

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21

Using Outliers in Freesurfer Segmentation Statistics to Identify Cortical Reconstruction Errors in
Structural Scans

Abigail B. Waters, awaters2@su.suffolk.edu, 703-975-3449 ^a

Ryan A. Mace, rmace@su.suffolk.edu, 410-917-8239 ^a

Kayle S. Sawyer, kslays@bu.edu, 617-875-5967 ^{b, c, d, e}

David A. Gansler (Corresponding Author), dgansler@su.suffolk.edu, 617-305-6397 ^a

^a Department of Psychology, Suffolk University, 73 Tremont Street, Boston, MA, USA

^b Department of Anatomy & Neurobiology, Boston University School of Medicine, Boston, MA,
USA

^c VA Boston Healthcare System, Boston, MA, USA

^d Athinoula A. Martinos Center for Biomedical Imaging, Massachusetts General Hospital,
Harvard Medical School, Charlestown, MA, USA,

^e Sawyer Scientific, LLC, Boston, MA, USA

Submission date: 8/21/2017

Counts: abstract (288); text (4192); references (27); tables (3); figures (2); pages (18)

Abstract

1
2 Introduction: Quality assurance (QA) is vital for ensuring the integrity of processed
3 neuroimaging data for use in clinical neurosciences research. Manual QA (visual inspection) of
4 processed brains for cortical surface reconstruction errors is resource-intensive, particularly with
5 large datasets. Several semi-automated QA tools use quantitative detection of subjects for editing
6 based on outlier brain regions. There were two project goals: (1) evaluate the adequacy of a
7 statistical QA method relative to visual inspection, and (2) examine whether error identification
8 and correction significantly impacts estimation of cortical parameters and established brain-
9 behavior relationships.

10 Methods: T1 MPRAGE images ($N = 530$) of healthy adults were obtained from the NKI-
11 Rockland Sample and reconstructed using Freesurfer 5.3. Visual inspection of T1 images was
12 conducted for: (1) participants ($n = 110$) with outlier values (z scores $\pm 3 SD$) for subcortical and
13 cortical segmentation volumes (outlier group), and (2) a random sample of remaining
14 participants ($n = 110$) with segmentation values that did not meet the outlier criterion (non-
15 outlier group).

16 Results: The outlier group had 21% more participants with visual inspection-identified errors
17 than participants in the non-outlier group, with a medium effect size ($\Phi = 0.22$). Nevertheless, a
18 considerable portion of images with errors of cortical extension were found in the non-outlier
19 group (41%). Sex significantly predicted error rate; men were 2.8 times more likely to have
20 errors than women. Although nine brain regions significantly changed size from pre- to post-
21 editing (with effect sizes ranging from 0.26 to 0.59), editing did not substantially change the
22 correlations of neurocognitive tasks and brain volumes ($ps > 0.05$).

1 Conclusions: Statistically-based QA, although less resource intensive, is not accurate enough to
2 supplant visual inspection. We discuss practical implications of our findings to guide resource
3 allocation decisions for image processing.

4 Keywords: quality assurance, automated segmentation statistics, reconstruction error, Freesurfer

5

1. Introduction

1 Accurate brain volume estimation is considered essential to research brain-behavior
2 relationships. Structural neuroimaging studies have found significant associations between
3 regional brain volumes and domains of cognitive functioning, including executive functioning
4 (Yuan & Raz, 2014), attention (Seidman, Valera, & Makris, 2005), and memory (Van Pettan,
5 2004). As neuroimaging increasingly emphasizes the use of large-scale datasets to assess these
6 relationships, ensuring the integrity and validity of data via quality assurance (QA; also called
7 quality control), is critical and must be done efficiently. There is an emergent need to understand
8 the effectiveness of QA methodology to help researchers make resource allocation decisions for
9 image processing.

11 Freesurfer (<http://surfer.nmr.mgh.harvard.edu/fswiki>) is a commonly used open source
12 software suite for automated processing of magnetic resonance imaging (MRI) data. Although
13 validation studies have shown that automated segmentation in Freesurfer is commensurate to
14 manual measurement (Fischl et al., 2002), post-reconstruction visual inspection of cortical “pial”
15 surface segmentation has become common practice to identify incorrect inclusion of non-brain
16 tissues (Desikan et al., 2010). Manual QA (i.e., visual inspection of imaging slices) used to
17 confirm the validity of reconstructed images is time- and resource-intensive, particularly for
18 large-scale neuroimaging datasets. QA methods vary greatly between studies (e.g. Chen et al.,
19 2015; Kaufmann et al., 2017, Ahmed et al., 2015); some studies use manual inspection, semi-
20 automated QA, or a combination of both methods. This variation may differentially influence
21 neuroimaging data used in research analyses.

22 Semi-automated QA methods for reconstructed images are needed to address limited
23 resources in neuroimaging research. Both the QA toolkit in the Freesurfer software suite (Koh,
24 Lee, Pacheco, Pappu, & Vinke, 2017) and the Mindcontrol web application (Keshavan et al.,

1 2017) use statistical analyses of cortical and subcortical regions to identify images that may need
2 manual correction. These regional volumetric measurements are automatically generated based
3 on the probabilistic location of structures (Fischl et al., 2002; Fischl et al., 2004), which due to
4 individual variation of anatomical features, may result in errors of cortical inclusion and
5 exclusion. Therefore, regions with surface reconstruction errors would have over- or under-
6 estimation of volumetric measurement resulting in statistical outliers.

7 The primary objective of this study was to investigate whether anomalies in cortical and
8 subcortical volumetric measurements are associated with reconstruction errors, thereby testing
9 the assumption that statistically-based methods can identify images with reconstruction errors. In
10 addition, we investigated whether participant characteristics (age and sex) impacted the odds of
11 reconstruction errors. This study focused on errors where the cortical surface extended into non-
12 brain tissues. Manual correction of the white matter surface using control points does not result
13 in significantly different volumetric measurements (McCarthy et al., 2015); therefore, errors in
14 delineation of white matter are unlikely to result in statistical outliers. We hypothesized that the
15 outlier images would have significantly more cortical surface errors than those not identified via
16 this method. We did not have a priori hypotheses regarding the influence of participant
17 characteristics.

