

# Participant followup rate can bias structural imaging measures in longitudinal studies

Richard Beare<sup>1,2</sup>, Gareth Ball<sup>1</sup>, Joseph Yuan-Mou Yang<sup>1,3,6,7</sup>, Chris Moran<sup>2</sup>, Velandai Srikanth<sup>2,4,5</sup>, Marc Seal<sup>1,3</sup>, and the Alzheimer's Disease Neuroimaging Initiative

<sup>1</sup>Developmental Imaging, Murdoch Children's Research Institute, Australia

<sup>2</sup>Monash University, Academic Unit, Peninsula Clinical School, Central Clinical School, Melbourne, Victoria, Australia

<sup>3</sup>Department of Paediatrics, University of Melbourne, Australia

<sup>4</sup>Peninsula Health, Department of Geriatric Medicine, Melbourne, Victoria, Australia

<sup>5</sup>University of Tasmania, Menzies Institute for Medical Research, Hobart, Australia

<sup>6</sup>Neuroscience Advanced Clinical Imaging Service (NACIS), Department of Neurosurgery, Royal Children's Hospital, Australia

<sup>7</sup>Neuroscience Research, Murdoch Children's Research Institute, Australia

Thur 11 Feb 2021

## Keywords

FreeSurfer longitudinal pipeline.

Cortical thickness

Longitudinal bias

Individual template

## Highlights

This study found that group differences in the number of acquisitions involved in construction of the individual templates used in the FreeSurfer longitudinal pipeline can cause group differences in estimated cortical thickness and brain volume.

## Abstract

Loss of participants to followup (dropout) in longitudinal studies is inevitable and leads to great difficulty in interpretation of statistical results if dropout is correlated with a study outcome or exposure. In this article we test whether there is an additional problem that must be considered in longitudinal imaging studies, namely whether dropout has an impact on the function of FreeSurfer, a popular software pipeline used to estimate important structural brain metrics. We find that the number of acquisitions per individual does have an impact on the function of the FreeSurfer longitudinal pipeline, and can induce group differences in brain metrics.

# 1 Introduction

Analysis and interpretation of data collected by longitudinal studies can be extremely difficult, especially when participant dropout is correlated with a study outcome or exposure in an indirect or unexpected way. For example, smokers are more likely to be lost to followup in longitudinal studies of aging than non-smokers and thus great care is needed when interpreting differences in rates of disease, or other measures, in those remaining in the study. Demographic, socio-economic and other risk factors that are relevant to study questions, can contribute to risk of loss of follow up in many other ways. Brain metrics derived from magnetic resonance (MR) imaging in longitudinal studies are an increasingly important part of studies of development, aging and disease. In this article we investigate whether biases attributable to dropout rates are detectable in MR structural image processing results. Such biases, if present, make interpretation of MR metrics even more complex.

The specific MR metric investigated in this analysis is the mean cortical thickness estimate generated by the FreeSurfer suite (Dale, Fischl, and Sereno 1999; Fischl and Dale 2000). FreeSurfer is a widely neuroimaging package (freely available from <http://surfer.nmr.mgh.harvard.edu/>) capable of many forms of analysis and one of the core functions it performs is segmentation and characterization, including vertex-wise thickness estimation, of the cortex. FreeSurfer also includes a longitudinal processing pipeline (Reuter et al. 2012) in which all acquisitions for a study participant are combined to create a participant-specific template that contributes to the estimation of cortical thickness from each acquisition. Great care has been taken in the development of the longitudinal pipeline to ensure that the procedure is not biased by any of the acquisitions. The template approach improves processing robustness in the presence of artifacts that might be present in some acquisitions, but is a potential source of bias at the study level via the number of acquisitions contributing to the template, which will be referred to as “template depth” in this article. There is potential for such biases to lead to over or under estimation of group effects if the number of acquisitions is correlated with group membership and the number of acquisitions influences the thickness estimate.

This study explores the issue of template depth causing biases in cortical thickness and supratentorial volume estimates produced by the longitudinal FreeSurfer pipeline, and trajectories of both, via simulated atrophy of scans from the Alzheimer’s Disease Neuroimaging Initiative (ADNI).

## 2 Code and data availability

All data used in this article is derived from the ADNI study and may be obtained from the referenced sources. Tools used to generate and process images and perform statistical analysis are publicly available and details are provided in Methods, below.

## 3 Methods

### 3.1 MRI data

Two subsets of data from the ADNI1 cohort were selected. *Set A* consisted of participants with MRI acquisitions at the 60 month followup. *Set B* consisted of 50 randomly selected participants from the normal cognition group with baseline scans who were not in Set A. Set B was used as a basis for simulated atrophy. ADNI was launched in 2003 with the primary aim of determining whether serial imaging and other biomarkers could be combined to measure progression of mild cognitive impairment and early Alzheimer’s disease. Written, informed consent was obtained from all participants. Full details of ADNI study design, ethics approvals and acquisition protocols are described at <http://www.adni-info.org>. This study used the 1.5T MP-RAGE acquisitions from ADNI1 that had been preprocessed using the image correction procedures implemented by ADNI.

### 3.2 FreeSurfer processing of Set A

Five estimates of baseline cortical thickness were obtained for each participant by carrying out five separate runs of the FreeSurfer longitudinal pipeline. Each run used a template with a different depth - i.e. constructed using a different number of acquisitions (ranging from 1 to 5).

Cortical thickness trajectories of Set A are unknown and cannot be assumed to be linear or consistent. Thus analysis of Set A was restricted to the effect of template depth on baseline thickness estimated using the FreeSurfer longitudinal pipeline.

FreeSurfer v6.0 was used for all experiments without manual editing.

### 3.3 Simulated MRI data

Scans from Set B were used as input to the Simul@trophy tool (Khanal et al. 2016; Khanal, Ayache, and Penneec 2016) (available from <https://inria-asclepios.github.io/simul-atrophy/>) to create longitudinal data with two forms of simulated atrophy:

1. Atrophy of cortical grey matter only.
2. Atrophy of white matter only.

Five time series, each containing six time-points, were created for each baseline scan for each class of atrophy. A different rate of atrophy, ranging from 0 to 4% (settings 00, 01, 02, 03 and 04) was used for each time series. Masks used to define brain regions experiencing atrophy (i.e a cortical mask and a white matter mask) were derived from the FreeSurfer parcellation of the original baseline scans. The white matter mask was eroded by 1 voxel to avoid leakage of simulated atrophy into the cortex.

The simulated data was thus constructed to have a known and constant rate of atrophy. All simulation data was derived from the baseline scan of each participant and was thus of similar quality to the baseline scan. The images produced with atrophy of 0 were the same for both atrophy patterns.

### 3.4 FreeSurfer processing of simulated data.

Six longitudinal FreeSurfer analyses, each with a different template depth, were carried out for each time series. All of the acquisitions contributing each template were processed longitudinally, rather than only the baseline scan, as was done with Set A.

Longitudinal FreeSurfer processing begins with cross sectional processing of all acquisitions and thus standard, cross-sectional, estimates of thickness were available as a result of this procedure. The mean cortical thickness and supratentorial volume estimates were used in the analysis.

## 4 Statistical Analysis

Trajectories of cortical thickness and supratentorial volume estimated with cross sectional processing for scans with simulated atrophy were analysed separately for each atrophy level and type via random-intercept multi-level models, with participant ID as the random effect. These models provide a group estimate of slope to validate the performance of simulation.

Effects of the longitudinal pipeline were investigated by comparing mean thickness and trajectories estimated using templates of different depth by using template depth as a group main effect in a model predicting thickness from time point, for example the trajectories estimated from a template depth of 3 were compared to those estimated from template depth 6. All possible template depth pairs (10) were explored for each atrophy type and level. The equivalent models based on cross sectional FreeSurfer estimates was estimated for comparison purposes. Template depth was also used as a continuous main effect and continuous interaction

term to investigate whether it was a predictor of mean thickness or trajectory. No correction for multiple comparisons was performed for the depth pair analysis as the intention was to explore the range of scenarios that may occur in practice. A p-value of below 0.05 was taken to indicate a group difference.

A similar multi-level structure was used to test whether the estimate of thickness at timepoint 1 in Set A was associated with template depth.

Cohen's F was used to characterize effect size.

Analysis was carried out in R version 4.0 (R Core Team 2020) and the lme4 package (Bates et al. 2015). Model formulae are provided in Supplementary Material Section A.2

## 5 Results

### 5.1 Demographics

The random sample of 50 (26 female) cognitively normal ADNI participants had mean age of 76 (sd=4.9). There were 57 participants (19 female) with followup scans at 60 months. Mean age was 75.2 (sd=5.5), 31 were cognitively normal at baseline and 26 had a diagnosis of mild cognitive impairment.

### 5.2 Visualization of cross sectional processing results

#### 5.2.1 Set A - ADNI acquisitions to 60 months

The trajectories for cross sectional processing of participants with acquisitions at 60 months is shown in Figure 1 and Supplementary Figure 15.

#### 5.2.2 Set B - Simulated atrophy of randomly selected, cognitively normal, ADNI participants

**5.2.2.1 Simulated cortical atrophy** Trajectories of supratentorial volume and mean cortical thickness produced by cross sectional FreeSurfer processing in the presence of simulated cortical atrophy are shown in Figure 16 and Supplementary Figure 2. Each panel illustrates a different atrophy setting.

An example of the surfaces extracted from two time points for the same participant with simulated cortical atrophy is illustrated in Supplementary Figure 14.

**5.2.2.2 Simulated WM atrophy** Trajectories of supratentorial volume and mean cortical thickness produced by cross sectional FreeSurfer processing in the presence of simulated white matter atrophy are shown in Figures 17 and Supplementary Figure 3.

### 5.3 Visualization of Set B longitudinal processing results

#### 5.3.1 Simulated cortical atrophy

Trajectories of mean cortical thickness and supratentorial volume produced by longitudinal FreeSurfer processing in the presence of simulated cortical atrophy are shown in Figure 4 and Supplementary Figure 18.

#### 5.3.2 Simulated WM atrophy

Trajectories of mean cortical thickness and supratentorial volume produced by longitudinal FreeSurfer processing in the presence of simulated WM atrophy are shown in Figure 5 and Supplementary Figure 19.

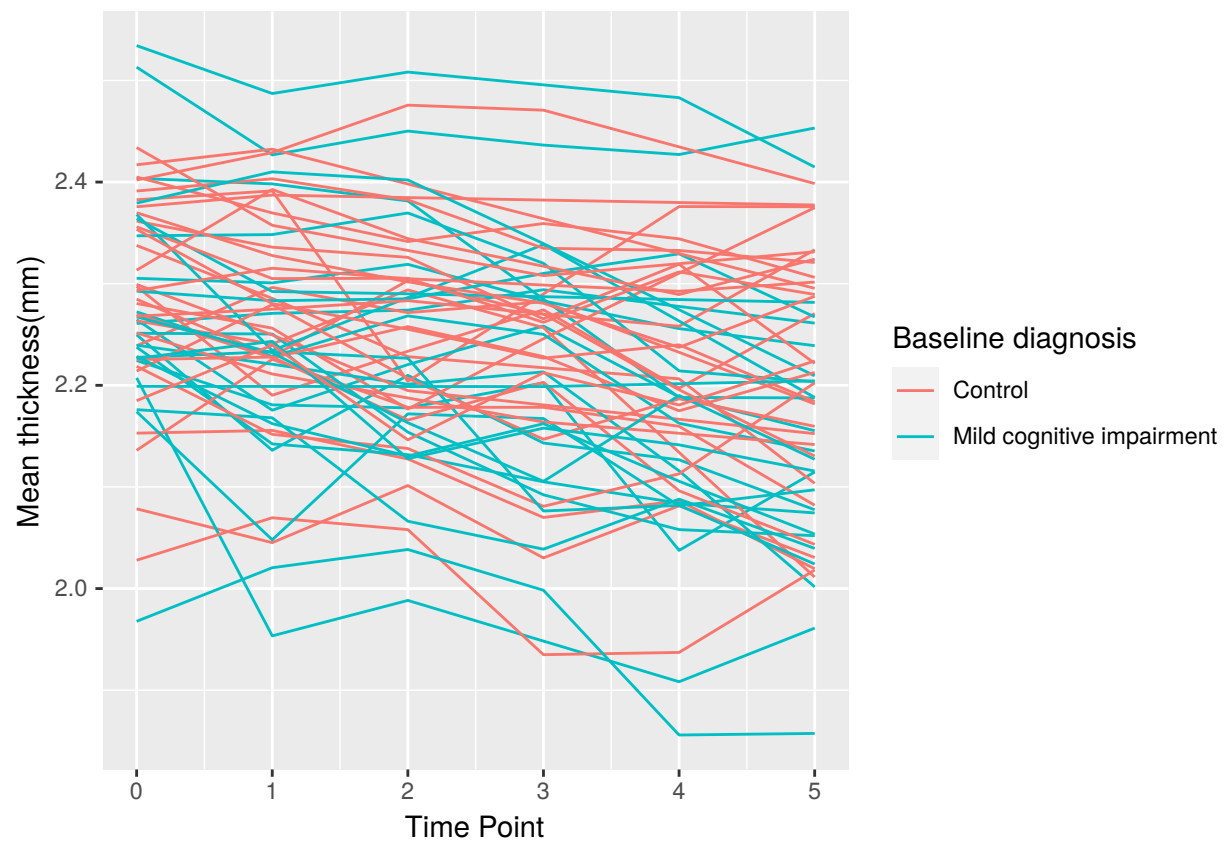


Figure 1: Mean cortical thickness estimated with cross sectional FreeSurfer for ADNI participants with scans at 60 months

