the T1-weighted susceptibility to blood flow [23, 51, 69] and MRI acquisition parameters such as field strength, site, resolution, scanner, longitudinal variation in scanner conditions, and pulse sequence [27, 40, 30] have been evaluated with FreeSurfer which has shown good reliability under various permutations of these conditions. Although we did not explicitly investigate the repeatability performance of the ANTs framework under such effects, the relatively good performance on the large and varied data (in terms of site, field strength, scanner, and acquisition sequence) used in this study provides confidence in its robustness to a variety of imaging conditions.

ANTs and FreeSurfer cortical thickness mean reliability are correlated across all regions (Pearson correlation = 0.44). Although our thickness reliability measurements represent the compound effect of registration, segmentation, anatomical labeling, and the thickness computation algorithm, this correlation suggests that these effects are non-random. That is, reliability measurements are influenced by characteristics intrinsic to the underlying neuroanatomy as represented in approximately one millimeter resolution volumetric T1-weighted MRI. Perhaps the least reliable region is entorhinal cortex (Region 4 in Figure 5) which has relatively small volume, is challenging to distinguish from surrounding structures [46], and is also relatively thin. Spatial variation in segmentation accuracy is known to relate to a structure’s volume and tissue characteristics and this has led to a body of research on both segmentation and acquisition protocols that are optimized for specific regions. Perhaps the most substantial work in MRI has focused on temporal lobe structures including the hippocampus. Both FreeSurfer and our own group have optimized protocols to address such concerns (<http://www.hippocampalsubfields.com/>). Given caveats associated with cost vs. benefit, our current results suggest that optimized protocols may be relevant for additional cortical regions.

4.2. Voxel/vertex-based analysis

One of the limitations of our evaluation was the limitation of comparative analysis to mean ROI thickness values defined by the 62 cortical regions of the DKT atlas. Quite common in the literature, however, are point-wise (vertex- or voxel-based) analyses [e.g., 9]. The ANTs pipeline described in this work is equally applicable to such studies. The only additional requirement is the specification of the normalization template. For this work we opted for the ROI analysis to avoid potential bias issues when navigating between surface and volume representations [33]. Future work will certainly explore such analyses.