18 A secondary objective was to determine whether the identification and correction of
19 cortical boundary errors significantly impact established brain-behavior relationships. We
20 investigated whether the association between neurocognitive measures and brain volumes would
21 significantly differ from pre- to post-editing. We expected that some volumes (i.e., precentral,
22 postcentral gyrus) would be more affected by errors than others, given previous literature
23 (Keshavan et al., 2017). We had no a priori hypothesis regarding the impact of editing on brain-

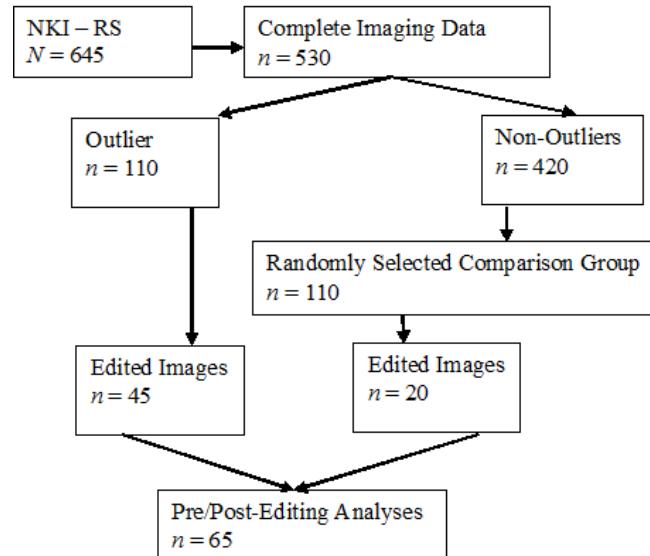
1 behavior relationships. We discuss the implications of our findings in the context of guiding
2 resource allocation decisions regarding neuroimaging QA and the potential influence on the
3 analysis of brain-behavior relationships.

4 **2. Methods**

5 **2.1 Participants**

6 De-identified phenotypic and neuroimaging data for 645 participants was made available
7 via the enhanced Nathan Kline Institute – Rockland Sample (NKI-RS), an open-access, cross
8 sectional, community sample (Nooner, et al., 2012). Rockland County’s economic and ethnic
9 demographics are representative of the United States census (U.S. Census Bureau, 2009), making
10 the NKI-RS generalizable to the U.S. population. Zip-code based recruitment was conducted via
11 mailings, flyers, and electronic advertisement. Data use agreement was accepted by NKI-RS and
12 data handling procedures were approved by the Institutional Review Board at Suffolk University.

13 The NKI-RS excludes participants if they were diagnosed with or reported severe
14 psychiatric disorders (bipolar disorder, schizophrenia disorder, schizoaffective disorder), severe
15 developmental disorders (autism spectrum disorders, intellectual disabilities), current suicidal or
16 homicidal ideation, severe cerebral trauma (stroke, moderate to severe traumatic brain injury,
17 transient ischemic attack in the past two years), severe neurodegenerative disorders (Parkinson’s
18 disease, Huntington’s disease, dementia), a history of substance dependence in the past two years
19 (except cannabis), a lifetime history of psychiatric hospitalization, current pregnancy, or MRI
20 contraindications. Participants were excluded from these analyses if T1-weighted structural
21 images were incomplete or had significant image artifacts. A total of 530 adult participants from
22 the NKI-RS who met eligibility criteria and had complete scan data were included in this study
23 (Figure 1 describes participant flow for study sample).



1

2 *Figure 1.* Flowchart of Nathan Kline Institute – Rockland Sample (NKI-RS) participants and the
3 study sample.

4 **2.2 Imaging.**

5 Structural images were acquired using a 3T Siemens Trio scanner (T1 MPRAGE, voxel
6 size = 1.0 x 1.0 x 1.0 mm, 176 slices, echo time = 2.52 ms, repetition time = 1900 ms, field of
7 view = 250 mm). MPRAGE data obtained from the NKI-RS dataset are available in their raw
8 form. DICOM data were converted to the mgz format and MPRAGE images were auto-
9 reconstructed in Freesurfer 5.3. Each image was mapped into standard morphological space
10 (MNI305; Collins, 1994), and volumes were generated for cortical and subcortical regions based
11 on the Desikan-Killiany Atlas (Desikan et al., 2006), which included white matter, gray matter,
12 and other anatomical features.

13 **2.3 Outlier identification and comparison group.**

14 Two participant subsamples were identified from the 530 participants that met study
15 inclusion criteria. First, Freesurfer-automated segmentation statistics for subcortical and cortical
16 segmentations were standardized. Next, all participants with z-scores 3 *SD* above or below the
17 mean for one or more normalized brain volume labels estimated by Freesurfer (66 total

1 automated segmentation brain region variables) were identified as the *outlier group* ($n = 110$).
 2 This statistically-based identification method was designed to assess the underlying assumption
 3 that segmentation statistics can be used to identify incorrectly reconstructed images, which is the
 4 basis of many semi-automated QA techniques (Kaufmann et al., 2017, Koh, Lee, Pacheco,
 5 Pappu, & Vinke, 2017).

6 The 3 *SD* cut-off was chosen because 2 *SD* cut-off was too broad for practical reasons (it
 7 identified 311 of 530 images as outliers) and would most likely result in low specificity. If
 8 outlier measurements predict error rates, we hypothesized that the more extreme outliers would
 9 have greater specificity. Finally, a random sample of participants that did not meet the outlier
 10 criterion (i.e., no standardized segmentation brain volumes ± 3 *SD*) were selected via random
 11 number generation for the *non-outlier group* ($n = 110$), which served as a comparison group.
 12 Table 1 presents the demographic characteristics for both groups, which represent the final
 13 sample for data analysis ($N = 220$).