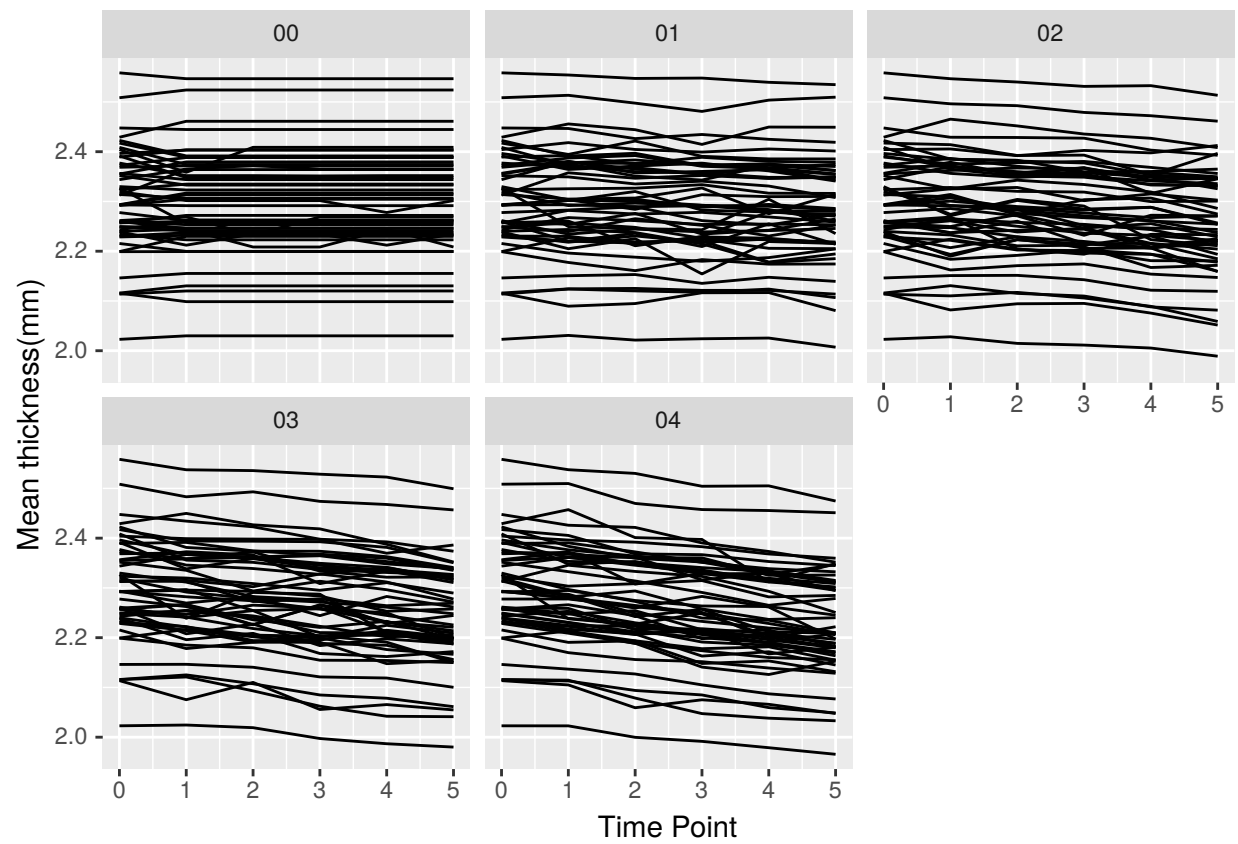


Figure 2: Mean cortical thickness estimated with cross sectional FreeSurfer in presence of simulated cortical atrophy. Each panel corresponds to a different level of simulated atrophy.

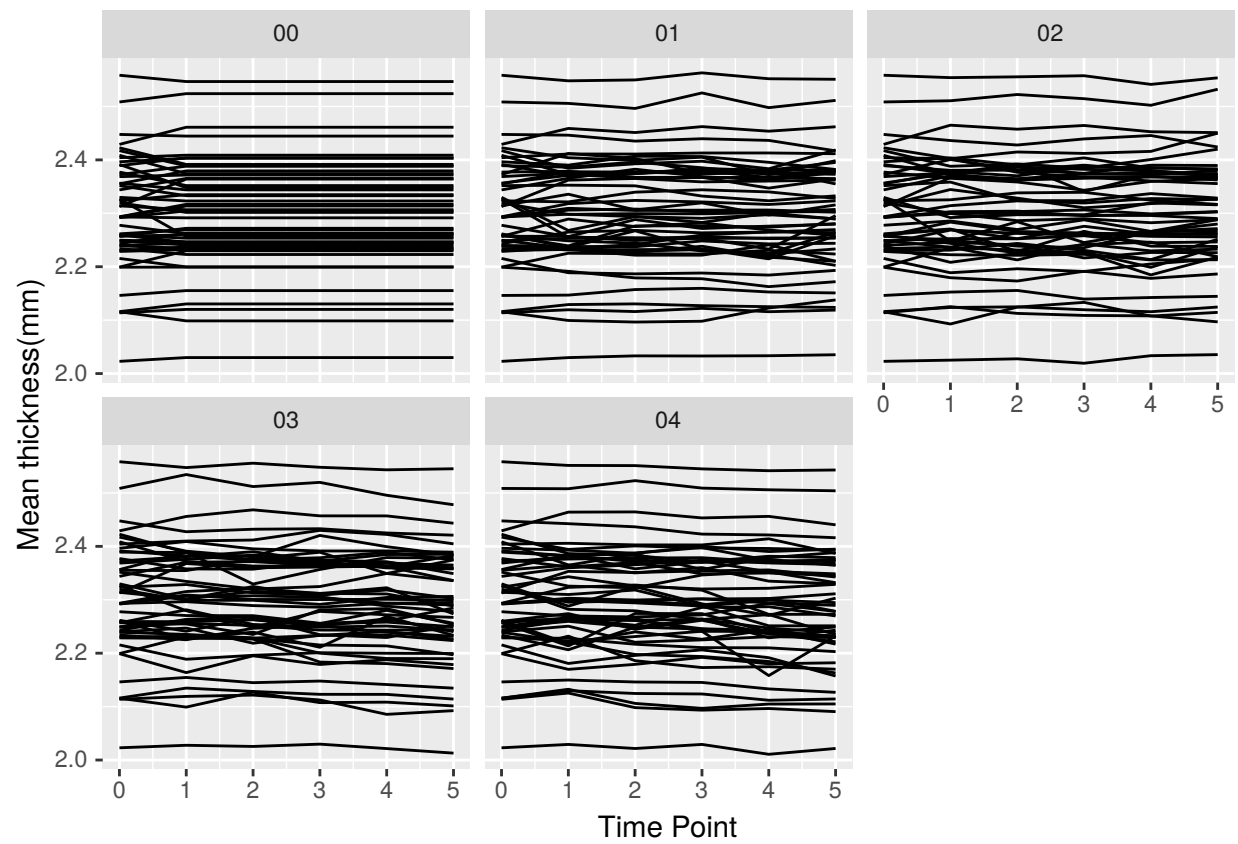


Figure 3: Mean cortical thickness estimated with cross sectional FreeSurfer in presence of simulated WM atrophy. Each panel corresponds to a different level of simulated atrophy.

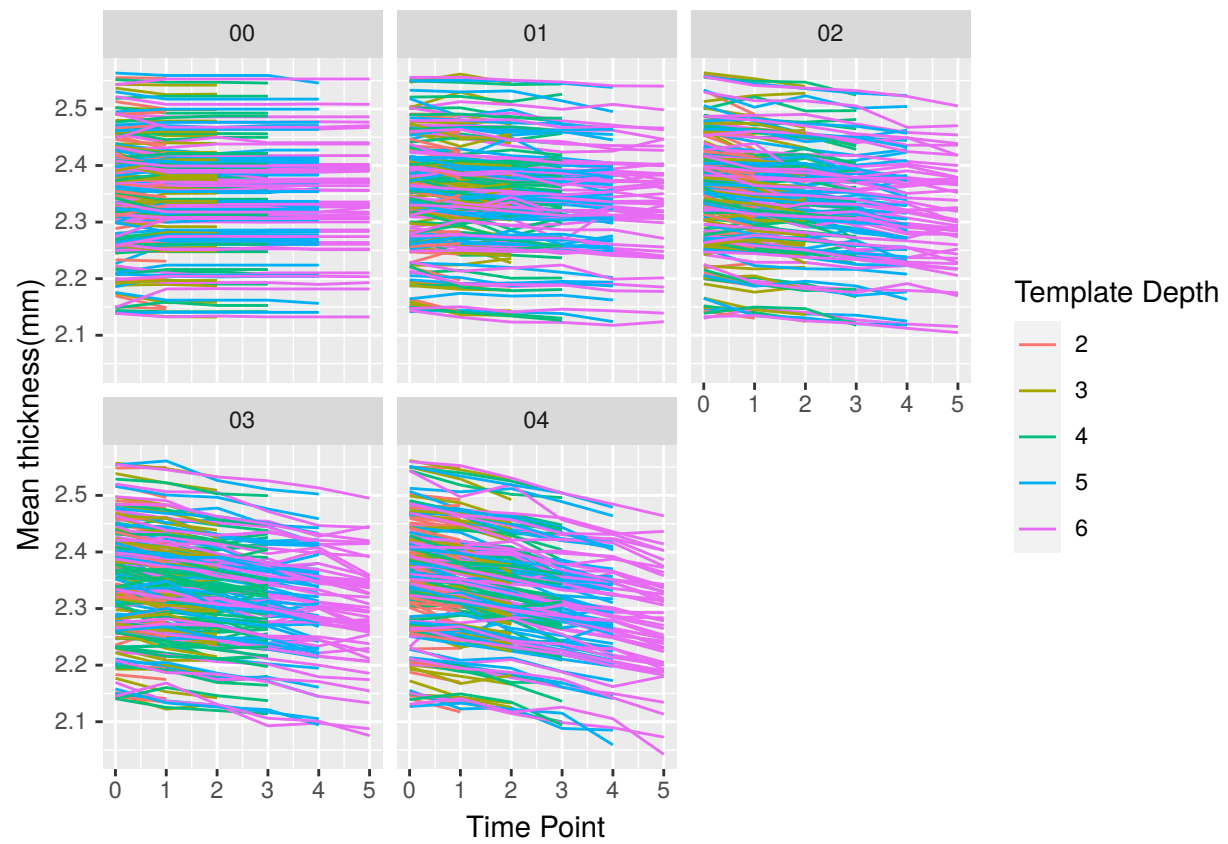


Figure 4: Cortical thickness estimated with longitudinal pipeline and varying numbers in template in presence of simulated cortical atrophy. Each panel corresponds to a different level of simulated atrophy.