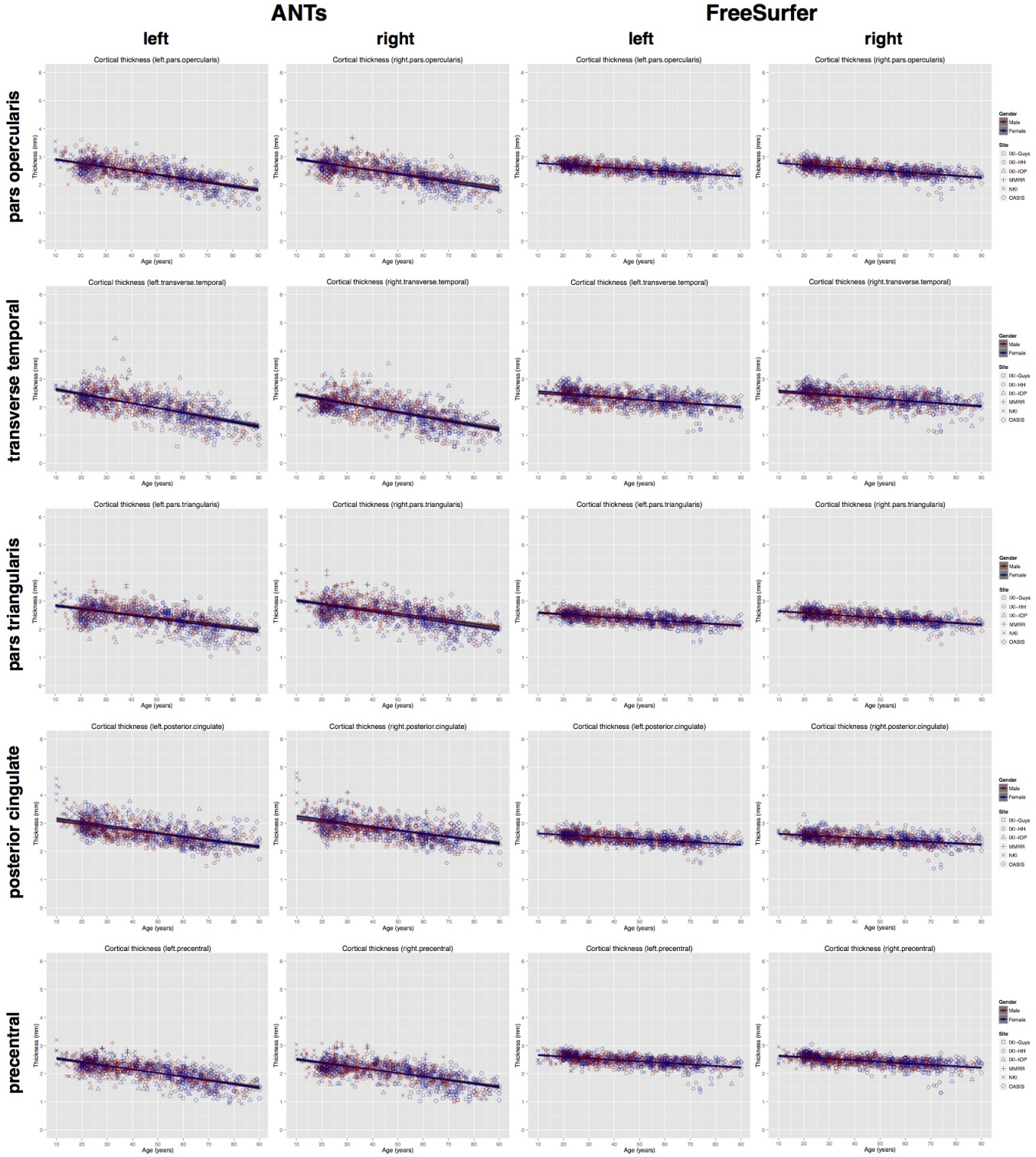


Figure 8: Age vs. thickness plots for cortical regions that are most relevant in age prediction. These are the most discriminative regions across both methods as determined by random forest importance measurements (cf Figure 7). Note that all regional plots for both ANTs and FreeSurfer are available online (see Appendix).

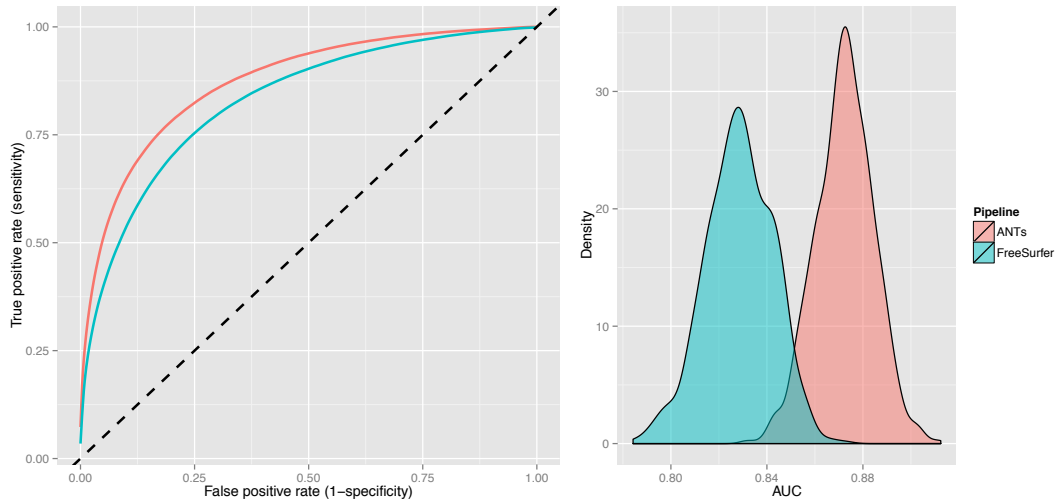


Figure 9: Average ROC curve and corresponding AUC distributions for gender prediction using ANTs and FreeSurfer thickness values. Values were averaged for 1000 permutations resulting in mean values of $ANTS_{AUC} = 0.87$ and $FreeSurfer_{AUC} = 0.83$ ($p < 10^{-16}$).