14

15 Table 1

16 Basic Demographics by Group

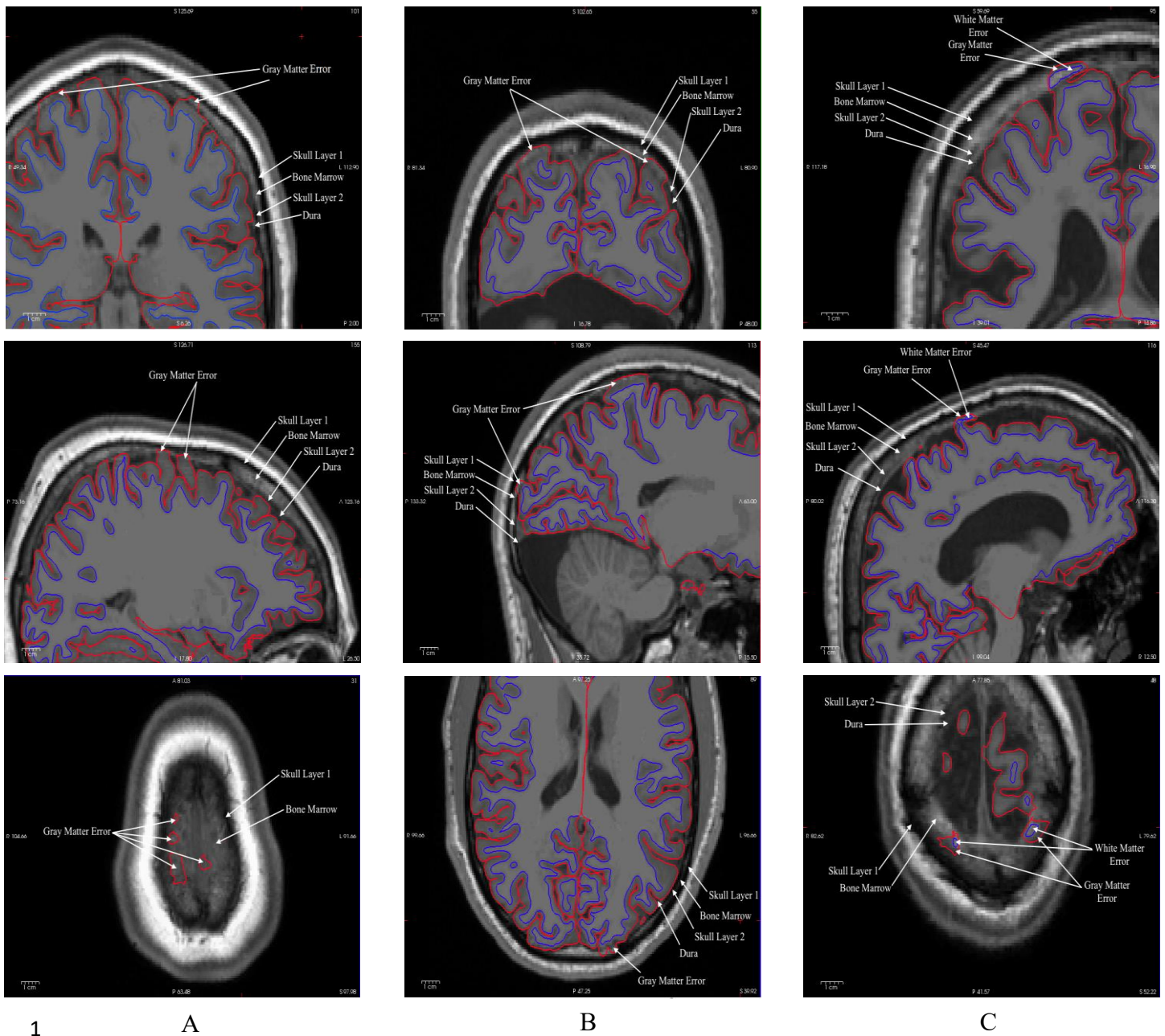
	Outlier group ($n = 110$)		Non-Outlier group ($n = 110$)		Total Sample ($N = 220$)	
	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>
Age	53.57	18.92	43.90	17.19	48.74	18.67
	%		%		%	
Sex						
Male	50		29.1		39.5	
Female	50		70.9		60.5	
Race						
Native Am.	0.9		0		0.5	
Asian	3.6		6.4		5.0	
Black	14.5		18.3		16.5	
White	79.1		73.4		76.6	
Other	1.8		1.8		1.4	

17

1 **2.4 Error identification.**

2 Because almost all of the identified outlier regions were ≥ 3 *SD* above the mean (only
3 8.2% of outliers were found at least 3 *SD* below the mean), investigators hypothesized that the
4 statistical method was best suited to identify erroneous inclusion of skull and dura into cortical
5 volumes. Therefore, error identification focused on errors of cortical boundary extension (See
6 Figure 2 for several example images). All brain images were visually inspected using Freeview
7 to identify the presence or absence of errors in reconstruction. Graduate students (A.W. and
8 R.M.), who were trained in visual inspection by a senior author (K.S.) with extensive experience
9 in imaging science, served as raters.

10 Both graduate students performed manual visual inspection on each of the auto-
11 reconstructed brains for both outlier and non-outlier participants ($N = 220$) in Freeview. White
12 and gray matter surfaces were overlaid onto T1 scans (i.e., before skull stripping) to assist in
13 cortical identification. All slices were viewed in the coronal plane. Each image was coded as
14 having no errors (0) or projection errors (1). A projection error occurred when the cortical
15 boundary incorrectly extended outward more than approximately 3 voxels into the dura and/or
16 skull on at least one slice. The trained graduate students completed their ratings independently
17 and were blind to participant characteristics (including outlier or non-outlier group membership).
18 Both raters recorded their viewing time using a stopwatch.



1

A

B

C

2 *Figure 2.* Coronal, Sagittal, and Horizontal View of (A) Errors of cortical extension in the pre-
3 central/post-central region; (B) Errors of cortical extension in posterior regions; and (C) Errors of
4 cortical extension that include both white and grey matter.

5 On average, raters took 70 seconds to visually inspect and identify errors for each image.

6 This estimate does not include the time required to load scan data in Freeview (approximately 20

1 seconds), which depended on computer speed. Interrater reliability analyses indicated moderate
2 agreement (Cohen's Kappa = 0.65, 95% CI: 0.55 – 0.74; 82.4% agreement), with 181 of the 210
3 ratings consistent between both raters. For the 39 discordant ratings, consensus was reached with
4 both raters and a senior author (KS). Most discordant images (76.9%) were assigned consensus
5 ratings of no error. Twenty-four of these 39 images were assigned a rating of no error because
6 they did not meet full criteria for an error (> 3 voxels cortical extension), but contained some
7 partial or full voxel cortical extension (some for large areas of the cortical surface; see
8 Supplemental Figure 1). Due to variability in scan quality and anatomical features, raters
9 concluded that these images could not be conclusively categorized as errors of cortical extension.
10 These images were instead categorized as no error images. These ratings were combined with the
11 agreement ratings (181 images) and used in all subsequent analyses. In total, 138 of the 210
12 images had errors (65.7%).

13 **2.5 Group and error characteristics.**

14 Between group t-tests, chi-square tests, and logistic regression were used to determine if
15 group membership (outlier versus non-outlier) was significantly associated with reconstruction
16 error or other participant characteristics. Further analyses were performed to assess
17 characteristics associated with errors, regardless of group membership.

18 **2.6 Brain volumes and measures of neurocognitive functioning.**

19 Pearson correlations were used to test the associations between brain regions and
20 neurocognitive tasks. Correlation coefficients were compared between images pre- and post-
21 editing using a Fisher r-to-z transformation (Meng, Rosenthal, & Rubin, 1992) and by comparing
22 effect sizes (Cohen, 1988).

1 A subsample of images with errors (65 of 138 images) were manually edited (Savalia,
2 Agres, & Wig, 2015), and cortical regions were averaged across hemispheres for the purpose of
3 data reduction. Volumes were compared before and after editing (Table 2). Nine brain volumes
4 across all edited images showed significant decreases in volume from pre- to post-editing.