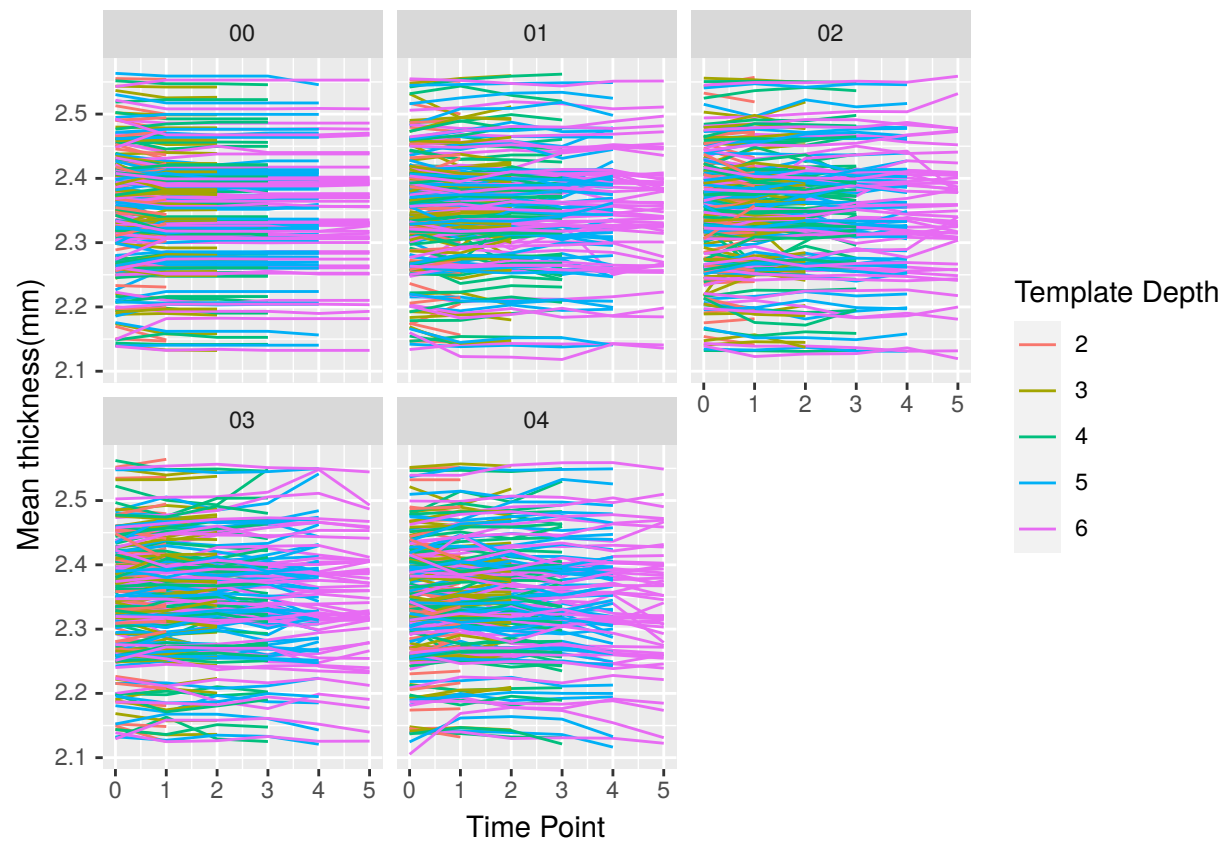


Figure 5: Cortical thickness estimated with longitudinal FreeSurfer pipeline and varying numbers in template in presence of simulated WM atrophy. Each panel corresponds to a different level of simulated atrophy.

## 5.4 Statistical Analysis

### 5.4.1 Comparison of template depth groupings

Group comparisons are presented in tabular form in Figures 6 to 13. Each figure panel corresponds to a different atrophy level. Tests involving the cross sectional pipeline are in the upper left of each panel while the lower right contains results of tests using data from the longitudinal pipeline. Each cell contains the p-value, beta coefficient and Cohen's f statistic resulting from a single group comparison. Red text indicates p-value below 0.05 and bold and italic text indicate negative and positive beta coefficients respectively.

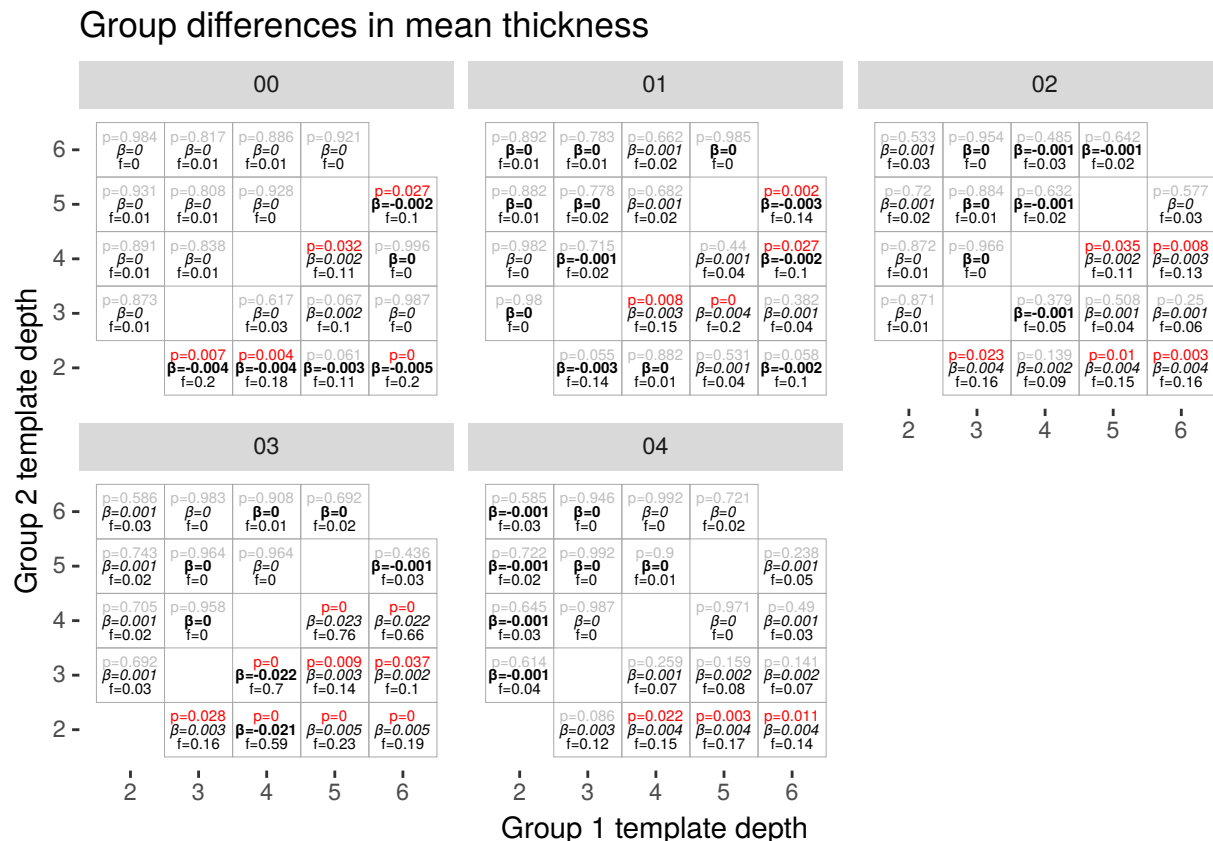


Figure 6: Group differences in cortical thickness for simulated GM atrophy.

**5.4.1.1 Cortical thickness** No template depth group differences were detected in mean GM thickness or trajectory (minimum p value 0.07) for either cortical atrophy or WM atrophy and cross section FreeSurfer processing

Differences in mean cortical thickness between template depth groups were detected ( $p < 0.05$ ) in 26 of the 50 comparisons (10 group pairs, 5 levels of atrophy) of simulated cortical atrophy processed the longitudinal pipeline (Figure 6).

The minimum p-value ( $p=1e-41$ ) occurred when comparing template depths 4 and 5 at atrophy level 03 and corresponded to a difference in means of 0.027mm (stderr 0.00149) in 2.33 (Cohen's  $F = 0.76$ ).

The sign of group differences was not consistent, i.e. more acquisitions did not imply thicker (or thinner estimates)

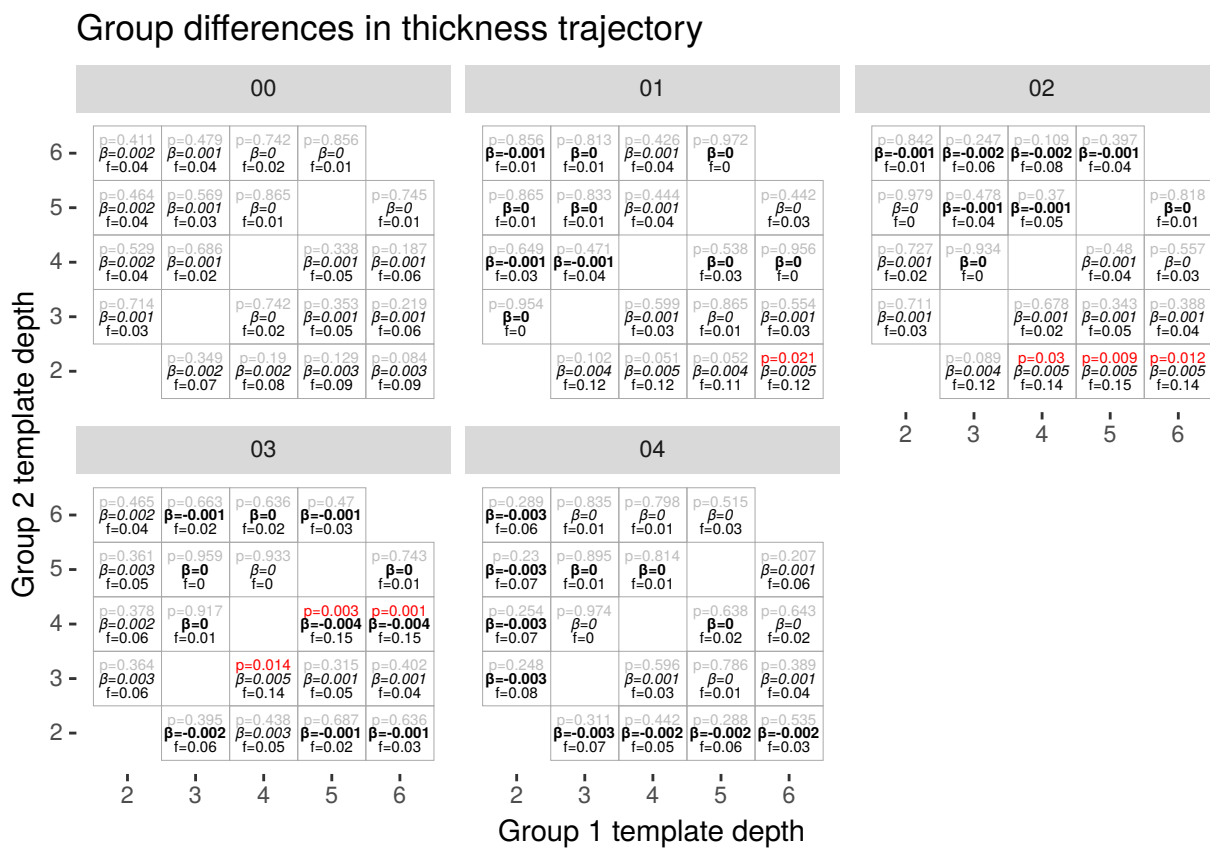


Figure 7: Group differences in trajectory of cortical thickness for simulated GM atrophy.

