1            Word count: 4,881
2            Abstract: 250
3            Methods: 2,271
4            Figures: 6
5            Extended Data Figures: 10
6            Supplemental Tables: 12

7

8 **Uncovering and mitigating bias in large, automated MRI analyses of brain development**

9

10 Safia Elyounssi[1,2]*, Keiko Kunitoki[1,2]*, Jacqueline A. Clauss[1,2], Eline Laurent[1,2], Kristina
11 Kane[1,2], Dylan E. Hughes[1,2,3], Casey E. Hopkinson[1,2], Oren Bazer[1,2], Rachel Freed Sussman[1,2],
12 Alysa E. Doyle[1,4], Hang Lee[5], Brenden Tervo-Clemmens[1], Hamdi Eryilmaz[1,2], Randy L.
13 Gollub[1,2], Deanna M. Barch[6], Theodore D. Satterthwaite[7,8,9], Kevin F. Dowling[1,10], Joshua L.
14 Roffman[1,2]
15 *Equal authorship

16

17 [1]Department of Psychiatry, Massachusetts General Hospital and Harvard Medical School
18 [2]Martinos Center for Biomedical Imaging, Massachusetts General Hospital
19 [3]Departments of Psychiatry & Biobehavioral Sciences, University of California, Los Angeles
20 [4]Center for Genomic Medicine, Massachusetts General Hospital
21 [5]Biostatistics Center, Massachusetts General Hospital and Harvard Medical School
22 [6]Department of Psychological and Brain Sciences, Washington University in St. Louis
23 [7]Department of Psychiatry, University of Pennsylvania Perelman School of Medicine
24 [8]Penn Lifespan and Neuroimaging Center, University of Pennsylvania Perelman School of
25 Medicine
26 [8]Penn-CHOP Lifespan Brain Institute
27 [10]Department of Psychiatry, University of Pittsburgh

28

29 Corresponding author:
30 Joshua L. Roffman MD, MMSc
31 Massachusetts General Hospital
32 149 13th St, Room 2616
33 Charlestown, MA 02129
34 617-724-1920
35 jroffman@partners.org

36

45

46 No authors declare any potential conflicts of interest.

**Abstract**

Large, population-based MRI studies of adolescents promise transformational insights into neurodevelopment and mental illness risk [1,2]. However, MRI studies of youth are especially susceptible to motion and other artifacts [3,4]. These artifacts may go undetected by automated quality control (QC) methods that are preferred in high-throughput imaging studies, [5] and can potentially introduce non-random noise into clinical association analyses. Here we demonstrate bias in structural MRI analyses of children due to inclusion of lower quality images, as identified through rigorous visual quality control of 11,263 T1 MRI scans obtained at age 9-10 through the Adolescent Brain Cognitive Development (ABCD) Study[6]. Compared to the best-rated images (44.9% of the sample), lower-quality images generally associated with decreased cortical thickness and increased cortical surface area measures (Cohen's d 0.14-2.84). Variable image quality led to counterintuitive patterns in analyses that associated structural MRI and clinical measures, as inclusion of lower-quality scans altered apparent effect sizes in ways that increased risk for both false positives and negatives. Quality-related biases were partially mitigated by controlling for surface hole number, an automated index of topological complexity that differentiated lower-quality scans with good specificity at Baseline (0.81-0.93) and in 1,000 Year 2 scans (0.88-1.00). However, even among the highest-rated images, subtle topological errors occurred during image preprocessing, and their correction through manual edits significantly and reproducibly changed thickness measurements across much of the cortex (d 0.15-0.92). These findings demonstrate that inadequate QC of youth structural MRI scans can undermine advantages of large sample size to detect meaningful associations.

71  **<u>Introduction</u>**

72

73      Magnetic resonance imaging (MRI) is widely used in clinical neuroscience research to study

74  neuroanatomical variation in healthy individuals as well as those with neuropsychiatric disease [7].