5 These nine volumes were used to assess brain-behavior relationships. Pearson correlation
6 coefficients were also used to estimate the associations between edited values and measures of
7 neurocognitive functioning to determine if error correction meaningfully influenced the outcome
8 of brain-behavior relationships. Tasks were chosen based on previous literature to assess the
9 most robust brain-behavior relationships (Yuan & Raz, 2014; Phan, Wager, Taylor, & Liberzon,
10 2002; Ward & Frackowiak, 2003).

11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

[Table 2 found on following page]

1 Table 2

2 Mean Decreases in Volumes (mm³) Pre- to Post-Editing (N = 65)

	<i>M</i>	<i>SD</i>	95% CI		<i>t</i>	<i>p</i>	<i>d</i>
			Lower	Upper			
Precentral	205.33	345.79	119.65	291.01	4.79	<0.001	0.59
Postcentral	221.45	461.53	107.09	335.82	3.87	<0.001	0.48
Superior Parietal	117.82	281.35	48.10	187.53	3.38	0.001	0.42
Posterior Cingulate	12.33	32.85	4.19	20.47	3.03	0.004	0.38
Inferior Parietal	132.58	386.33	36.85	228.30	2.77	0.007	0.34
Isthmus Cingulate	10.80	37.32	1.55	20.05	2.33	0.023	0.29
Precuneus	12.15	46.96	0.52	23.79	2.09	0.041	0.26
Superior Frontal	57.04	220.52	2.40	111.68	2.09	0.041	0.26
Caudal Middle Frontal	35.18	137.32	1.16	69.21	2.07	0.043	0.26
Fusiform	32.55	147.25	-3.93	69.04	1.78	0.079	0.22
Rostral Middle Frontal	50.85	240.63	-8.78	110.47	1.70	0.093	0.21
Parsorbitalis	5.97	28.50	-1.09	13.03	1.69	0.096	0.21
Rostral Anterior Cingulate	10.11	48.84	-1.99	22.21	1.67	0.100	0.21
Caudal Anterior Cingulate	6.65	32.81	-1.48	14.78	1.63	0.107	0.20
Inferior Temporal	49.22	249.59	-12.62	111.07	1.59	0.117	0.20
Lingual	28.65	146.77	-65.02	7.72	-1.57	0.121	0.20
Supramarginal	96.60	496.99	-26.55	219.75	1.57	0.122	0.19
Medial Orbitofrontal	15.61	82.56	-4.85	36.07	1.52	0.132	0.19
Parstriangularis	13.60	73.57	-4.63	31.83	1.49	0.141	0.18
Paracentral	6.96	44.67	-4.11	18.03	1.26	0.214	0.16
Middle Temporal	39.75	276.40	-28.73	108.24	1.16	0.251	0.14
Parahippocampal	6.92	53.18	-6.26	20.09	1.05	0.298	0.13
Temporal Pole	6.86	54.70	-6.69	20.42	1.01	0.316	0.13
Insula	8.93	89.67	-13.29	31.15	0.80	0.425	0.10
Pericalcarine	2.97	31.66	-4.88	10.81	0.76	0.452	0.09
Bankssts	3.61	47.20	-8.09	15.30	0.62	0.540	0.08
Frontal Pole	2.09	27.83	-8.99	4.80	-0.61	0.547	0.08
Parsopercularis	4.12	69.96	-13.22	21.45	0.47	0.637	0.06
Lateral Orbitofrontal	4.92	89.43	-27.08	17.24	-0.44	0.659	0.06
Superior Temporal	8.37	188.52	-38.34	55.08	0.36	0.722	0.04
Transversetemporal	0.33	12.04	-2.65	3.31	0.22	0.825	0.03
Lateral Occipital	3.73	153.26	-41.71	34.24	-0.20	0.845	0.02
Cuneus	0.89	46.45	-12.40	10.62	-0.16	0.877	0.02
Entorhinal	0.97	72.06	-16.89	18.83	0.11	0.914	0.01

3 *Note. Volumes in descending order by Cohen's d effect size.*

1 Participants included in these analyses had completed the Delis-Kaplan Executive
2 Functioning System (D-KEFS), the Wechsler Abbreviated Scale of Intelligence (WASI-II), and
3 the Penn Computerized Neuropsychological Test Battery (Penn CNB) during their involvement
4 in NKI-RS. All tests were administered per the guidelines presented in their respective
5 administration manuals by trained research assistants. All tests on the Penn CNB were part of a
6 larger computerized battery that participants completed on a desktop with trained research
7 assistant supervision.

8 **2.6.1 D-KEFS Trails.** The Trails subtest on the D-KEFS includes five conditions.
9 Completion time in seconds from the first, fourth, and fifth conditions were included in analyses.
10 These conditions assess visual attention, mental flexibility, and visuo-motor coordination,
11 respectively (Kaplan, Delis, & Kramer, 2001).

12 **2.6.2 WASI Perceptual Reasoning Index (PRI).** The WASI-II PRI was calculated using
13 the standard scores from two subtests: Block Design and Matrix Reasoning. The Block Design
14 subtest assesses spatial visualization, abstract conceptualization, and visuo-motor coordination.
15 The Matrix Reasoning subtest assesses non-verbal fluid reasoning (Wechsler, 2011).

16 **2.6.3 Penn CNB Conditional Exclusion Task (PCET).** The PCET was developed to
17 assess mental flexibility and abstraction (Gur et al., 2001).

18 **2.6.4 Penn CNB Emotion Recognition Task (ERT).** The ERT is a measure of emotion
19 recognition along a continuum of expression intensity (Gur et al., 2001).

20 **3. Results**

21 **3.1 Error characteristics and statistical method**

22 The rate of reconstruction errors was 62.7% in the outlier group and 40.9% in the non-
23 outlier group (Table 3). There were significant differences in error rates between the non-outlier

1 and outlier groups, ($X^2(1, 220) = 10.49, p = 0.001$), with a medium effect size ($\Phi = 0.22$). Based
 2 on the odds ratio, the outlier images were 1.69 times more likely to have a projection error than
 3 the randomly selected images. Raters identified 45 out of 220 images (20.5%) in need of further
 4 reconstruction that were not detected by the statistical method (i.e., false negative). The
 5 statistical method had an accuracy of 60.9%, a sensitivity of 60.5%, and a specificity of 59.1%.

6

7 Table 3

8 Crosstabulation of Group Status and Error Rate

	Outlier Group (Predicted Positive)	Non-Outlier Group (Predicted Negative)	Total
Error (Positive)	69 (62.7%; 60.5%)	45 (40.9%; 39.5%)	114 (51.8%)
No Error (Negative)	41 (37.3%; 38.7%)	65 (29.5%; 61.3%)	106 (48.2%)
Total	110 (50%)	110 (50%)	220

9 *Note.* * $p \leq .001$. Numbers in parentheses indicate: (column %; row %). There were significant
 10 differences in error rates between the non-outlier and outlier groups, ($X^2(1, 220) = 10.49, p =$
 11 0.001), with a medium effect size ($\Phi = 0.22$).