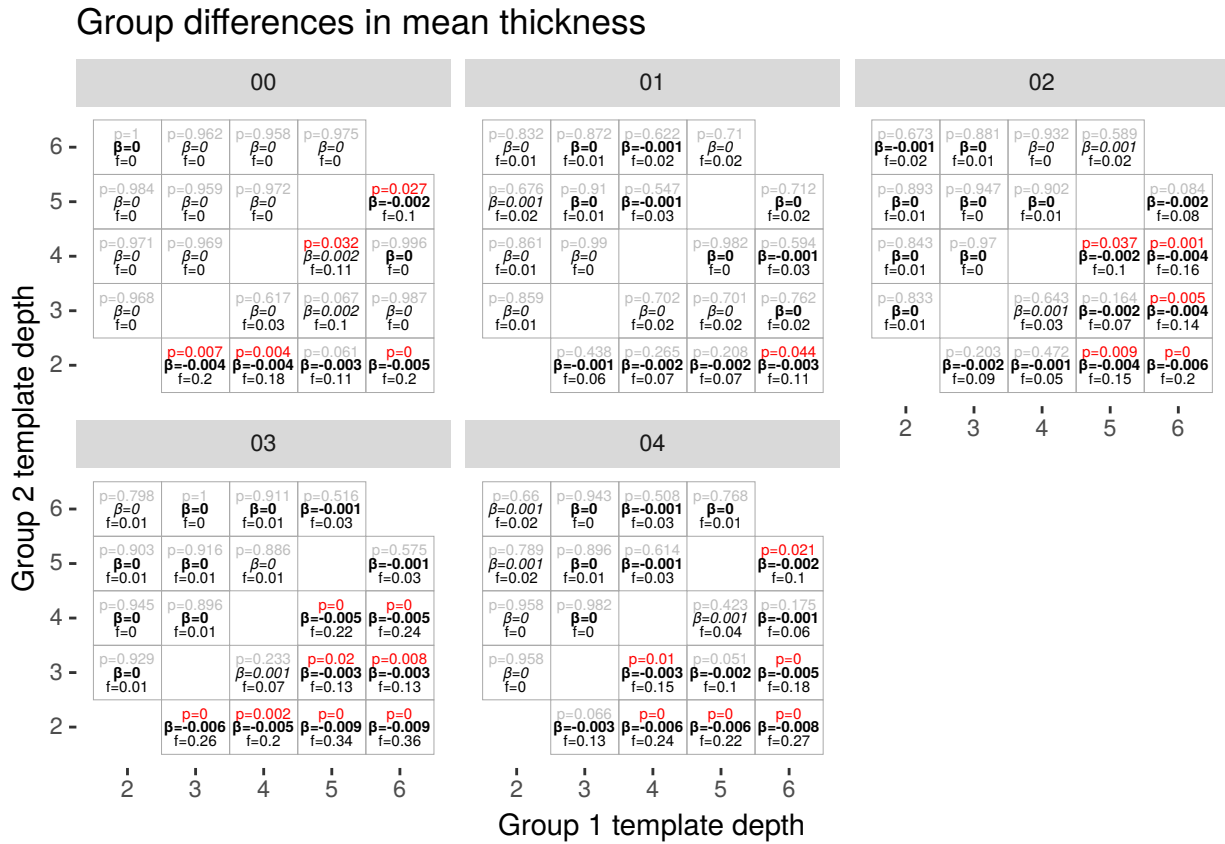


Figure 8: Group differences in cortical thickness for simulated WM atrophy.

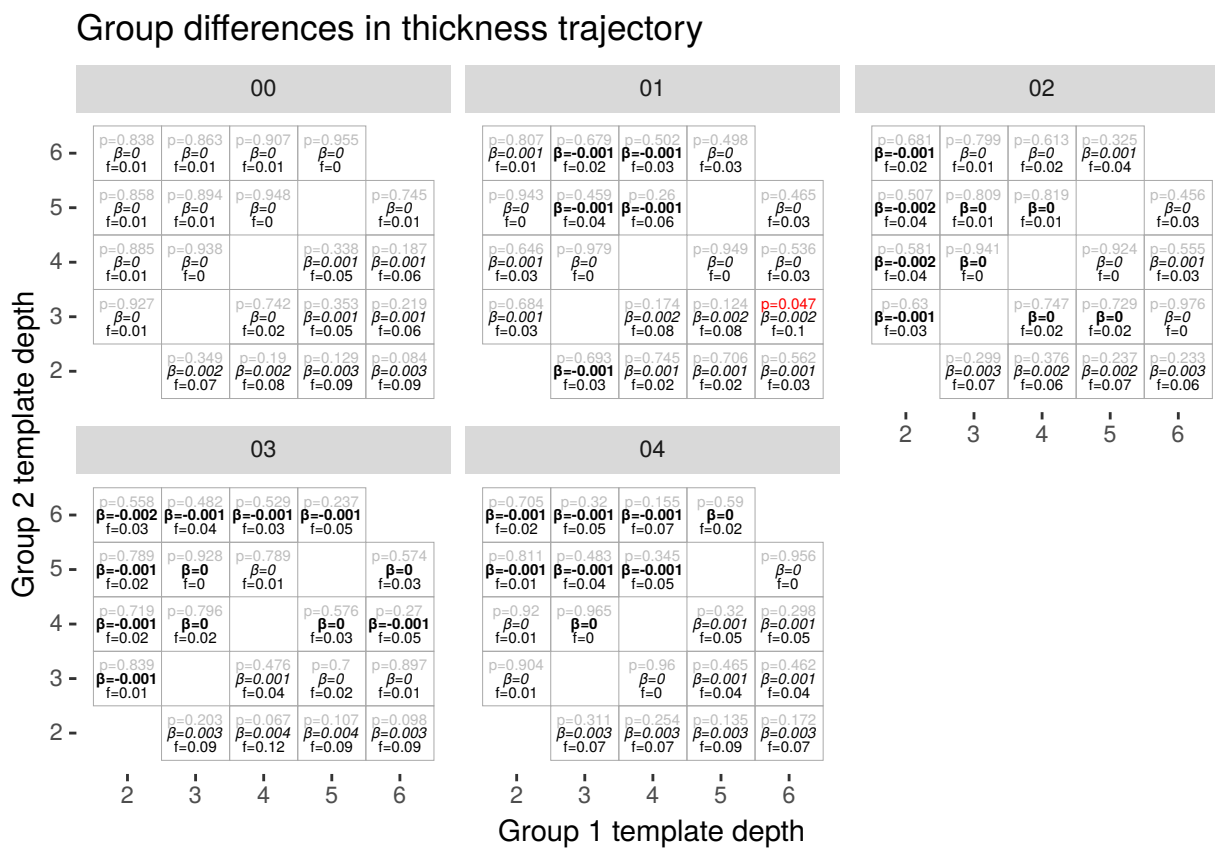


Figure 9: Group differences in trajectory of cortical thickness for simulated WM atrophy.

Differences in trajectories of cortical thickness between template groups were observed in 7 of the 50 comparisons (Figure 7).

The minimum p-value ( $p=0.00012$ ) occurred when comparing template depths 4 and 6 at atrophy level 03. Group difference in trajectory =  $-0.00374\text{mm/TP}$  and the slope of base group =  $-0.008$  (Cohen's  $F = 0.15$ ). (almost 50% difference)

The sign of group differences in trajectory of cortical thickness was not consistent.

Difference in mean cortical thickness between template groupings were detected in 25 of the 50 comparisons of simulated WM atrophy and the longitudinal processing pipeline (Figure 8). Minimum p value =  $1e-11$ .

The minimum p-value ( $1e-11$ ) occurred when comparing template depths 2 and 6 at atrophy level 03. Group difference in mean =  $-0.00882\text{mm}$ (stderr 0.00131) in 2.35 (Cohen's  $F = 0.36$ )

The sign of the group difference was negative in 24 of 25 cases, meaning a lower thickness estimate for deeper templates.

Only 1 difference in trajectory detected for simulated WM atrophy ( $p = 0.047$ ) occurred when comparing templates depth 3 and 6 at atrophy level 01 (Figure 9). Group difference in trajectory =  $0.00229$  and the slope of base group =  $-0.00177$  (Cohen's  $F = 0.10$ ).

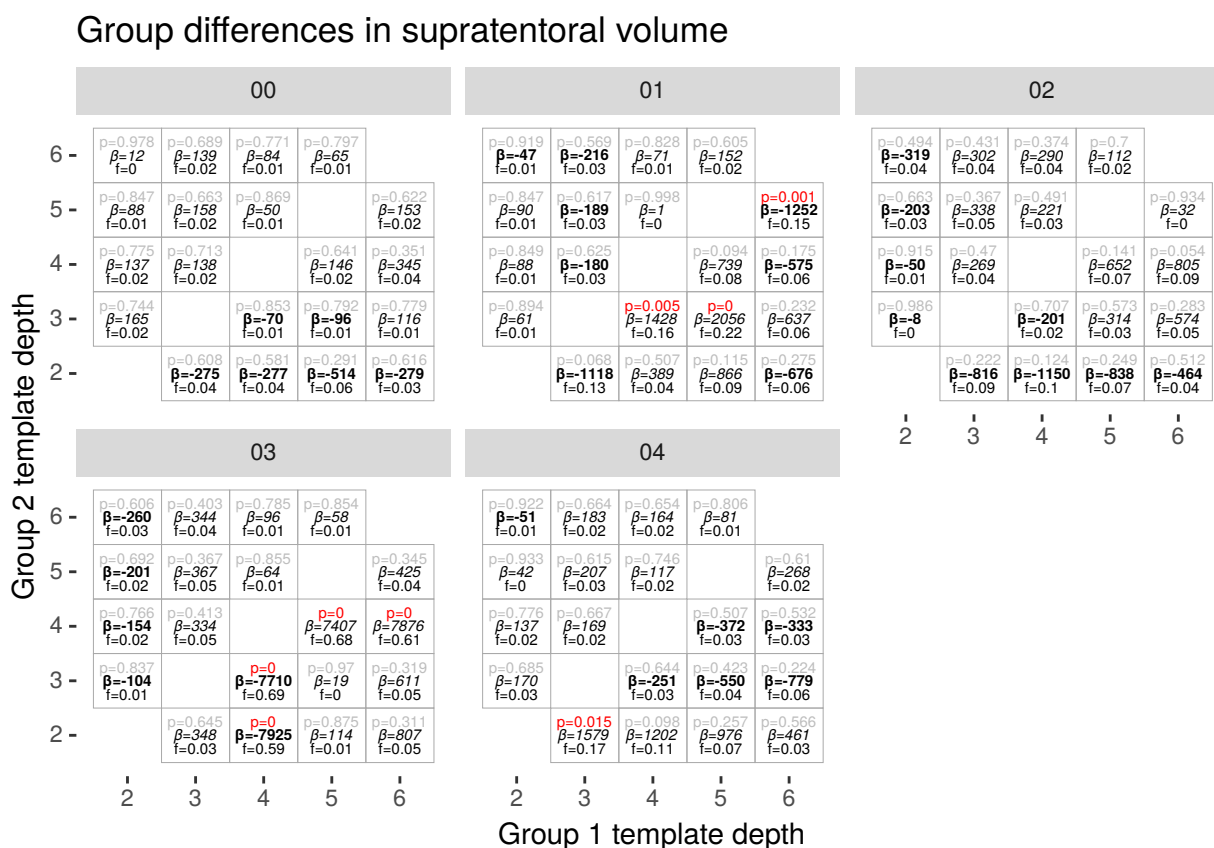


Figure 10: Group differences in supratentorial volume for simulated GM atrophy.

**5.4.1.2 Supratentorial volume** Several significant differences in trajectory of supratentorial volume estimated from cross sectional processing were detected: 3 vs 5 time points ( $p=0.015$ ) and 3 vs 6 time points ( $p=0.005$ ) for cortical atrophy level 02, and 3 vs 5 time points ( $p=0.046$ ) and 3 vs 6 time points ( $p=0.027$ ) for cortical atrophy level 03. See upper left triangle of Figure 11.