4.3. Age and gender prediction

Although repeatability between ANTs and FreeSurfer is comparable, such measures are not as useful in determining the utility of the measuring software. That is the reason we used a training and testing paradigm to evaluate how well both frameworks produce measurements capable of predicting demographics which are well-known to correlate with cortical thickness. Additionally, these demographic measures are probably some of the easiest and most reliably obtained of all possible demographic measures used for this type of assessment.

Previous research has used predictive modeling for comparing cortical thickness algorithms. For example, in [10], classification of healthy, semantic dementia, and progressive non-fluent aphasia categories using regional cortical thickness values was used to determine the predictive modeling capabilities of different cortical thickness processing protocols in 101 subjects. However, differential diagnosis of dementia [44] is not as straightforward as obtaining a subject’s age or gender and regressing that against cortical thickness; the latter constitute biological relationships that have been well-studied and reported in the literature.

For age prediction, we used both a linear model (due to its general ubiquity) and a random forest model (a non-parametric model to contrast with the linear approach) which showed overall good performance. Also, the linear and random forest models have the advantage of being interpretable—that is, the models reveal the specific predictors that are most valuable which makes comparison with previous age versus thickness assessments possible.

For example, in [28], 322 T1-weighted MRI of healthy adults with an age range of [20, 85] were used, in part, to characterize the relationship between age and cortical thickness using FreeSurfer and a similar linear modeling approach. Significant findings for age were reported in the “precentral gyrus, medial parts of the superior frontal gyrus, DMPFC, and rostral middle frontal cortex.” Based on the cortical parcellation provided by

the DKT atlas, we also saw similar strong effects in the precentral gyrus (cf Figure 7).

This study was limited to a cross-sectional investigation thus limiting extrapolations of ANTs performance to longitudinal data unlike recent FreeSurfer extensions which accommodate longitudinal data [48, 30]. Also, some users may choose to segment and register with ANTs and subsequently employ any alternative (e.g., surface-based) method for thickness estimation. Further work is needed by independent authors working on established pipelines to better compare surface-based and volume-based thickness reliability and accuracy across different populations, age ranges, and with longitudinal protocols.

4.4. Computation time

Computation time for the registration and segmentation components of the ANTs pipeline are substantial but are not significantly worse than those of FreeSurfer. It is likely that nearly as reliable results can be obtained in much less time for many of the subjects in this study. However, our interest in maximizing reliability and quality led us to employ parameters in the registration, segmentation, and bias correction that are as robust as possible to differences in head position, the presence of large deformations between template and target brains and substantial inhomogeneity or other artifacts in the image content itself.

5. Conclusions

Imaging biomarkers such as cortical thickness play an important role in neuroscience research. Extremely useful to researchers are robust software tools for generating such biomarkers. In this work we detailed our open source offering for estimating cortical thickness directly from T1 images and demonstrated its utility on a large collection of public brain data from multiple databases acquired at multiple sites. To our knowledge, this study constitutes the largest collection of cortical thickness data processed in a single study. We anticipate

that public availability of our tools and extensive tuning on the specified cohorts will prove useful to the larger research community. In this work, we only explored a portion of the potentially interesting investigations possible with these data. Since all of the data are publicly available, further work can be easily pursued by us or by other interested groups.

Appendix

Available resources are listed in Table 4 with their corresponding addresses. Examples and data for all scripts described in the manuscript are also available for download. This should enable interested researchers to duplicate the results in this work.