75  Structural (T1-weighted) MRI scans (sMRI) provide reliable, individual-level indices of cortical

76  thickness, surface area, and volume, and enable registration of other brain imaging data (such as

77  functional MRI and PET) to anatomical templates that facilitate group-level analyses [8]. In

78  accordance with neurodevelopmental models of mental illness, large-scale brain MRI studies of

79  children and adolescents offer potential to elaborate neural signatures of emergent

80  psychopathology [1,2]. Such insights could be harnessed in efforts to develop improved early

81  recognition and treatment, outcomes that may help ameliorate the current youth mental health

82  crisis [9]. As such, the US National Institutes of Health and other funding agencies have invested

83  heavily in longitudinal MRI studies of adolescent brain development, such as the ongoing ABCD

84  Study [6].

85

86      Recent work has underscored the need for thousands of participants in such clinical MRI

87  studies, as within-group variation is considerable and effect sizes for relationships between

88  psychopathology and MRI indices tend to be small [2]. Further, MRI scans of children and

89  adolescents are particularly susceptible to artifact due to participant motion within the scanner [3,4].

90  An unanswered question concerns whether large sample size – e.g. in studies involving

91  thousands of participants – sufficiently compensates for errant sMRI measurements arising from

92  inclusion of poorer quality images. Alternatively, smaller studies have suggested the possibility

3

93    that visible motion artifact results not only in random noise but in bias [3,4], which (again) may or

94    may not be sufficiently offset by inclusion of more participants.

95

96       A related question concerns the adequacy of automated quality control (QC) measures,

97    applied during scan acquisition, processing, or analysis, to identify or adjust for poor quality

98    images in large sMRI studies of children.  Notably, unlike for functional MRI, head motion is

99    less routinely quantified as part of sMRI analyses, and its effects on sMRI measurements have

100    been less well studied – although some prior work has associated induced or measured motion

101    with bias in sMRI estimates [3,10].  Newer sMRI sequences, including those deployed on Siemens

102    and GE magnets in ABCD [11,12], have incorporated real-time motion correction protocols that re-

103    acquire data immediately after significant motion is detected.  Whether this feature mitigates

104    artifact sufficiently well to prevent bias in large-scale studies remains uncertain.  Image

105    preprocessing software can provide automated QC metrics, such as the overall "pass/fail" rating

106    in the FreeSurfer processing stream. [13] This metric is used by ABCD in conjunction with raw

107    data screens and clinical (radiology) evaluations to provide an overall recommendation on

108    whether to include images in analyses.  However, routine automated QC measures have shown

109    inconsistent sensitivity to detect artifact identified by manual (visual) QC ratings of sMRI scans

110    in youth[5,14].

111

112       As such, a final consideration – one especially pertinent to large-scale studies such as ABCD,

113    which is collecting 6 sets of MRI scans over 10 years from >10,000 youth participants – is the

114    added value of manual QC of postprocessed sMRI scans, and of the even more time- and

115    resource-intensive process of manual cortical edits[15,16], to minimize artifact-related errors.

116    Depending on image quality, manual edits of a single scan can take a skilled technician as few as

117    30 minutes to as long as several days to complete.  While the utility of manual edits in

118    identifying case-control differences in pediatric sMRI studies has been questioned [17–19], their

119    importance to accurately detecting subtle neurodevelopmental differences among youth is

120    evident in other studies[20,21].

121

122        Here we conducted in-depth, manual QC assessments of >12,000 sMRI images obtained at

123    Baseline (age 9-10) and Year 2 follow-up (age 11-12) from ABCD study participants.  We then

124    characterized the impact of poorer-quality scans on the fidelity of sMRI measurements (cortical

125    thickness, surface area, and volume) and the on reliability sMRI-clinical associations.  Further, in

126    light of efficiency considerations, we evaluated the sensitivity of automated QC to detect poorer-

127    quality scans, and contrasted the effectiveness and reliability of several automated and manual

128    error mitigation strategies that varied by labor intensity.