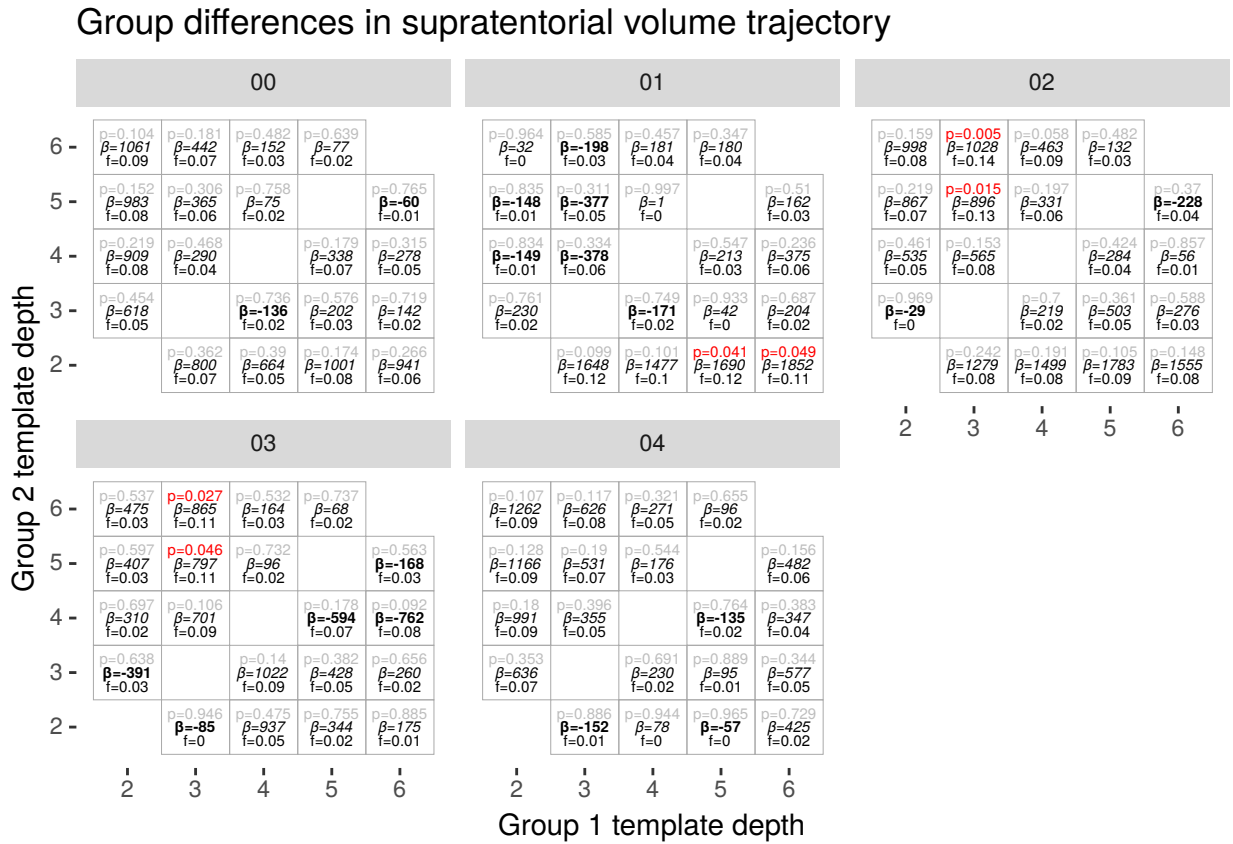


Figure 11: Group differences in trajectory of supratentorial volume for simulated GM atrophy.

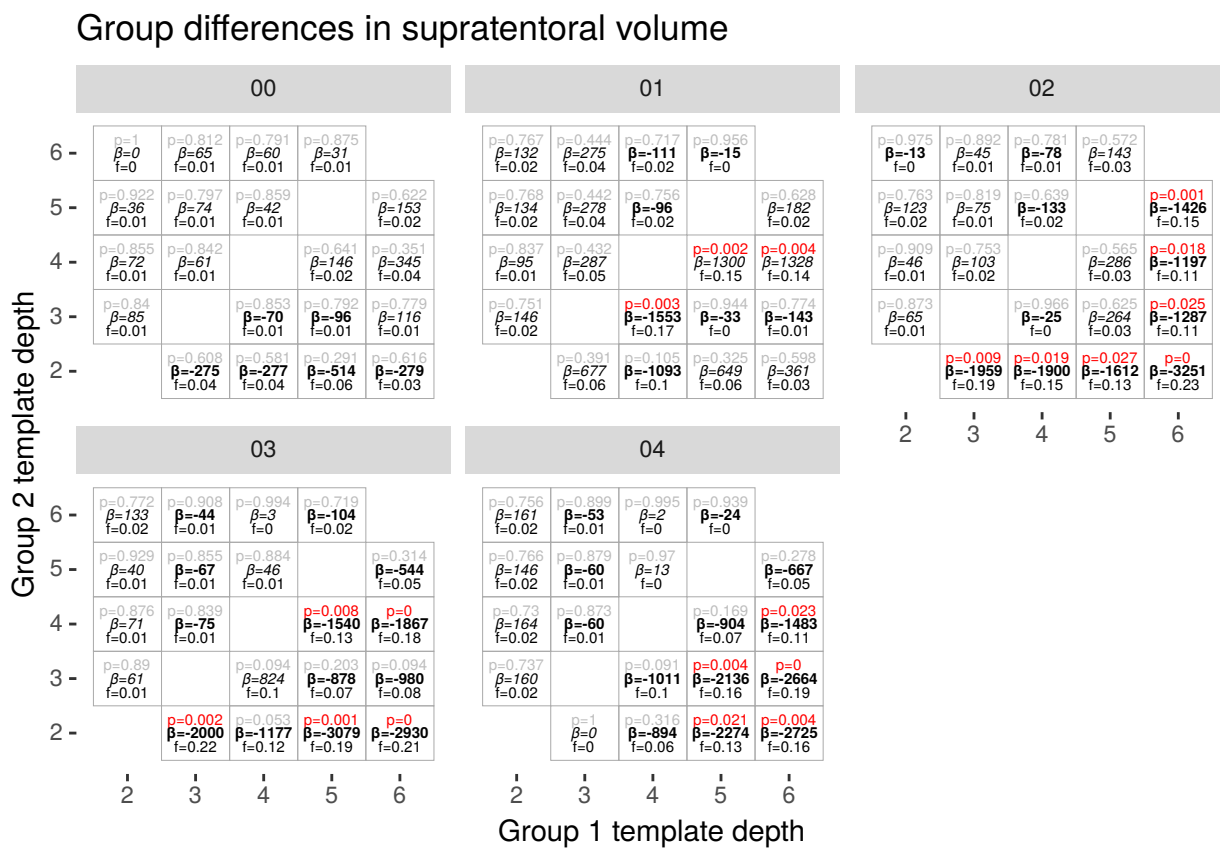


Figure 12: Group differences in supratentorial volume for simulated WM atrophy.



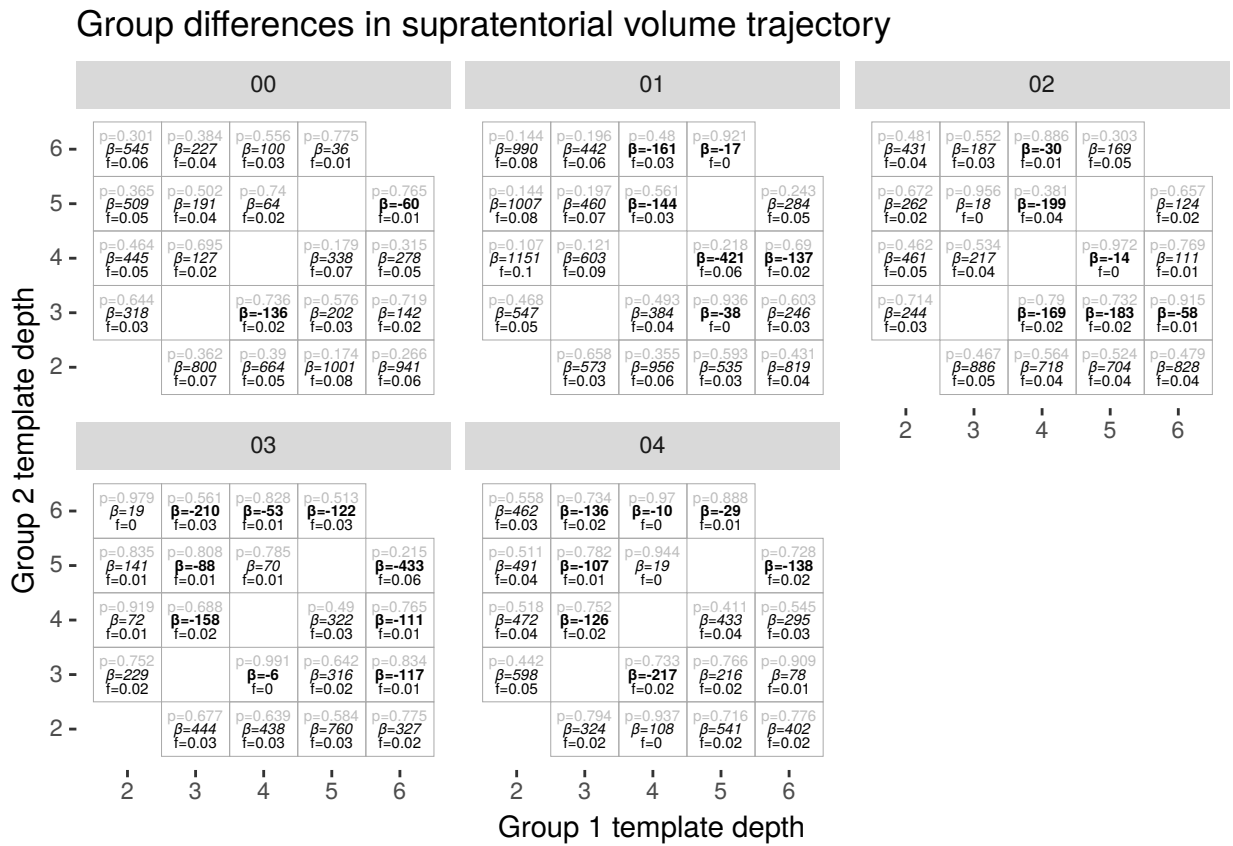


Figure 13: Group differences in trajectory of supratentorial volume for simulated WM atrophy.

Differences in mean supratentorial volume between template depth groups were detected ( $p < 0.05$ ) in 8 of the 50 comparisons of simulated cortical atrophy processed the longitudinal pipeline (Figure 10).

The minimum p-value ( $p=2e-34$ ) occurred when comparing template depths 4 and 5 at atrophy level 03. The group difference in mean was 7407(stderr 549) in 898576 (Cohen's F = 0.68).

The sign of group differences was not consistent.

Marginally significant differences in trajectory of supratentorial volume were detected in 2 of 50 comparisons of simulated cortical atrophy (Figure 11). The minimum p-value ( $p = 0.041$ ) occurred when comparing template depths 2 and 5 at atrophy level 01. The group difference was 1690(stderr 825) and the slope of base group = -2785. (Cohen's F = 0.12)

Sign of group differences was positive for both cases.

Difference in mean supratentorial volume between template depth groups were detected in 20 of the 50 comparisons of simulated WM atrophy and the longitudinal processing pipeline (Figure 12).

The minimum p value ( $p=2.97e-5$ ) occurred when comparing template depths 2 and 6 at atrophy level 02. The group difference in mean was -3251 (stderr 768) in 906529 (Cohen's F = 0.23).

The sign was negative in 18 of the 20 cases, indicating lower volume for the deeper template group.

No differences in trajectory between template groups were detected for simulated WM atrophy (Figure 13).

## 5.4.2 Template depth as a predictor

**5.4.2.1 Set A** Template depth was a significant predictor of mean cortical thickness in Set A ( $p=1e-9$ ), coefficient 0.0026(stderr 0.0004) and a Cohen's F statistic of 0.42.