References

- [1] Ashburner, J., Friston, K. J., Jul 2005. Unified segmentation. *Neuroimage* 26 (3), 839–51.
- [2] Avants, B. B., Klein, A., Tustison, N. J., Woo, J., Gee, J. C., 2010. Evaluation of an open-access, automated brain extraction method on multi-site multi-disorder data. *Human Brain Mapping*.
- [3] Avants, B. B., Tustison, N. J., Song, G., Cook, P. A., Klein, A., Gee, J. C., Feb 2011. A reproducible evaluation of ANTs similarity metric performance in brain image registration. *Neuroimage* 54 (3), 2033–44.
- [4] Avants, B. B., Tustison, N. J., Wu, J., Cook, P. A., Gee, J. C., Dec 2011. An open source multivariate framework for n -tissue segmentation with evaluation on public data. *Neuroinformatics* 9 (4), 381–400.
- [5] Avants, B. B., Yushkevich, P., Pluta, J., Minkoff, D., Korczykowski, M., Detre, J., Gee, J. C., Feb 2010. The optimal template effect in hippocampus studies of diseased populations. *Neuroimage* 49 (3), 2457–66.
- [6] Bouix, S., Martin-Fernandez, M., Ungar, L., Nakamura, M., Koo, M.-S., McCarley, R. W., Shenton, M. E., Jul 2007. On evaluating brain tissue classifiers without a ground truth. *Neuroimage* 36 (4), 1207–24.
- [7] Boyes, R. G., Gunter, J. L., Frost, C., Janke, A. L., Yeatman, T., Hill, D. L. G., Bernstein, M. A., Thompson, P. M., Weiner, M. W., Schuff, N., Alexander, G. E., Killiany, R. J., DeCarli, C., Jack, C. R., Fox, N. C., ADNI Study, Feb 2008. Intensity non-uniformity correction using N3 on 3-T scanners with multichannel phased array coils. *Neuroimage* 39 (4), 1752–62.
- [8] Breiman, L., 2001. Random forests. In: *Machine Learning*, pp. 5–32.
- [9] Chung, M. K., Robbins, S. M., Dalton, K. M., Davidson, R. J., Alexander, A. L., Evans, A. C., May 2005. Cortical thickness analysis in autism with heat kernel smoothing. *Neuroimage* 25 (4), 1256–65.
- [10] Clarkson, M. J., Cardoso, M. J., Ridgway, G. R., Modat, M., Leung, K. K., Rohrer, J. D., Fox, N. C., Ourselin, S., Aug 2011. A comparison of voxel and surface based cortical thickness estimation methods. *Neuroimage* 57 (3), 856–65.
- [11] Collins, D. L., Neelin, P., Peters, T. M., Evans, A. C., 1994. Automatic 3D intersubject registration of MR volumetric data in standardized Talairach space. *J Comput Assist Tomogr* 18 (2), 192–205.
- [12] Dale, A. M., Fischl, B., Sereno, M. I., Feb 1999. Cortical surface-based analysis. i. segmentation and surface reconstruction. *Neuroimage* 9 (2), 179–94.
- [13] Das, S. R., Avants, B. B., Grossman, M., Gee, J. C., Apr 2009. Registration based cortical thickness measurement. *Neuroimage* 45 (3), 867–79.
- [14] Davatzikos, C., Bryan, N., 1996. Using a deformable surface model to obtain a shape representation of the cortex. *IEEE Trans Med Imaging* 15 (6), 785–95.
- [15] Desikan, R. S., Ségonne, F., Fischl, B., Quinn, B. T., Dickerson, B. C., Blacker, D., Buckner, R. L., Dale, A. M., Maguire, R. P., Hyman, B. T., Albert, M. S., Killiany, R. J., Jul 2006. An automated labeling system for subdividing the human cerebral cortex on mri scans into gyral based regions of interest. *Neuroimage* 31 (3), 968–80.
- [16] Dickerson, B. C., Bakkour, A., Salat, D. H., Feczko, E., Pacheco, J., Greve, D. N., Grodstein, F., Wright, C. I., Blacker, D., Rosas, H. D., Sperling, R. A., Atri, A., Growdon, J. H., Hyman, B. T., Morris, J. C., Fischl, B., Buckner, R. L., Mar 2009. The cortical signature of Alzheimer’s disease: regionally specific cortical thinning relates to symptom severity in very mild to mild AD dementia and is detectable in asymptomatic amyloid-positive individuals. *Cereb Cortex* 19 (3), 497–510.
- [17] Dogdas, B., Shattuck, D. W., Leahy, R. M., Dec 2005. Segmentation of skull and scalp in 3-D human MRI using mathematical morphology. *Hum Brain Mapp* 26 (4), 273–85.
- [18] Du, A.-T., Schuff, N., Kramer, J. H., Rosen, H. J., Gorno-Tempini, M. L., Rankin, K., Miller, B. L., Weiner, M. W., Apr 2007. Different regional patterns of cortical thinning in Alzheimer’s disease and frontotemporal dementia. *Brain* 130 (Pt 4), 1159–66.