129

130    **Results**

131

132    *Manual quality control (MQC) ratings in Baseline scans*

133

134    The ABCD study enrolled 11,875 participants, age 9 or 10 at Baseline, across 22 U.S. sites.

135    Participant race and ethnicity mirrored those of the U.S., and enrollment was enriched for

136    multiple births and siblings from multiple pregnancies [22]. Structural MRI (sMRI) scans were

137    obtained from participants on 3T Siemens, Philips, or GE magnets as described in **Methods** and

138    by Casey and colleagues.[23] Minimally processed T1 volumes were available from the NIMH

139    Data Archive (NDA) for all but 160 participants. After removing those marked as requiring

140    clinical consultation, T1 volumes for the remaining scans were downloaded from the NDA and

141    pre-processed in FreeSurfer version 7.1. While several processing streams are available to

142    process and analyze sMRI data, the present analyses used FreeSurfer software for two reasons:

143    first, existing, tabulated region-of-interest sMRI analyses available through the NIMH Data

144    Archive and widely used in published ABCD analyses were conducted wither FreeSurfer; and

145    second, FreeSurfer offers manual cortical edit capabilities. Following training and calibration

146    (**Methods**), a single research coordinator (S.E.) who was blind to subject-level information then

147    viewed each MRI volume individually. This approach was chosen because it eliminated

148    concerns over inter-rater reliability, which has previously been shown to be modest (~0.75) when

149    including multiple tiers of sMRI QC,[5] but could alternatively be assessed for intra-rater

150    reliability (e.g., drift in ratings over time) and for triangulation with automated QC measures.

151    During manual review an additional 740 scans were removed from further consideration due to

152    the presence of cysts $>1$ cm$^3$, and 228 were omitted from the main analyses due to segmentation

153    errors and related signal dropout that persisted after a second round of preprocessing (**Figure**

154    **1a**).

155

156        The remaining 10,295 T1 scans received manual quality control (MQC) ratings.  Ratings were

157    based on overall appearance of the entire T1 volumes, as follows: "1" (requiring minimal edits,

158    n=4,630, 45.0%), "2" (requiring moderate edits, n=4,063, 39.5%), "3" (requiring substantial

159    edits, n=1,383, 13.4%), or "4" (unusable, n=219, 2.1%) (**Figure 1b,c**).  We have uploaded these

160    MQC ratings to the NDA (see **Data Availability**).  Demographic, clinical, and scanner

161    characteristics of participants stratified by MQC group are described in **Table S1a**.  Individuals

162    with higher quality scans tended to be slightly older and female, also demonstrated less

163    externalizing psychopathology and total symptoms on the Child Behavior Checklist (CBCL).

164    Scan quality also differed by scanner manufacturer; notably, the mean MQC rating for images

165    from Philips magnets (1.34, 95% CI 1.29-1.38), which were not subject to real-time motion

166    correction, was more favorable than those for Siemens (1.71, 95% CI 1.69-1.73) and GE (1.96,

167    95% CI 1.93-1.99), which did include this feature (p's<.0001, after controlling for age, gender,

168    and psychopathology).  MQC ratings were stable over the sequence of scan evaluations after

169    controlling for each of these factors (see **Extended Data Figure 1a**, **Table S2**), and their

170    distribution did not change in sensitivity analyses that included the 228 scans with segmentation

171    errors (**Extended Data Figure 2, Table S1b**).