### 5.4.2.2 Set B

**5.4.2.2.1 Cortical atrophy** Template depth was a predictor of mean cortical thickness at atrophy levels of 02 ( $p=0.002$ , Cohen's F = 0.1), 03( $p<0.00001$ , Cohen's F = 0.15) and 04( $p=0.002$ , Cohen's F = 0.10) and of supratentorial volume for atrophy level 03 ( $p=0.002$ , Cohen's F = 0.10). It was also a predictor of cortical thickness trajectory ( $p=0.004$ , Cohen's F = 0.09) and supratentorial volume trajectory ( $p=0.0001$ , Cohen's F = 0.13) at an atrophy level of 03. Coefficients were positive - deeper template corresponding to thicker cortex or larger volume.

**5.4.2.2.2 WM atrophy** Template depth was a predictor of mean cortical thickness and supratentorial volume for atrophy levels 02, 03 and 04 ( $p < 0.0001$ , Cohen's F = 0.1, 0.15, 0.1). Coefficients were negative.

Template depth was positively associated with cortical thickness trajectory for atrophy level 01 ( $p=0.02$ ).

## 6 Discussion

In this study we have investigated the impact of *template depth*, the number acquisitions used to construct individual templates, on measures produced by the FreeSurfer longitudinal pipeline. It is rare for longitudinal neuroimaging studies to have a balanced design, with the same number of acquisitions for every participant. Participants may drop out of studies for many reasons, some of which may be directly or indirectly related to study questions or exposures. This is particularly common in studies of aging or chronic/progressive disease that may limit participants' ability to complete demanding study protocols, leading to lower rates of follow-up in the sicker or frailer portions of the study cohort. Other biases in followup-rates caused by, for example socio-economic-status, have the potential to affect studies of typical development or healthy

controls. One can imagine diverse potential causes - for example a group-dependent frequency of dental braces could lead to a group difference in number of usable scans acquired during adolescence while differing rates of employment may lead to increased or decreased attendance, depending on study type. Difference in followup rates of study groups will lead to group differences in template depth. It is therefore important to understand the effect of template depth on derived measures.

This article has investigated the problem via two sets of data derived from the ADNI study. We created several simulated longitudinal datasets derived from real scans using a previously described atrophy simulator (Set B). The datasets were constructed using atrophy rates that remained constant for the course of the simulation (i.e constant over time) and with the same rate and pattern for each individual. We investigated two patterns of atrophy, one cortical specific and one white matter specific. The visualization of cross sectional processing results suggests that simul@tropy, with a cortical atrophy pattern, produced rates of cortical thinning that remained constant over the 6 simulated time points and were linearly related to the atrophy setting (can provide a plot of this). The WM atrophy pattern did not produce detectable cortical thinning in cross sectional processing results at lower atrophy levels, but did begin to influence cortical thickness measures in a detectable manner at atrophy settings of 03 and 04. This is likely to be due to the spatial regularization required to produce realistic physical models and partial volume effects at the GM-CSF boundary. The highest rate of cortical thinning detected with the WM atrophy pattern was slightly lower than that detected with the lowest nonzero cortical atrophy pattern.

These simulated longitudinal acquisitions represent an ideal data set for FreeSurfer processing, as the non brain tissue was not modified by the atrophy simulator and the contrast to noise and bias field inhomogeneity effects remained constant within a participant.

Significant differences in mean cortical thickness and trajectories between template-depth groups were more frequently detected in the results of the FreeSurfer longitudinal pipeline than the standard cross-sectional pipeline. The differences that were detected in the cross sectional estimates were marginally significant and involved supratentorial volume rather than cortical thickness.

Differences in mean thickness were detected more frequently than differences in thickness trajectory. Signs of mean thickness difference appear to depend on the type of atrophy, with deeper templates associated with lower thickness in the white matter atrophy simulation but no consistent relationship present in the cortical atrophy simulation. Only one marginally significant difference in thickness trajectory was detected in the white matter atrophy simulation - but this occurred at an atrophy level at which the thickness trajectory was zero in the cross sectional data. Differences in trajectory were detected in some comparisons of cortical atrophy, and the signs of difference were not consistent.

These results are difficult to interpret. It is clear that group differences in mean cortical thickness, trajectory of cortical thickness, mean supratentorial volume and trajectory of supratentorial volume can be induced via the choice of template depth used in longitudinal image processing. Differences in mean thickness appear to be more common than differences in trajectory, and were more likely to occur in situations with white matter atrophy, supported by our finding that template depth was a significant predictor of mean thickness in the white matter atrophy experiment. We also detected group differences in our zero atrophy cases. Supratentorial volume measures exhibited a similar behaviour.

The largest difference in cortical thickness that we detected as a result of group differences in template depth was 0.027mm. Accuracy and test-retest reliability of the FreeSurfer, for the purposes of sample size calculations, have been investigated previously (Liem et al. 2015; Pardoe et al. 2013). These articles examined the problem in the context of a vertex-wise analysis with varying levels of smoothing and a mix of acquisition models, with (Liem et al. 2015) using repeated scans and the longitudinal FreeSurfer pipeline while (Pardoe et al. 2013) cross sectional studies and group comparisons. The analysis by (Liem et al. 2015) suggests that a difference of 0.053mm can be detected between groups 20 if a large smoothing kernel (20mm) is used. The power calculation tool they provide suggests that a 0.027mm difference is detectable in a pairwise vertex based analysis with N=76 and 20mm smoothing. Studies examining global mean thickness, i.e. a very large smoothing, like ours, and more than two timepoints can expect a higher sensitivity. The magnitude of thickness differences resulting from differences in template depth is thus large enough to be of potential concern in some studies.

Template depth was not a predictor of trajectory of cortical thickness or supratentorial volumes in any of our tests.

The effect of template depth on mean thickness was also observed in real longitudinal data from the ADNI study. Repeated analysis of baseline scans using templates of different depth, constructed from real (not simulated) follow up scans, demonstrated that template depth was a predictor of mean cortical thickness, and that the effect size was large ( $f=0.42$ ).

The Cohen's  $f$  effect sizes for significant group differences were typically in the range 0.1-0.2 (small) for cortical thickness in the presence of GM atrophy, while moderate effects (over 0.25) were observed in the presence of WM atrophy. Occasional large effect sizes ( $>0.6$ ) were observed for tests of supratentorial volume in the presence of GM atrophy. Effect sizes for prediction of thickness or supratentorial volume from template depth were typically small.

Trajectories are often the property of interest in longitudinal studies, and fortunately our study suggests that the impact of difference in template depth on trajectories of thickness is not as strong or consistent as the impact on mean thickness. However we have demonstrated that detectable differences can be induced. Group differences in mean thickness, while typically not the main aim of longitudinal studies, are a frequently reported and are an important component of group characteristics. We have detected a strong association between template depth and mean thickness in both real data and a simulated white matter atrophy data set.

We can speculate on possible underlying causes of the bias we have detected. The tendency of templates created via iterative registration procedures to be larger than the average size of the contributing images is known as *template drift* and is well documented (Abdollahi et al. 2014; Van Essen et al. 2011; Buckner et al. 2004). The individual templates used by the FreeSurfer longitudinal pipeline are created via an iterative process, and may be susceptible to template drift. The relationship between template drift and the bias we have detected is currently under investigation. Another possibility is that the small number of samples involved in creating a template leads to discretization effects, where adding a single sample may lead to a sudden change in template appearance and thus pipeline behaviour.

The FreeSurfer longitudinal pipeline is very effective at improving segmentation reproducibility (Jovicich et al. 2013) and maximizing the information that can be extracted from any one scan based on information from other acquisitions of the same individual at different study timepoints. Great care has been taken to ensure that the templates used are not biased by one acquisition. The potential for bias that we have discussed is real, and is currently best addressed by careful monitoring of group differences in distribution of template depth, and relationships between template depth and study exposures, in the same way that other, non imaging, study measures must be carefully monitored for potential sources of bias. The option of statistically adjusting for template depth may be valid in some studies, but will need to be considered carefully due to correlation with other study variables, such as age. Another option may be to use a fixed template depth, for example in a study with a mix of 3 and 4 scans per individual an attractive option could be to use the first 3 scans to construct the individual template, gaining most of the improved robustness. However the current FreeSurfer implementation requires that scans must be included in the template in order to use the longitudinal pipeline.

*Limitations.* This is largely a simulation study with longitudinal data derived from a set of real baseline scans. Simulation is the simplest way to provide known patterns and rates of change. However the limitation is that the variation between acquisitions is lower than might be expected with real data, with scanner drift and changes in noise characteristics, orientation, bias field inhomogeneity and non brain tissue. The FreeSurfer longitudinal pipeline improves segmentation reliability in the presence of such inter-acquisition differences. The absence of such effects in our simulations may have strengthened the signal due to template depth compared to what may be expected in a real study. However, we were able to detect the association between mean thickness and template depth in a real study.

## 7 Conclusion

There is an association between template depth and estimates of mean cortical thickness and supratentorial volume produced by the FreeSurfer longitudinal pipeline. This study has detected the association in both real and simulated data sets. We did not detect an association between template depth and cortical thickness trajectories, but were able to detect group differences in trajectory when the groups differed by template depth.

These findings illustrate yet another factor that must be carefully considered when interpreting the results of longitudinal studies, as differences in followup rates can have direct impact on the results of processing pipelines.

## 8 Acknowledgements

This research was conducted within the Developmental Imaging research group, Murdoch Childrens Research Institute and the Peninsula Clinical School, Monash University. It was supported by the Murdoch Childrens Research Institute, the Royal Children’s Hospital, Monash University and the Victorian Government’s Operational Infrastructure Support Program. The project was generously supported by the Royal Children’s Hospital Foundation devoted to raising funds for research at The Royal Children’s Hospital.

## 9 References

- Abdollahi, Rouhollah O., Hauke Kolster, Matthew F. Glasser, Emma C. Robinson, Timothy S. Coalson, Donna Dierker, Mark Jenkinson, David C. Van Essen, and Guy A. Orban. 2014. “Correspondences Between Retinotopic Areas and Myelin Maps in Human Visual Cortex.” *NeuroImage* 99: 509–24. <https://doi.org/https://doi.org/10.1016/j.neuroimage.2014.06.042>.
- Bates, Douglas, Martin Mächler, Ben Bolker, and Steve Walker. 2015. “Fitting Linear Mixed-Effects Models Using lme4.” *Journal of Statistical Software* 67 (1): 1–48. <https://doi.org/10.18637/jss.v067.i01>.
- Buckner, Randy L., Denise Head, Jamie Parker, Anthony F. Fotenos, Daniel Marcus, John C. Morris, and Abraham Z. Snyder. 2004. “A Unified Approach for Morphometric and Functional Data Analysis in Young, Old, and Demented Adults Using Automated Atlas-Based Head Size Normalization: Reliability and Validation Against Manual Measurement of Total Intracranial Volume.” *NeuroImage* 23 (2): 724–38. <https://doi.org/https://doi.org/10.1016/j.neuroimage.2004.06.018>.
- Dale, Anders M, Bruce Fischl, and Martin I Sereno. 1999. “Cortical Surface-Based Analysis: I. Segmentation and Surface Reconstruction.” *Neuroimage* 9 (2): 179–94.
- Fischl, Bruce, and Anders M Dale. 2000. “Measuring the Thickness of the Human Cerebral Cortex from Magnetic Resonance Images.” *Proceedings of the National Academy of Sciences* 97 (20): 11050–5.
- Jovicich, Jorge, Moira Marizzoni, Roser Sala-Llonch, Beatriz Bosch, David Bartrés-Faz, Jennifer Arnold, Jens Benninghoff, et al. 2013. “Brain Morphometry Reproducibility in Multi-Center 3T Mri Studies: A Comparison of Cross-Sectional and Longitudinal Segmentations.” *NeuroImage* 83: 472–84. <https://doi.org/https://doi.org/10.1016/j.neuroimage.2013.05.007>.
- Khanal, Bishesh, Nicholas Ayache, and Xavier Pennec. 2016. “Simulating Realistic Synthetic Longitudinal Brain Mris with Known Volume Changes.” *NeuroImage*, 12.
- Khanal, Bishesh, Marco Lorenzi, Nicholas Ayache, and Xavier Pennec. 2016. “A Biophysical Model of Brain Deformation to Simulate and Analyze Longitudinal Mris of Patients with Alzheimer’s Disease.” *NeuroImage* 134: 35–52.