- [19] Fischl, B., Dale, A. M., Sep 2000. Measuring the thickness of the human cerebral cortex from magnetic resonance images. *Proc Natl Acad Sci U S A* 97 (20), 11050–5.
- [20] Fischl, B., Salat, D. H., Busa, E., Albert, M., Dieterich, M., Haselgrove, C., van der Kouwe, A., Killiany, R., Kennedy, D., Klaveness, S., Montillo, A., Makris, N., Rosen, B., Dale, A. M., Jan 2002. Whole brain segmentation: automated labeling of neuroanatomical structures in the human brain. *Neuron* 33 (3), 341–55.
- [21] Fischl, B., Sereno, M. I., Dale, A. M., Feb 1999. Cortical surface-based analysis. ii: Inflation, flattening, and a surface-based coordinate system. *Neuroimage* 9 (2), 195–207.
- [22] Fischl, B., van der Kouwe, A., Destrieux, C., Halgren, E., Ségonne, F., Salat, D. H., Busa, E., Seidman, L. J., Goldstein, J., Kennedy, D., Caviness, V., Makris, N., Rosen, B., Dale, A. M., Jan 2004. Automatically parcellating the human cerebral cortex. *Cereb Cortex* 14 (1), 11–22.
- [23] Franklin, T. R., Wang, Z., Shin, J., Jagannathan, K., Suh, J. J., Detre, J. A., O’Brien, C. P., Childress, A. R., 2013. A VBM study demonstrating ‘apparent’ effects of a single dose of medication on T1-weighted MRIs. *Brain Structure and Function* 218 (1), 97–104.
- [24] Gernsbacher, M. A., Jan 2007. Presidential column: The eye of the beholder. *Observer* 20 (1).
- [25] Gronenschild, E. H. B. M., Habets, P., Jacobs, H. I. L., Mengelers, R., Rozendaal, N., van Os, J., Marcelis, M., 2012. The effects of FreeSurfer version, workstation type, and Macintosh operating system version on anatomical volume and cortical thickness measurements. *PLoS One* 7 (6), e38234.
- [26] Haier, R. J., Karama, S., Leyba, L., Jung, R. E., 2009. MRI assessment of cortical thickness and functional activity changes in adolescent girls following three months of practice on a visual-spatial task. *BMC Res Notes* 2, 174.
- [27] Han, X., Jovicich, J., Salat, D., van der Kouwe, A., Quinn, B., Czanner, S., Busa, E., Pacheco, J., Albert, M., Killiany, R., Maguire, P., Rosas, D., Makris, N., Dale, A., Dickerson, B., Fischl, B., Aug 2006. Reliability of mri-derived measurements of human cerebral cortical thickness: the effects of field strength, scanner upgrade and manufacturer. *Neuroimage* 32 (1), 180–194.
URL <http://dx.doi.org/10.1016/j.neuroimage.2006.02.051>
- [28] Hogstrom, L. J., Westlye, L. T., Walhovd, K. B., Fjell, A. M., Nov 2013. The structure of the cerebral cortex across adult life: age-related patterns of surface area, thickness, and gyrification. *Cereb Cortex* 23 (11), 2521–30.
- [29] Jones, S. E., Buchbinder, B. R., Aharon, I., Sep 2000. Three-dimensional mapping of cortical thickness using Laplace’s equation. *Hum Brain Mapp* 11 (1), 12–32.
- [30] Jovicich, J., Marizzoni, M., Sala-Llloch, R., Bosch, B., Bartrés-Faz, D., Arnold, J., Benninghoff, J., Wiltfang, J., Roccatagliata, L., Nobili, F., Hensch, T., Tränkner, A., Schönknecht, P., Leroy, M., Lopes, R., Bordet, R., Chanoine, V., Ranjeva, J.-P., Didic, M., Gros-Dagnac, H., Payoux, P., Zoccatelli, G., Alessandrini, F., Beltramello, A., Bargalló, N., Blin, O., Frisoni, G. B., The PharmaCog Consortium, May 2013. Brain morphometry reproducibility in multi-center 3T MRI studies: A comparison of cross-sectional and longitudinal segmentations. *Neuroimage*.
- [31] Jubault, T., Gagnon, J.-F., Karama, S., Pito, A., Lafontaine, A.-L., Evans, A. C., Monchi, O., Mar 2011. Patterns of cortical thickness and surface area in early Parkinson’s disease. *Neuroimage* 55 (2), 462–7.
- [32] Kim, J. S., Singh, V., Lee, J. K., Lerch, J., Ad-Dab’bagh, Y., MacDonald, D., Lee, J. M., Kim, S. I., Evans, A. C., Aug 2005. Automated 3-D extraction and evaluation of the inner and outer cortical surfaces using a Laplacian map and partial volume effect classification. *Neuroimage*