12

13 Outliers were found in 55 of 66 segmentation volumes. Outliers in the ventricles and
 14 surface hole regions were most frequent, as 56.4% of images in the outlier group were identified
 15 using statistical outliers in the ventricular, surface hole, vessel, and CSF volumes alone (13 of 66
 16 brain region). To verify that $\pm 3 SD$ criterion used in outlier identification was not overly
 17 stringent, chi-square tests were performed on the 220 error-rated images using the less
 18 conservative 2.5 SD and 2 SD cut-offs. Both the 2.5 SD cut-off ($X^2(1, 220) = 10.49, p = 0.001$)
 19 and the 2 SD cut-off ($X^2(1, 220) = 10.49, p = 0.001$) showed significant differences in error rates

1 between outlier images and images that did not meet cut-off criteria. However, the proportion of
2 images with errors that did not meet cut-off criterion remained relatively consistent (41.4% and
3 40% respectively) as stringency was reduced. The proportion of images with errors for the
4 outlier group decreased (60.3% and 56.3% respectively), even as the group size expanded to 121
5 for 2.5 *SD* and 160 for 2 *SD*. Reducing the stringency did not improve sensitivity or specificity
6 of the statistical method.

7 Age ($t(218) = -3.97, p < .001$) and sex ($X^2(1, 220) = 10.06, p = 0.002$) were significantly
8 different between the two groups, with the outlier group 9.7 years older on average and 21%
9 more male. A series of logistic regression analyses were used to examine whether participant
10 characteristics were associated with error rate across groups. Age and sex were entered into the
11 model as predictor variables for projection errors (criterion). The model was statistically
12 significant and explained an estimated 6% (Cox and Snell R^2) to 9% (Nagelkerke R^2) of the
13 variance in projection errors ($X^2 = 14.91, p = 0.001$). Sex was the only significant predictor of
14 error rate (67.8% of men had errors, compared to 41.6% of women), with women being 66% less
15 likely to have projection errors than men ($OR = .34, 95\% CI = .19 - .59$).

16 **3.2 Brain volumes and measures of neurocognitive functioning**

17 A Kruskal-Wallis one-way analysis of variance did not yield significant gender or group
18 effects on error size, as measured by changes in volume post-editing (i.e., the ranks were
19 unchanged by editing). Therefore, all edited images were analyzed together across groups to
20 assess brain-behavior relationships. Of the 34 brain regions, nine (26.5%) showed significant
21 decreases in volume from pre- to post-editing: the inferior parietal, superior parietal, posterior
22 cingulate, isthmus cingulate, precuneus, precentral, postcentral, superior frontal and caudal
23 middle frontal volumes (See Table 2). Cohen's d ranged from 0.26 to 0.59, with all but one

1 meeting criteria for a small effect size ($0.20 < d < 0.50$). None of the associations between
2 selected ROIs and the neurocognitive tasks significantly changed from pre- to post-editing, as
3 estimated by the Fisher z transformation. Pre-editing Pearson's r -values ranged from -0.35 to
4 0.24, while post-editing r -values ranged from -0.35 to 0.23; all correlations meeting Cohen's d
5 criteria for a small to no effect size. The change in r from pre- to post-editing ranged from -.02 to
6 .08. Of 54 correlations, one resulted in changes of effect size categorization: PCET and the
7 isthmus cingulate volume met Cohen's d criteria for a small effect size before editing ($r = .20$)
8 and did not meet criteria after editing ($r = .19$). However, the pre- and post-editing correlations
9 between PCET and the isthmus cingulate volume did not approach significance.

10 **4. Discussion**

11 The principal findings were that: (1) cortical reconstruction error rates were higher in a
12 group identified by a statistical outlier QA method; (2) reconstruction error rates were too
13 prevalent in a randomly identified non-outlier group to conclude that identification by the
14 statistical outlier method alone was effective; and (3) while post-editing volume estimations were
15 significantly lower in a number of instances, these putatively more accurate volume estimations
16 did not meaningfully impact the outcome of a brain-behavior analysis. To our knowledge, this is
17 the first study investigating the effects of QA error correction on associations between brain and
18 neurocognitive measures. Contrary to our hypothesis that errors would only affect surface
19 structures, we found significant decreases in both lateral and midline structures from pre- to post-
20 editing. Regarding participant characteristics, we found that while both age and sex were
21 associated with statistical outlier status. Only sex was significantly associated with error
22 occurrence.

1 Although the use of T2-weighted scans in conjunction with T1 MPRAGE improve
2 segmentation accuracy (Lichy et al, 2005; available with FreeSurfer v5.3 and above), there is
3 still a need for consistent and efficient QA methodology for reconstruction of T1 scans alone.
4 Previous research on statistically-based QA of automated reconstruction has noted the necessity
5 of establishing a link between summary statistics and cortical extension errors. Similarities in
6 signal intensity between cortex and dura on T1 MPRAGE scans make distinguishing the two
7 difficult, even when using visual inspection (Viviani et al., 2017). This difficulty was apparent
8 during this study, as trained graduate student raters with comparable experience had only
9 “adequate” interrater reliability. Given the variability in QA methodology, statistically based QA
10 presents the potential for a more standardized approach (Keshavan et al., 2017).

11 However, our use of statistical outliers in segmentation statistics to identify images with
12 reconstruction errors was insufficient to identify all images with errors. Given the frequency of
13 errors in the non-outlier group, one could expect approximately 40% of errors to go unidentified.
14 These missed errors would be biased toward younger, female brains. Contrary to our hypothesis,
15 statistically-based identification did not reliably identify inflated values resulting from cortical
16 extension errors. Instead, images in the outlier groups are most often identified by inflated
17 ventricle size or other measures of atrophy. These measures of atrophy are most likely age-
18 related, given the significantly higher age in the outlier group. Perhaps statistically based
19 methods are not identifying errors themselves, but rather anatomical features associated with
20 poor registration and reconstruction in the Freesurfer automated pipeline. Poor registration may
21 also account for the higher proportion of errors among male brains, given that men generally
22 have larger ventricular volumes (Lenroot et al., 2007). It is possible that regression based QA
23 diagnostics tailored by subject variables (e.g. age, gender) could increase the error detection rate.

1 It is unclear from our analysis if male brains have a higher incidence of errors because of poor
2 registration. Alternatively, the statistical method may be confounded by unknown additional
3 factors that account for both error rates and participant characteristics.