Liem, Franziskus, Susan Méridat, Ladina Bezzola, Sarah Hirsiger, Michel Philipp, Tara Madhyastha, and Lutz Jäncke. 2015. “Reliability and Statistical Power Analysis of Cortical and Subcortical Freesurfer Metrics in a Large Sample of Healthy Elderly.” *NeuroImage* 108: 95–109. <https://doi.org/https://doi.org/10.1016/j.neuroimage.2014.12.035>.

Pardoe, Heath R., David F. Abbott, Graeme D. Jackson, and The Alzheimer’s Disease Neuroimaging Initiative. 2013. “Sample Size Estimates for Well-Powered Cross-Sectional Cortical Thickness Studies.” *Human Brain Mapping* 34 (11): 3000–3009. <https://doi.org/https://doi.org/10.1002/hbm.22120>.

R Core Team. 2020. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.

Reuter, Martin, Nicholas J Schmansky, H Diana Rosas, and Bruce Fischl. 2012. “Within-Subject Template Estimation for Unbiased Longitudinal Image Analysis.” *Neuroimage* 61 (4): 1402–18.

Van Essen, David C., Matthew F. Glasser, Donna L. Dierker, John Harwell, and Timothy Coalson. 2011. “Parcellations and Hemispheric Asymmetries of Human Cerebral Cortex Analyzed on Surface-Based Atlases.” *Cerebral Cortex* 22 (10): 2241–62. <https://doi.org/10.1093/cercor/bhr291>.

## A Supplementary Material

### A.1 Simulated atrophy

Figure 14 illustrates the surfaces extracted by the cross sectional FreeSurfer pipeline at two time points for a participant with simulated cortical atrophy. This example illustrates time points 1 and 5 and an atrophy level of 04.

### A.2 Statistical model formulae

#### A.2.1 Group comparisons

The following lme4 formulae were used to evaluate group differences in brain measure (mean cortical thickness and supra tentorial volume). TimePoint was a continuous variable and TemplateDepthGroup was a categorical variable with two levels (for example levels corresponding to template depths of 3 and 6).

The first model was used to investigate group differences in mean, while second model, with the interaction term, was used to investigate differences in trajectory.

```
lmer(BrainMeasure ~ TimePoint + TemplateDepthGroup + (1|ID))  
lmer(BrainMeasure ~ TimePoint * TemplateDepthGroup + (1|ID))
```

#### A.2.2 Template depth as a predictor

The formula used to investigate the relationship between template depth and mean thickness in Set A (the real ADNI data) was:

```
lmer(MeanThickness ~ TemplateDepth + (1|ID))
```

The dataset for which this model was fitted contained multiple estimates of baseline thickness, each corresponding to a different template depth, for each participant. TemplateDepth was a continuous variable

The relationship between template depth and brain measure in Set B, the simulated atrophy data, was investigated with the following models. The models were fitted separately to each simulated atrophy level and both TimePoint and TemplateDepth were continuous variables. The model with an interaction term was used to investigate the relationship between template depth and trajectory of brain measure.



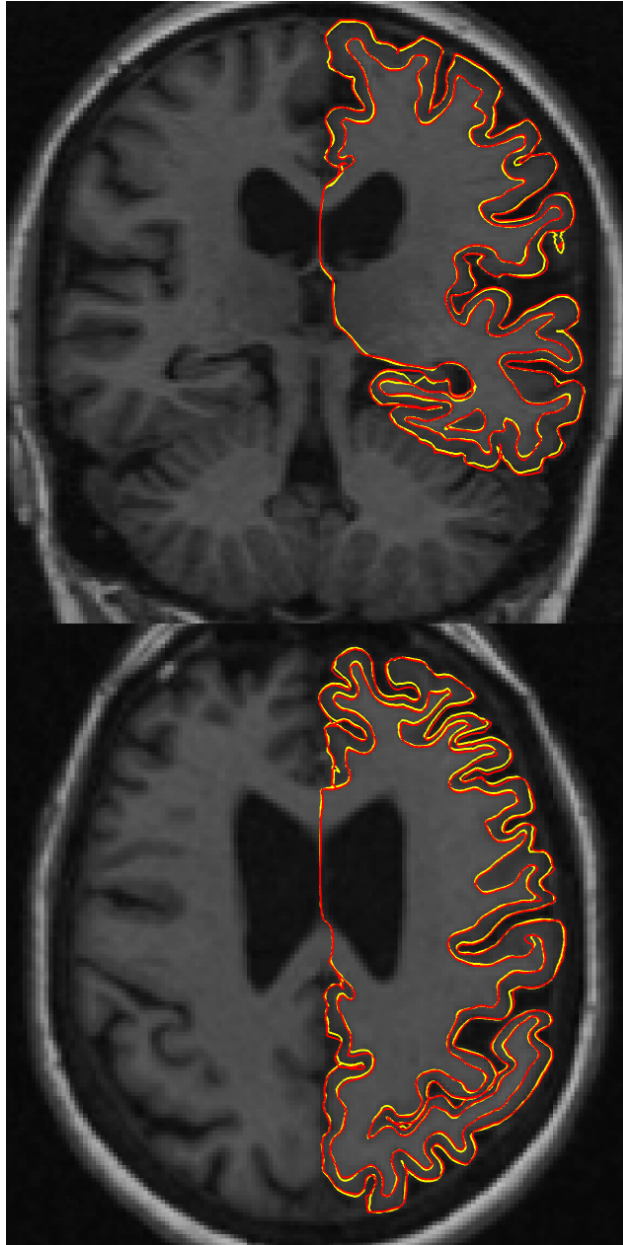


Figure 14: FreeSurfer WM and pial surfaces corresponding to simulated cortical atrophy at two time points. The red and yellow surfaces correspond to the earlier and later time points and the background image is from the earlier timepoint

```
lmer(BrainMeasure ~ TimePoint + TemplateDepth + (1|ID))  
lmer(BrainMeasure ~ TimePoint * TemplateDepth + (1|ID))
```

### A.3 Visualization of supratentorial volume

#### A.3.1 Cross sectional processing, Set A

Trajectories of supratentorial volume for Set A, estimated using the cross sectional FreeSurfer pipeline are shown in Supplementary Figure 15.

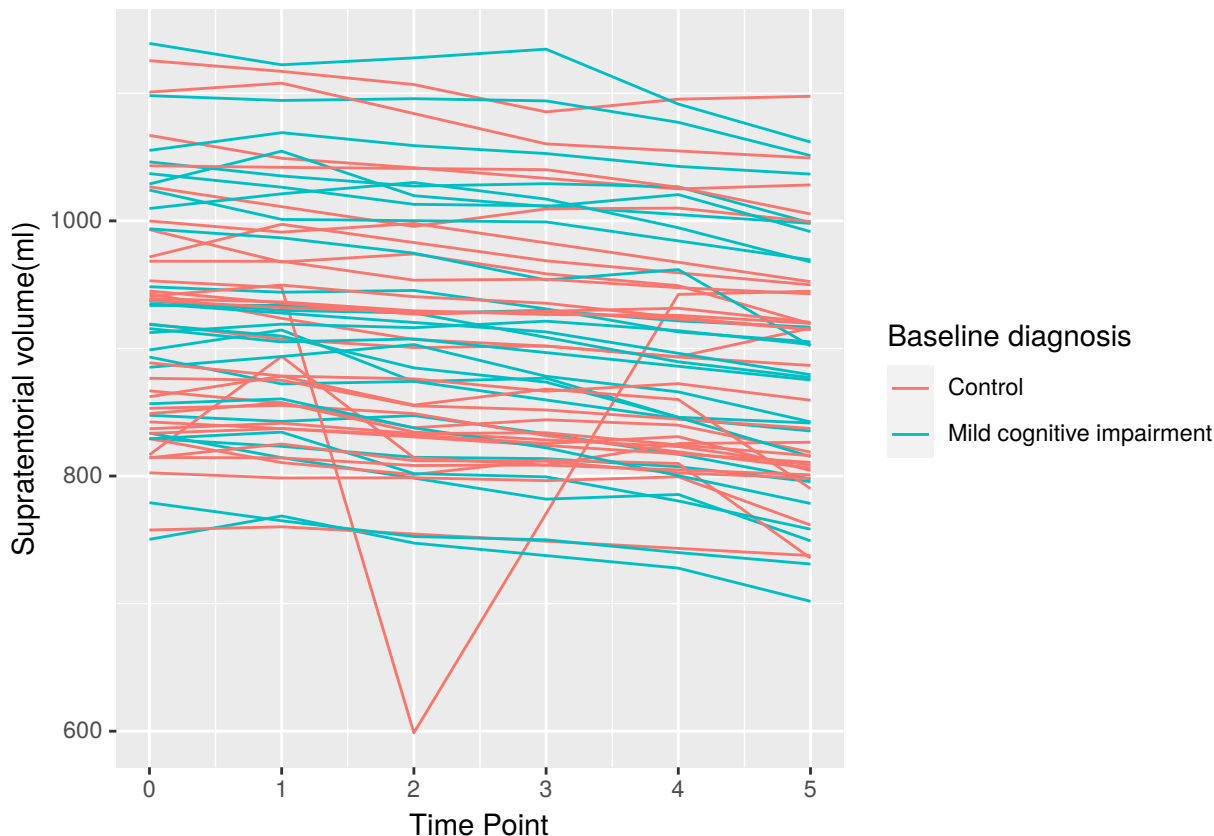


Figure 15: Supratentorial volume estimated with cross sectional FS for ADNI participants with scans at 60 months

#### A.3.2 Cross sectional processing, Set B

Trajectories of supratentorial volume for Set B with simulated cortical atrophy and simulated WM atrophy, estimated using the cross sectional FreeSurfer pipeline are shown in Supplementary Figures 16 and 17.

#### A.3.3 Longitudinal Processing Results, Set B

Trajectories of supratentorial volume for Set B with simulated cortical atrophy and simulated WM atrophy, estimated using the longitudinal FreeSurfer pipeline are shown in Supplementary Figures 18 and 19.



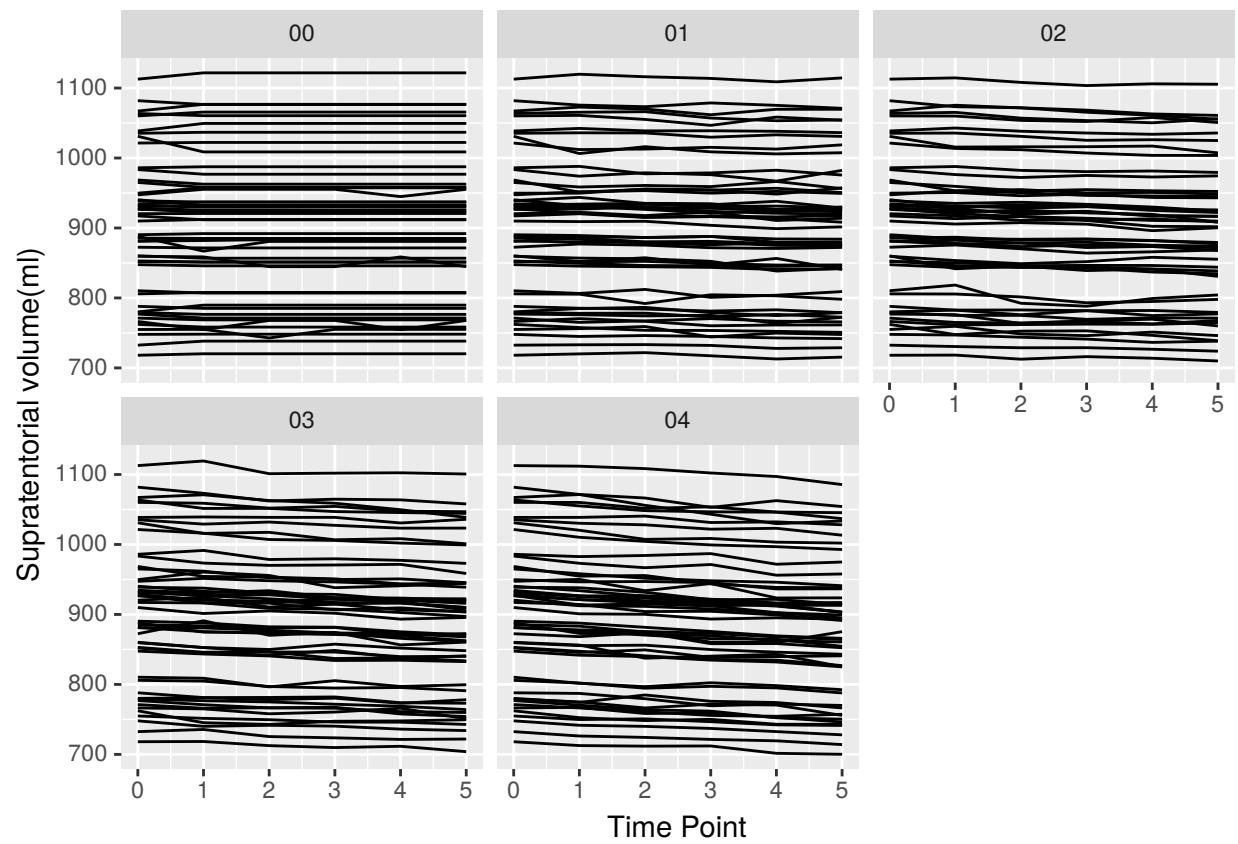


Figure 16: Supratentorial volume estimated with cross sectional FS in presence of simulated cortical atrophy. Each panel corresponds to a different level of simulated atrophy.