Table 4: Resources used in this work.

Packages	
ANTs	http://stnava.github.io/ANTs
FreeSurfer	http://surfer.nmr.mgh.harvard.edu
Available scripts and examples	
<code>antsBrainExtraction.sh</code>	https://github.com/ntustison/antsBrainExtractionExample
<code>antsAtroposN4.sh</code>	https://github.com/ntustison/antsAtroposN4Example
<code>antsCorticalThickness.sh</code>	https://github.com/ntustison/antsCorticalThicknessExample
<code>antsMultivariateTemplateConstruction.sh</code>	https://github.com/ntustison/TemplateBuildingExample
<code>antsMalfLabeling.sh</code>	https://github.com/ntustison/MalfLabelingExample
Analysis scripts	https://github.com/ntustison/KapowskiChronicles
Public data	
MindBoggle101	http://mindboggle.info/data.html
Cohort templates and priors	http://figshare.com/articles/ANTs_ANTsR_Brain_Templates/915436
IXI	http://biomedic.doc.ic.ac.uk/brain-development
MMRR	http://www.nitrc.org/projects/multimodal
NKI	http://fcon_1000.projects.nitrc.org
OASIS	http://www.oasis-brains.org
MICCAI 2012 Workshop on Multi-Atlas Labeling	https://masi.vuse.vanderbilt.edu/workshop2012/index.php