4 Additionally, we observed a disjunction between the criterion stringency and the
5 detection of errors. Had the outliers been related to the errors themselves, the proportion of errors
6 would increase as stringency decreased. The degree to which atypical anatomical features
7 contribute to poor registration, and by extension errors in reconstruction, seems to be restricted to
8 the most extreme outliers. Refining the statistical criterion did not increase the specificity or
9 sensitivity of identifying statistical outliers with existing methods. Statistically based methods
10 were biased toward brains with age-related atrophy, as indicated by the outlier frequency among
11 measures of atrophy and the age differences between groups. Studies that use this method alone
12 to identify images for surface editing may introduce confounding variables into their data. This
13 may be especially salient for research in aging or degenerative diseases, but further research is
14 warranted.

15 Our findings regarding the precentral and postcentral gyri are consistent with the patterns
16 of cortical misclassification found during the development of Mindcontrol (Keshavan et al.,
17 2017), where these regions were identified as being more prone to errors of cortical extension.
18 We identified additional cortical areas where errors are more likely to affect the volumetric
19 estimation. Editing resulted in significant decreases of volume in broader aspects of the parietal,
20 frontal, and cingulate regions. Although the focus of this study was on errors of cortical
21 extension (i.e., errors in the boundary between gray matter, dura, and skull), mid-line structures
22 (i.e., posterior and isthmus cingulate) significantly changed from pre- to post-editing. We found
23 several errors of cortical extension during visual inspection that included both gray and white

1 matter where bone marrow was included in white matter estimates (Figure 1). Because of the
2 intensity of bone marrow, Freesurfer was more likely to characterize these as white matter and
3 adjust the gray matter boundary accordingly.

4 Despite the limitations of the statistically based QA, manual inspection may be too
5 resource intensive for large datasets. For the NKI Rockland Sample of 530 subjects, it would
6 have taken approximately 10.5 hours to identify all images that require editing to correct pial
7 surfaces. Accounting for the time it took for raters to load the image (approximately 20 seconds),
8 it would have taken 13.25 hours. This does not include the considerable time (approximately 1 to
9 2 hours) it would have taken to edit these images, which we would estimate to be approximately
10 50% of the sample given error rates of scan reconstructions in this study, for a total time of
11 approximately 411 hours for 530 scans.

12 Directing QA efforts toward brain regions most commonly affected by reconstruction
13 error may lessen the resource commitment to manual inspection. Although errors may be present
14 in other brain regions, they are unlikely to result in significant changes pre- to post-editing. This
15 trade-off between time and likelihood of error occurrence can be considered when prioritizing
16 error correction efforts. While there were statistically significant decreases in some brain regions
17 post-editing (26.5%), editing resulted in practically insubstantial decreases in volume (0.1% to
18 2.3%). This is consistent with previous research on the effects of control points to manually edit
19 segmentation between white and gray matter (McCarthy et al., 2015). This suggests that manual
20 intervention does not produce incremental utility for cortical surface boundaries either.

21 Editing images did not significantly impact the relationship between brain volumes and
22 neurocognitive measures. Based on our findings, we expect that allowing these errors to exist as
23 noise in the dataset may decrease statistical power but not confound results. Although error

1 occurrence is more likely in male brains, errors do not significantly impact associations with
2 neurocognitive variables and usually do not significantly change volumetric estimations.
3 Depending on the goals of individual studies and availability of resources for QA, researchers
4 may find that the costs of visual inspection outweigh potential benefits of manual intervention in
5 large datasets. Further research is needed to determine whether these results are replicable for
6 other brain-behavior relationships, clinical samples, or for other imaging techniques (i.e., fMRI,
7 PET).

8 There were several limitations in this study. Interrater reliability for error identification
9 was only at the adequate level. Although consensus ratings were reached, this finding reflects the
10 subjective nature of error identification in QA. There is no consensus regarding the threshold of
11 error size for determining images that need manual edits. Error identification is further
12 complicated by variation in anatomical features and scan quality. Viewing surfaces overlaid onto
13 T1 images before skull stripping was essential to identifying errors and should be incorporated
14 into QA processes whenever feasible.

15 Errors were defined in this study as extensions of three or more voxels outside the
16 cortical boundary to limit the effects of variability in scan quality that made the barrier between
17 dura and cortex unclear. There were 24 images which had partial or full voxel extensions
18 (typically 1 voxel) that did not meet criteria for an error (Supplemental Figure 1). These
19 extensions often continued across multiple slides over the top of the cortex. Further research is
20 needed to explore the effect of small but pervasive extensions of cortex.

21 This study utilized open-access structural data from healthy, community dwelling older
22 adults. Given the bias in identification toward images with inflated measures of atrophy, these
23 findings cannot be generalized to clinical populations with increased cortical deterioration.

1 Future research is needed to examine the impact of pathological structural changes on quality of
2 reconstruction and error rates. Additionally, identification of specific anatomical features
3 associated with poor reconstruction may help to increase sensitivity and specificity for
4 statistically based QA. Future studies could use a receiver operator characteristic (ROC) curve
5 analysis to explore cutoffs, based on a balance of sensitivity and specificity, for number of total
6 outliers in identifying images with reconstruction error. Given the sex bias found in our studies,
7 ROC analyses should be conducted for men and women separately.

8 **Conclusions**

9 This study highlights the limited incremental utility of correcting errors of cortical
10 extension to assess the relationships between brain volumes and neurocognitive measures.
11 Utilizing statistically based methods alone can introduce confounds by differentially identifying
12 older, male brains for editing. This finding is especially important for researchers utilizing large-
13 scale datasets, given the resource commitment to manual QA intervention.

14

15 **Conflicts of Interest**

16 The authors declare that they have no conflicting interests. This research did not receive any
17 specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

18 **Acknowledgements**

19 The authors would like to acknowledge the following people and organizations for their
20 contributions:

21 Douglas Greve at the MGH/HST Athinoula A. Martinos Center for Biomedical Imaging for his
22 comments and consultation.