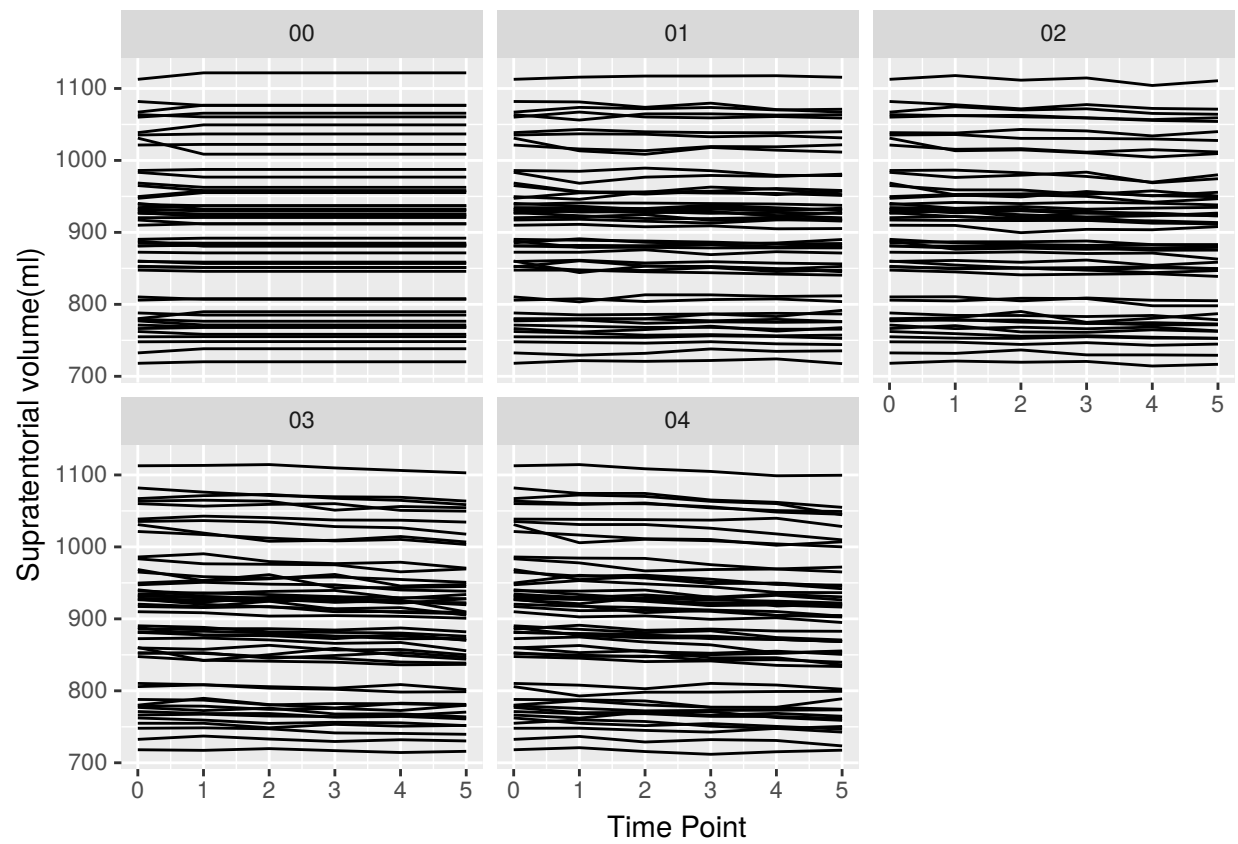


Figure 17: Supratentorial volume estimated with cross sectional FS in presence of simulated WM atrophy. Each panel corresponds to a different level of simulated atrophy.

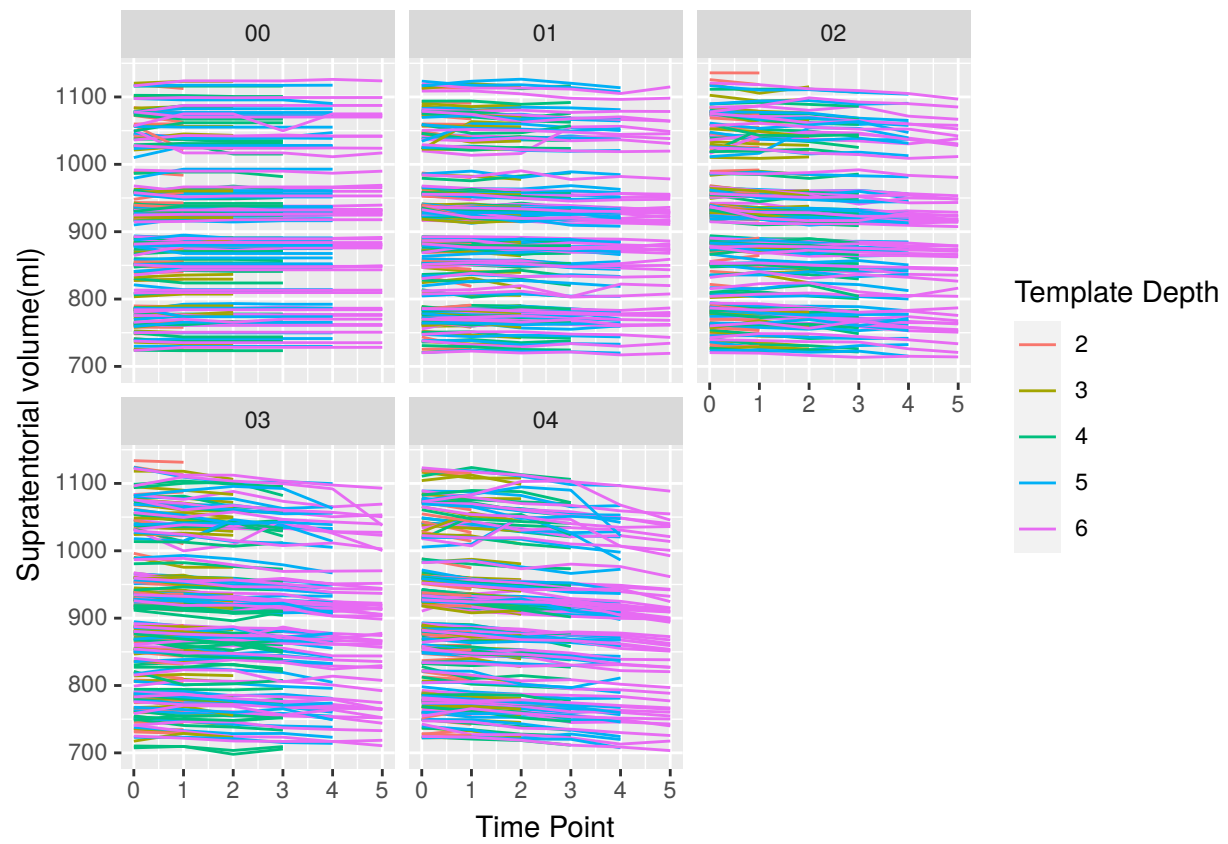


Figure 18: Supratentorial volume estimated with longitudinal FS in presence of simulated cortical atrophy. Each panel corresponds to a different level of simulated atrophy.

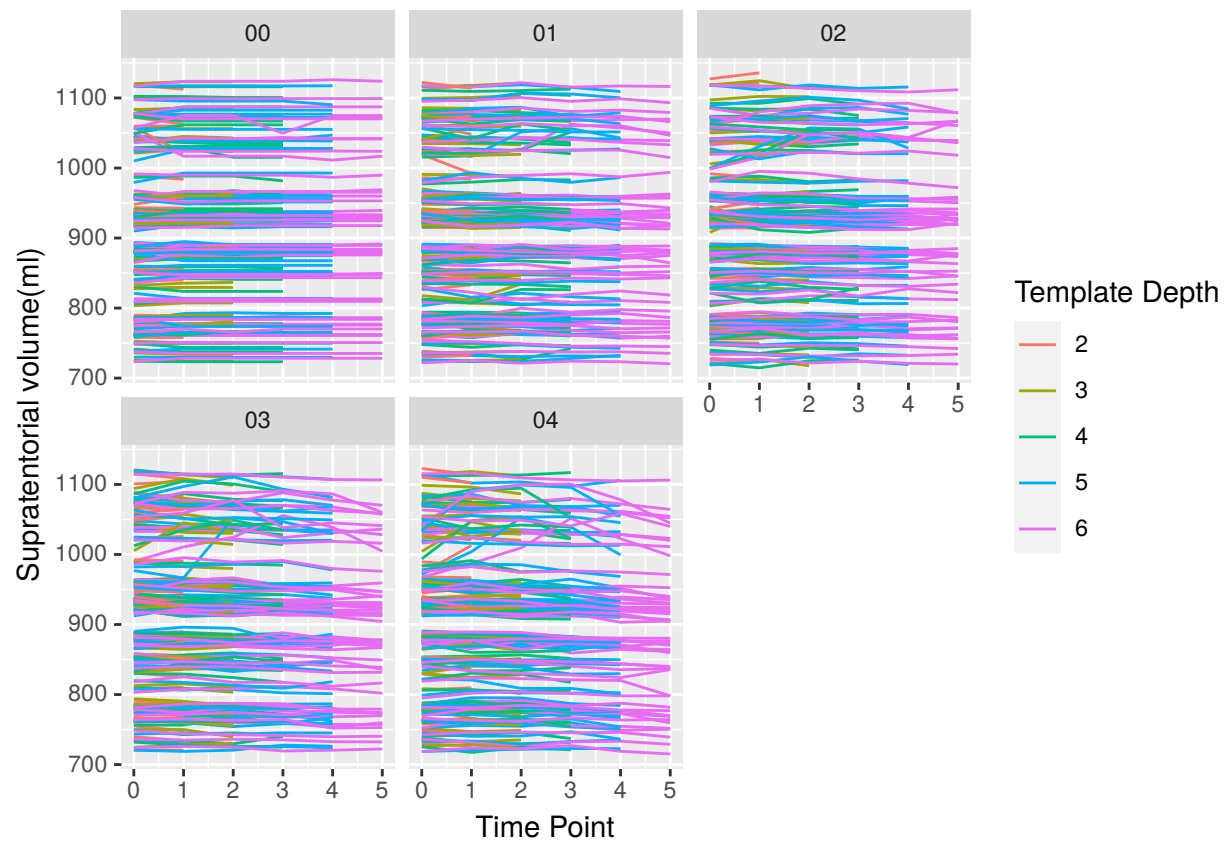


Figure 19: Supratentorial volume estimated with longitudinal FS in presence of simulated WM atrophy. Each panel corresponds to a different level of simulated atrophy.