172

173        All scans had also received automated quality control ratings (pass/fail), available as part of

174    the ABCD NDA.  Of the 10,295 scans with MQC ratings, all but 325 were designated as

175    recommended for use; these 325 fell disproportionately within higher MQC groups (comprising

7

176    0.4% of MQC=1 scans, 1.4% of MQC=2, 10.6% of MQC=3, and 48.9% of MQC=4) but this

177    designation missed numerous poorer-quality scans. Subsequent analyses including these 10,295

178    scans were all adjusted for participant age, gender, total intracranial volume, study site, and

179    scanner manufacturer; region-of-interest (ROI)-based analyses were further covaried by family

180    ID to control for effects of participant relatedness.

181

182    *Associations between MQC ratings and cortical structure, and comparison with surface hole*

183    *number (SHN)*

184

185       Automated measures of cortical thickness, surface area, and volume are commonly used to

186    identify case-control differences or as predictors of dimensional measures (e.g.,

187    psychopathology) in psychiatric neuroimaging research [8].  We next determined the extent to

188    which MQC ratings associate with variance in these measures, as determined by FreeSurfer.

189    MQC ratings associated linearly with reduced thickness across much of the cortical mantle

190    (**Figure 2a**), with increased cortical surface area in lateral/superior and reduced surface area in

191    medial/inferior regions (**Figure 2b**), and heterogeneous effects on cortical volume (**Figure 2c**).

192    Pairwise comparisons of best quality (MQC=1) versus lower quality (MQC=2, 3, and 4) images

193    demonstrated increasingly strong effects on each structural index as QC ratings worsened, with

194    moderate to strong effect sizes noted in numerous cortical regions (see also **Table S3a, b, c** for

195    effects of MQC rating differences in each of the 68 cortical ROI defined by the Desikan-Killiany

196    Atlas).  For example, comparison of cortical thickness values between MQC=1 versus MQC=2,

197    3, and 4 yielded a total of 39, 55, and 61 ROI (of 68), respectively, with statistically significant

198    differences (FDR q <.05).  Regions demonstrating stronger effects of poor quality control on

8

199 thickness included, but also extended beyond, those identified as showing similar effects in a

200 previous, smaller study of adolescent and adult participants (n=1,840) [5], in consistent directions

201 (e.g., increased thickness in numerous lateral ROIs, decreased thickness in medial occipital and

202 posterior cingulate cortices). Subcortical volumes also differed significantly based on MQC

203 rating, with higher ratings generally associated with smaller volumes (**Table S4**).

204

205     We next compared the performance of an automated QC measure, the surface hole number

206 (SHN), to manual (MQC) ratings. SHN reflects the Euler number, which measures continuity of

207 tessellated images (e.g., those that contain continuous triangular structures, as do FreeSurfer-

208 generated maps of the cortical surface, see **Methods**) based on the sum of the vertices and faces

209 subtracted by the number of faces. Higher SHN have predicted worse manual quality control

210 ratings in previous MRI studies and have been proposed as an automated quality control index

211 for use in high-throughput neuroimaging studies, outperforming other measures (such as signal-

212 to-noise ratio and motion during functional MRI scans conducted during the same scan session)

213 [5,16]. We calculated SHN for each available Baseline and Year 2 scan using FreeSurfer 7.0 and

214 have uploaded the data to the NDA (see **Data Availability**). SHN increased in tandem with

215 MQC ratings (rho=0.59; mean SHN differed between all MQC level pairs, p≤1.02E-121), and

216 linear associations of SHN with differences in cortical thickness, surface area, and volume

217 (**Figure 3a,b,c**) closely resembled those of MQC (**Figure 2**). Distribution of SHN values among

218 MQC groups was stable over the temporal sequence of MQC evaluations (**Extended Data**

219 **Figure 1b**).

220

221     We then examined whether including SHN as an additional covariate mitigated effects of

222     variable scan quality on sMRI indices, as defined by differences in measurements between

223     MQC=1 and MQC=2, 3, and 4 respectively (**Table S3d,e,f; Figure 3d,e,f**).  Depending on the

224     specific comparison (MQC=1 vs. 2, 3, or 4), inclusion of SHN reduced the effect size (Cohen's

225     d) of manual quality control-related differences in cortical thickness by 42 to 59%; reductions in

226     effect size for cortical surface area ranged from 39 to 57%, and for cortical volume from 16 to

227     62% (**Table S5**). Meaningful effects of SHN correction can also be demonstrated by comparing

228     the number of ROIs that showed statistically significant (FDR, q<.05) effects of MQC ratings

229     before versus after including SHN as a covariate.  For example, among 39 ROIs exhibiting

230     differences in cortical thickness between MQC=1 and MQC=2 before covarying for SHN, 17 fell

231     out of significance after covarying for SHN, while 1 ROI became newly significant.