- 1 The NKI-Rockland Sample Initiative for providing the data used in these analyses (data
- 2 collection funded through NIMH BRAINS R01MH094639-01).
- 3 The Suffolk University Psychology Department for their support of doctoral students and David
- 4 Gansler's Lab.
- 5

References

- 1
2 Ahmed, B., Brodley, C. E., Blackmon, K. E., Kuzniecky, R., Barash, G., Carlson, C., & Thesen,
3 T. (2015). Cortical feature analysis and machine learning improves detection of “MRI-
4 negative” focal cortical dysplasia. *Epilepsy & Behavior*, 48, 21-28.
- 5 Chen, X., Liang, S., Pu, W., Song, Y., Mwansisya, T. E., Yang, Q., & Xue, Z. (2015). Reduced
6 cortical thickness in right Heschl’s gyrus associated with auditory verbal hallucinations
7 severity in first-episode schizophrenia. *BMC psychiatry*, 15(1), 152.
- 8 Cohen, J. (1988). Statistical power analysis for the behavioral sciences. Hillsdale, NJ: *Lawrence*
9 *Earlbaum Associates*, 2.
- 10 Collins, D. L. (1994). *3D Model-based segmentation of individual brain structures from*
11 *magnetic resonance imaging data* (Doctoral dissertation, McGill University).
- 12 Delis, D. C., Kaplan, E., & Kramer, J. H. (2001). *Delis-Kaplan executive function system (D-*
13 *KEFS)*. Psychological Corporation.
- 14 Desikan, R. S., Segonne, F., Fischl, B., Quinn, B. T., Dickerson, B. C., Blacker, D., Buckner, R.
15 L., Dale, A. M., Maguire, R. P., & Hyman, B. T. (2006). An automated labeling system
16 for subdividing the human cerebral cortex on MRI scans into gyral based regions of
17 interest. *NeuroImage*, 31(3), 968–980.
- 18 Fischl, B., Salat, D.H., Busa, E., Albert, M., Dieterich, M., Haselgrove, C., van der Kouwe, A.,
19 Killiany, R., Kennedy, D., Klaveness, S., Montillo, A., Makris, N., Rosen, B., Dale, A.M.
20 (2002). Whole brain segmentation: automated labeling of neuroanatomical structures in
21 the human brain. *Neuron* 33(3), 341-355.
- 22 Fischl, B., van der Kouwe, A., Destrieux, C., Halgren, E., Segonne, F., Salat, D.H., Busa, E.,
23 Seidman, L.J., Goldstein, J., Kennedy, D., Caviness, V., Makris, N., Rosen, B., & Dale,

- 1 A.M., (2004). Automatically parcellating the human cerebral cortex. *Cereb Cortex*, 14(1),
2 11-22.
- 3 Gur, R. C., Richard, J., Hughett, P., Calkins, M. E., Macy, L., Bilker, W. B., & Gur, R. E.
4 (2010). A cognitive neuroscience-based computerized battery for efficient measurement
5 of individual differences: standardization and initial construct validation. *Journal of*
6 *Neuroscience Methods*, 187(2), 254-262.
- 7 Kaufmann, L. K., Baur, V., Hänggi, J., Jäncke, L., Piccirelli, M., Kollias, S., & Milos, G. (2017).
8 Fornix Under Water? Ventricular Enlargement Biases Forniceal Diffusion Magnetic
9 Resonance Imaging Indices in Anorexia Nervosa. *Biological Psychiatry: Cognitive*
10 *Neuroscience and Neuroimaging*, 2(5), 430-437.
- 11 Keshavan, A., Datta, E., McDonough, I., Madan, C. R., Jordan, K., & Henry, R. G. (2017).
12 Mindcontrol: A web application for brain segmentation quality control. *NeuroImage*.
- 13 Koh, D., Lee, S., Pacheco, J., Pappu, V., & Vinke, L. (2017, January 19). *Freesurfer QA Tools*.
14 Retrieved from <https://surfer.nmr.mgh.harvard.edu/fswiki/QATools>.
- 15 Lenroot, R. K., Gogtay, N., Greenstein, D. K., Wells, E. M., Wallace, G. L., Clasen, L. S., ... &
16 Thompson, P. M. (2007). Sexual dimorphism of brain developmental trajectories during
17 childhood and adolescence. *Neuroimage*, 36(4), 1065-1073.
- 18 Lichy, M. P., Wietek, B. M., Mugler III, J. P., Horger, W., Menzel, M. I., Anastasiadis, A., ... &
19 Schick, F. (2005). Magnetic resonance imaging of the body trunk using a single-slab, 3-
20 dimensional, T2-weighted turbo-spin-echo sequence with high sampling efficiency
21 (SPACED) for high spatial resolution imaging: initial clinical experiences. *Investigative*
22 *radiology*, 40(12), 754-760.

- 1 McCarthy, C. S., Ramprashad, A., Thompson, C., Botti, J. A., Coman, I. L., & Kates, W. R.
2 (2015). A comparison of FreeSurfer-generated data with and without manual
3 intervention. *Frontiers in Neuroscience*, 9.
- 4 Meng, X. L., Rosenthal, R., & Rubin, D. B. (1992). Comparing correlated correlation
5 coefficients. *Psychological bulletin*, 111(1), 172-175.
- 6 Moore, T. M., Reise, S. P., Gur, R. E., Hakonarson, H., & Gur, R. C. (2015). Psychometric
7 properties of the Penn Computerized Neurocognitive Battery. *Neuropsychology*, 29(2),
8 235-246.
- 9 Nooner, K. B., Colcombe, S., Tobe, R., Mennes, M., Benedict, M., Moreno, A., & Sikka, S.
10 (2012). The NKI-Rockland sample: a model for accelerating the pace of discovery
11 science in psychiatry. *Frontiers in Neuroscience*, 6, 152.
- 12 Phan, K. L., Wager, T., Taylor, S. F., & Liberzon, I. (2002). Functional neuroanatomy of
13 emotion: a meta-analysis of emotion activation studies in PET and fMRI. *NeuroImage*,
14 16(2), 331-348.
- 15 Savalia, N. K., Agres, P. F., & Wig, G. S. (2015). Processing & editing overview. *The Center for*
16 *Vital Longevity*.
- 17 Seidman, L. J., Valera, E. M., & Makris, N. (2005). Structural brain imaging of attention-
18 deficit/hyperactivity disorder. *Biological psychiatry*, 57(11), 1263-1272.
- 19 U.S. Census Bureau. (2009). Census Data. US Department of Health and Human Services,
20 Washington, D.C.
- 21 Van Petten, C. (2004). Relationship between hippocampal volume and memory ability in healthy
22 individuals across the lifespan: review and meta-analysis. *Neuropsychologia*, 42(10),
23 1394-1413.

- 1 Viviani, R., Pracht, E. D., Brenner, D., Beschoner, P., Stingl, J. C., & Stöcker, T. (2017).
2 Multimodal MEMPRAGE, FLAIR, and R_2^* Segmentation to Resolve Dura and Vessels
3 from Cortical Gray Matter. *Frontiers in neuroscience*, *11*.
- 4 Ward, N. S., & Frackowiak, R. S. J. (2003). Age-related changes in the neural correlates of
5 motor performance. *Brain*, *126*(4), 873-888.
- 6 Wechsler, D. (2011). *WASI-II: Wechsler abbreviated scale of intelligence--*. Psychological
7 Corporation.
- 8 Yuan, P., & Raz, N. (2014). Prefrontal cortex and executive functions in healthy adults: a meta-
9 analysis of structural neuroimaging studies. *Neuroscience & Biobehavioral Reviews*, *42*,
10 180-192.