- 27 (1), 210–21.
- [33] Klein, A., Ghosh, S. S., Avants, B., Yeo, B. T. T., Fischl, B., Ardekani, B., Gee, J. C., Mann, J. J., Parsey, R. V., May 2010. Evaluation of volume-based and surface-based brain image registration methods. *Neuroimage* 51 (1), 214–20.
- [34] Klein, A., Tourville, J., 2012. 101 labeled brain images and a consistent human cortical labeling protocol. *Front Neurosci* 6, 171.
- [35] Kovacevic, J., 2006. From the editor-in-chief. *IEEE Trans Imag Process* 15 (12).
- [36] Kuperberg, G. R., Broome, M. R., McGuire, P. K., David, A. S., Eddy, M., Ozawa, F., Goff, D., West, W. C., Williams, S. C. R., van der Kouwe, A. J. W., Salat, D. H., Dale, A. M., Fischl, B., Sep 2003. Regionally localized thinning of the cerebral cortex in schizophrenia. *Arch Gen Psychiatry* 60 (9), 878–88.
- [37] Landman, B. A., Huang, A. J., Gifford, A., Vikram, D. S., Lim, I. A. L., Farrell, J. A. D., Bogovic, J. A., Hua, J., Chen, M., Jarso, S., Smith, S. A., Joel, S., Mori, S., Pekar, J. J., Barker, P. B., Prince, J. L., van Zijl, P. C. M., Feb 2011. Multi-parametric neuroimaging reproducibility: a 3-t resource study. *Neuroimage* 54 (4), 2854–66.
- [38] Luders, E., Narr, K. L., Thompson, P. M., Rex, D. E., Woods, R. P., Deluca, H., Jancke, L., Toga, A. W., Apr 2006. Gender effects on cortical thickness and the influence of scaling. *Hum Brain Mapp* 27 (4), 314–24.
- [39] Luders, E., Sánchez, F. J., Tosun, D., Shattuck, D. W., Gaser, C., Vilain, E., Toga, A. W., Aug 2012. Increased cortical thickness in male-to-female transsexualism. *J Behav Brain Sci* 2 (3), 357–362.
- [40] Lüsebrink, F., Wollrab, A., Speck, O., Apr 2013. Cortical thickness determination of the human brain using high resolution 3T and 7T MRI data. *Neuroimage* 70, 122–31.
- [41] MacDonald, D., Kabani, N., Avis, D., Evans, A. C., Sep 2000. Automated 3-D extraction of inner and outer surfaces of cerebral cortex from MRI. *Neuroimage* 12 (3), 340–56.
- [42] Magnotta, V. A., Andreasen, N. C., Schultz, S. K., Harris, G., Cizadlo, T., Heckel, D., Nopoulos, P., Flaum, M., Mar 1999. Quantitative in vivo measurement of gyrification in the human brain: changes associated with aging. *Cereb Cortex* 9 (2), 151–60.
- [43] Mazziotta, J. C., Toga, A. W., Evans, A., Fox, P., Lancaster, J., Jun 1995. A probabilistic atlas of the human brain: theory and rationale for its development. The International Consortium for Brain Mapping (ICBM). *Neuroimage* 2 (2), 89–101.
- [44] Neary, D., Snowden, J., Mann, D., Nov 2005. Frontotemporal dementia. *Lancet Neurol* 4 (11), 771–80.
- [45] Otsu, N., 1979. A threshold selection method from gray-level histograms. *EEE Trans. Sys., Man., Cyber.* 9 (1), 62–66.
- [46] Price, C. C., Wood, M. F., Leonard, C. M., Towler, S., Ward, J., Montijo, H., Kellison, I., Bowers, D., Monk, T., Newcomer, J. C., Schmalfuss, I., Sep 2010. Entorhinal cortex volume in older adults: reliability and validity considerations for three published measurement protocols. *J Int Neuropsychol Soc* 16 (5), 846–55.
- [47] Raji, C. A., Ho, A. J., Parikshak, N. N., Becker, J. T., Lopez, O. L., Kuller, L. H., Hua, X., Leow, A. D., Toga, A. W., Thompson, P. M., Mar 2010. Brain structure and obesity. *Hum Brain Mapp* 31 (3), 353–64.
- [48] Reuter, M., Schmansky, N. J., Rosas, H. D., Fischl, B., Jul 2012. Within-subject template estimation for unbiased longitudinal image analysis. *Neuroimage* 61 (4), 1402–18.
- [49] Rosas, H. D., Hevelone, N. D., Zaleta, A. K., Greve, D. N., Salat, D. H., Fischl, B., Sep 2005. Regional cortical thinning in preclinical Huntington disease and its relationship to cognition. *Neurology* 65 (5), 745–7.
- [50] Rosas, H. D., Liu, A. K., Hersch, S., Glessner, M., Ferrante, R. J., Salat, D. H., van der Kouwe, A., Jenkins, B. G., Dale, A. M., Fischl, B., Mar 2002. Regional and progressive thinning of the cortical ribbon in Huntington’s disease. *Neurology* 58 (5), 695–701.
- [51] Salgado-Pineda, P., Delaveau, P., Falcon, C., Blin, O., 2006. Brain t1 intensity changes after levodopa administration in healthy subjects: a voxel-based morphometry study. *British journal of clinical pharmacology* 62 (5), 546–551.
- [52] Scott, M. L. J., Bromiley, P. A., Thacker, N. A., Hutchinson, C. E., Jackson, A., Apr 2009. A fast, model-independent method for cerebral cortical thickness estimation using MRI. *Med Image Anal* 13 (2), 269–85.
- [53] Ségonne, F., Dale, A. M., Busa, E., Glessner, M., Salat, D., Hahn, H. K., Fischl, B., Jul 2004. A hybrid approach to the skull stripping problem in MRI. *Neuroimage* 22 (3), 1060–75.
- [54] Shattuck, D. W., Sandor-Leahy, S. R., Schaper, K. A., Rottenberg, D. A., Leahy, R. M., May 2001. Magnetic resonance image tissue classification using a partial volume model. *Neuroimage* 13 (5), 856–76.
- [55] Shrout, P. E., Fleiss, J. L., Mar 1979. Intraclass correlations: uses in assessing rater reliability. *Psychol Bull* 86 (2), 420–8.
- [56] Sled, J. G., Zijdenbos, A. P., Evans, A. C., Feb 1998. A nonparametric method for automatic correction of intensity nonuniformity in MRI data. *IEEE Trans Med Imaging* 17 (1), 87–97.
- [57] Smith, S. M., 1996. Flexible filter neighbourhood designation. In: *Proc. 13th Int. Conf. on Pattern Recognition*. Vol. 1. pp. 206–212.
- [58] Smith, S. M., Nov 2002. Fast robust automated brain extraction. *Hum*

Brain Mapp 17 (3), 143–55.