232

233     We then used SHN data in concert with MQC ratings to develop and assess the reliability of a

234     tiered, automated sMRI QC rubric to classify the quality of individual scans.  This rubic assigned

235     scans to 4 levels akin to the MQC groups, but based exclusively on SHN-based thresholds, so

236     that these ratings could be applied even in the absence of manual QC.  **Figure 3g** displays the

237     distribution of SHN among MQC groups.  Using receiver operating characteristic (ROC) curve

238     analyses, we derived 3 optimized SHN thresholds to isolate poorer-quality scans (**Figure 3h**).

239     The most conservative threshold eliminated scans with MQC ratings of 2 or higher, based on an

240     SHN cutoff of 29.5 (sensitivity=0.81; **Figure 3i**).  The next threshold eliminated scans with

241     MQC ratings of 3 or higher, based on a SHN cutoff of 36.5 (sensitivity=0.81; **Figure 3j**). The

242     most liberal threshold eliminated scans with MQC ratings of 4, based on an SHN cutoff of 62.5

243     (sensitivity=0.93, **Figure 3k**).

244

245     These 3 thresholds defined 4 SHN groups (tiers A-D), that in turn associated with linear

246     effects on sMRI indices (**Extended Data Figure 3**).  The linear effects of SHN tiers closely

247     approximated the linear effects of MQC groupings (**Figure 2**), as well as those of continuous

248     SHN values (**Figure 3a,b,c**).  Still, MQC and SHN each accounted for distinct variance in scan

249     quality as seen in **Extended Data Figure 4**.  In a sensitivity analysis, inclusion of scans with

250     FreeSurfer segmentation errors (n=228) did not substantially alter either the distribution of SHN

251     across MQC ratings or optimal boundaries between SHN tiers in ROC analyses (**Table S6**).

252

253     ***SHN tiers as predictors of MQC in Year 2 follow-up scans***

254

255     Evaluation of Year 2 scans from ABCD enabled us to test the reliability of SHN tiers derived

256     from Baseline scans.  A total of 6,941 minimally processed Year 2 T1 volumes were available

257     through the ABCD Data Archive, after removing those that did not meet inclusion criteria for

258     Baseline analysis; see **Extended Data Figure 5**.  Following preprocessing in FreeSurfer 7.0 and

259     extraction of SHN, 1,000 sMRI volumes were semi-randomly selected such that they included

260     (1) a range of scan quality, operationalized by ensuring a mix of tiers A, B, C, and D; and (2) a

261     distribution of magnet types (Siemens, Philips, GE) that was equivalent to the analyzed Baseline

262     sample.  Of note, Year 2 scans showed better overall quality than Baseline scans, with 83.9%

263     falling into SHN tier A (**Extended Data Figure 6a**; compare to **Figure 3g**, where 57.3% of

264     Baseline scans fell into tier A).  Group characteristics of SHN tiers A to D in the Year 2 sample

265     are described in **Table S7**.  One scan was discarded due to presence of a large cyst.  Only 168

266     Year 2 scans fell within SHN tier D, all of which were used for the analysis.  The selected Year 2

11

267    scans underwent MQC ratings by 2 trained and calibrated research coordinators (500 scans

268    randomly disbursed to each of K.A.K. and E.L; see **Methods**), using the same method as

269    Baseline MQC ratings.

270

271        **Table S8** describes the performance of SHN tiers in predicting MQC ratings for the 999 Year