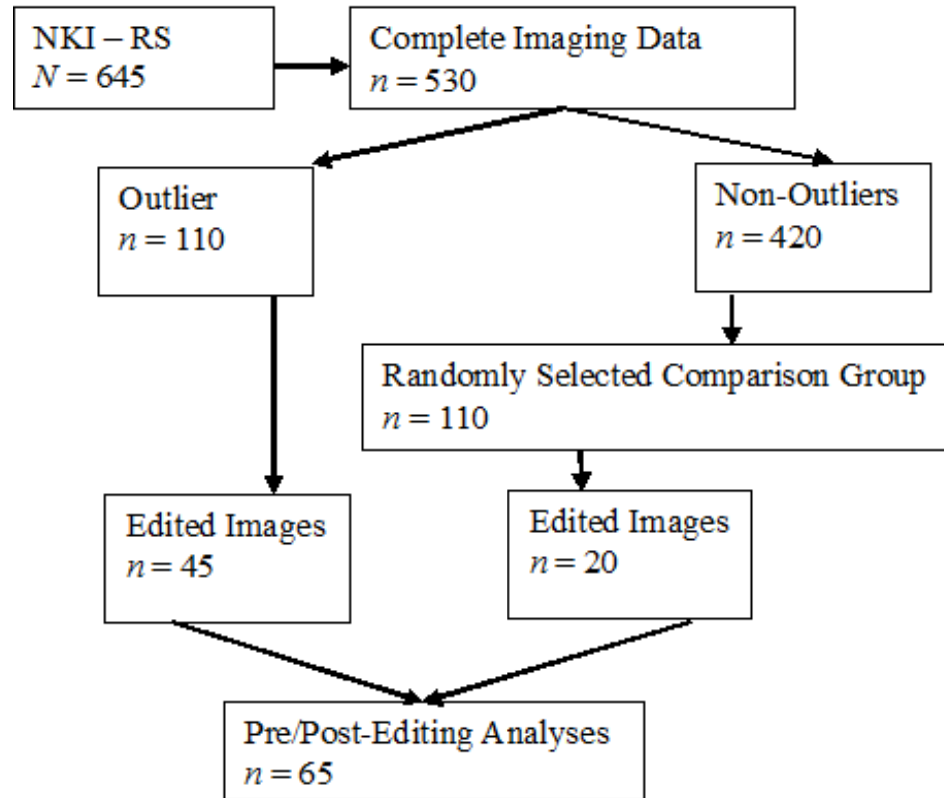


Figure 1. Flowchart of Nathan Kline Institute – Rockland Sample (NKI-RS) participants and the study sample.

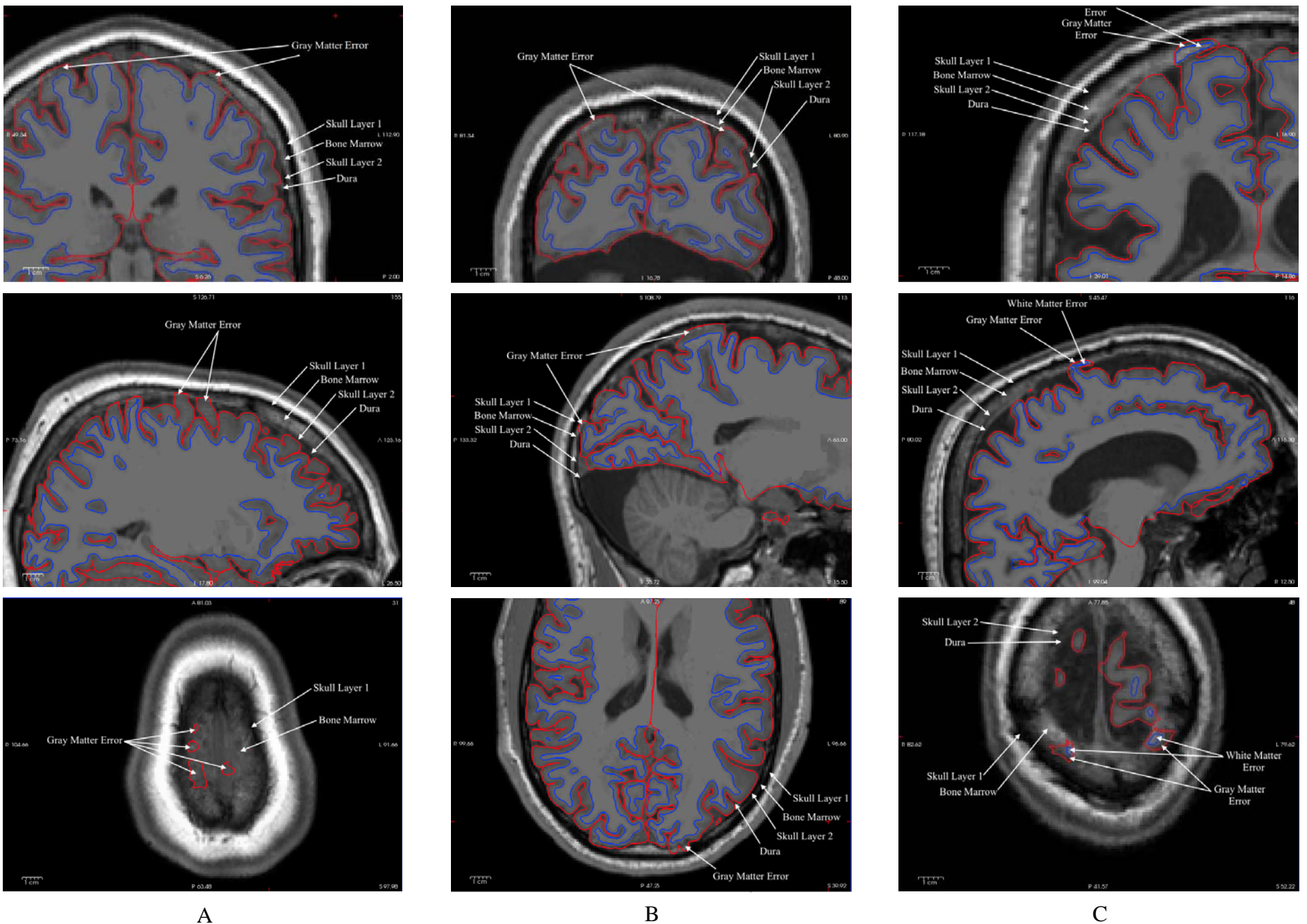


Figure 2. Coronal, Sagittal, and Horizontal View of (A) Errors of cortical extension in the pre-central/post-central region; (B) Errors of cortical extension in posterior regions; and (C) Errors of cortical extension that include both white and grey matter.