- [59] Smith, S. M., Jenkinson, M., Woolrich, M. W., Beckmann, C. F., Behrens, T. E. J., Johansen-Berg, H., Bannister, P. R., De Luca, M., Drobnjak, I., Flitney, D. E., Niazy, R. K., Saunders, J., Vickers, J., Zhang, Y., De Stefano, N., Brady, J. M., Matthews, P. M., 2004. Advances in functional and structural MR image analysis and implementation as FSL. *Neuroimage* 23 Suppl 1, S208–19.
- [60] Talairach, J., Tournoux, P., 1988. Co-planar stereotaxic atlas of the human brain: 3-Dimensional proportional system—An approach to cerebral imaging. Thieme.
- [61] Tustison, N. J., Avants, B. B., Cook, P. A., Zheng, Y., Egan, A., Yushkevich, P. A., Gee, J. C., Jun 2010. N4ITK: improved N3 bias correction. *IEEE Trans Med Imaging* 29 (6), 1310–20.
- [62] Tustison, N. J., Johnson, H. J., Rohlfing, T., Klein, A., Ghosh, S. S., Ibanez, L., Avants, B., 2013. Instrumentation bias in the use and evaluation of scientific software: Recommendations for reproducible practices in the computational sciences. *Frontiers in Neuroscience* 7 (162).
- [63] Vachet, C., Hazlett, H. C., Niethammer, M., Oguz, I., Cates, J., Whitaker, R., Piven, J., Styner, M., February 2011. Group-wise automatic mesh-based analysis of cortical thickness. In: Benoit M. Dawant, D. R. H. (Ed.), *SPIE Medical Imaging: Image Processing*.
- [64] Walhovd, K. B., Storsve, A. B., Westlye, L. T., Drevon, C. A., Fjell, A. M., Nov 2013. Blood markers of fatty acids and vitamin d, cardiovascular measures, body mass index, and physical activity relate to longitudinal cortical thinning in normal aging. *Neurobiol Aging*. URL <http://dx.doi.org/10.1016/j.neurobiolaging.2013.11.011>
- [65] Wang, H., Suh, J. W., Das, S. R., Pluta, J., Craige, C., Yushkevich, P. A., 2013. Multi-atlas segmentation with joint label fusion. *IEEE Trans Pattern Analysis and Machine Intelligence*.
- [66] Ward, B. D., 1999. Intracranial segmentation. Tech. rep., Medical College of Wisconsin, <http://afni.nimh.nih.gov/pub/dist/doc/3dIntracranial.pdf>.
- [67] Wei, G., Zhang, Y., Jiang, T., Luo, J., 2011. Increased cortical thickness in sports experts: a comparison of diving players with the controls. *PLoS One* 6 (2), e17112.
- [68] Wilkinson, G. N., Rogers, C. E., 1973. Symbolic description of factorial models for analysis of variance. *Journal of the Royal Statistical Society. Series C (Applied Statistics)* 22 (3), 392–399.
- [69] Yamasue, H., Abe, O., Kasai, K., Suga, M., Iwanami, A., Yamada, H., Tochigi, M., Ohtani, T., Rogers, M. A., Sasaki, T., et al., 2007. Human brain structural change related to acute single exposure to sarin. *Annals of neurology* 61 (1), 37–46.
- [70] Yezzi, Jr, A. J., Prince, J. L., Oct 2003. An Eulerian PDE approach for computing tissue thickness. *IEEE Trans Med Imaging* 22 (10), 1332–9.
- [71] Zeng, X., Staib, L. H., Schultz, R. T., Duncan, J. S., Oct 1999. Segmentation and measurement of the cortex from 3-D MR images using coupled-surfaces propagation. *IEEE Trans Med Imaging* 18 (10), 927–37.
- [72] Zhang, Y., Brady, M., Smith, S., Jan 2001. Segmentation of brain MR images through a hidden Markov random field model and the expectation-maximization algorithm. *IEEE Trans Med Imaging* 20 (1), 45–57.