272    2 scans.  The SHN tiers effectively filtered out scans with higher MQC ratings, with sensitivity

273    ranging from 0.87 (for differentiating scans rated 2 and higher from those rated 1) to 1.00 (for

274    differentiating scans rated 4 from those rated lower).  **Extended Data Figure 6b** shows the

275    distribution of MQC ratings within each SHN tier.  **Extended Data Figure 7** indicates the effect

276    of SHN tiers on sMRI indices across all 6,941 Year 2 scans (most of which had not received

277    MQC ratings); comparison to **Extended Data Figure 3** affirms that SHN tiers reproducibly

278    tracked variance in scan quality, especially in regard to cortical thickness and surface area.

279

280    ***Scan quality and risk for error in applied sMRI analyses***

281

282        sMRI measures are frequently explored for associations with clinical and developmental data.

283    The ABCD Study provides an unprecedented opportunity in this regard, with multiple imaging

284    and clinical measurements obtained within the same youth participants over a 10-year period.

285    However, given the tendency of poorer quality images to bias sMRI measurements among youth,

286    we next examined the extent to which unaccounted variance in scan quality might affect

287    associations between MRI and clinical indices.

288

289      As a positive control, we first considered a well-established relationship between age and

290    cortical thickness. Most of the cortex is known to thin linearly during adolescence, as seen in

291    smaller but well quality-controlled samples [24]. As points of reference, we compared age-

292    thickness effects in the SHN-corrected, MQC=1 sample (n=4,617, "ground truth") to those in the

293    full, non-corrected sample (n=10,257, "full non-QC-adjusted sample"). Significant age-

294    thickness relationships were readily observed, even cross-sectionally between ages 9.0 and 10.9,

295    within the full non-QC-adjusted sample (**Figure 4a**) – although note the considerably smaller

296    effect size of age on thickness compared to that of quality control ratings (**Figure 2a**). Despite

297    these smaller effects, age-thickness effects were sufficiently robust to be detected within the

298    smaller ground truth sample: among 68 cortical ROIs, significant (FDR q<.05) negative

299    associations were present in 59 regions, regardless of SHN adjustment (**Table S9**). Notably,

300    though, several of these ROI did *not* show significant age-thickness differences in the (larger)

301    full unadjusted sample – but then *regained* significance in the full sample after SHN adjustment.

302    As such, inclusion of SHN mitigated Type II error (i.e., false negatives) that would have

303    otherwise occurred in the full non-QC-adjusted sample, albeit for only a small number of

304    regions.

305

306      The risk of Type II error arising from non-quality-corrected images can also be appreciated in

307    **Figure 4b**, which plots effect sizes for the age-thickness relationship across all 68 ROIs. To

308    facilitate comparisons across MQC levels, ROIs were rank-ordered (left-to-right) by effect size

309    among the 1-rated scans. Effect sizes generally diminished as poorer quality images were

310    iteratively included (2s, then 3s, then 4s) in the analysis. These results echo a prior, smaller

13

311      analysis (n=1,598, mean age=15.0), wherein poorer quality scans associated with blunted effects

312      of age on cortical thickness [5].

313

314      Next, we considered a more exploratory relationship between dimensional psychopathology

315      and cortical volume. Several groups have reported inverse associations between CBCL scales

316      and cortical volume, including using ABCD data [25,26]. In a recent study focused on genetic and

317      neurodevelopmental underpinnings of psychopathology in ABCD [27], among the broadband

318      CBCL scales (total, internalizing, externalizing) we identified externalizing symptoms

319      (CBCLext) as most strongly related to cortical volumes at Baseline, after taking into account

320      both MQC ratings and SHN.

321

322      In the full, non-QC-adjusted sample, CBCLext scores showed a diffuse, inverse relationship

323      with volume across the cortical mantle (**Figure 4c**), although effect sizes were smaller than for

324      the age-thickness relationship (**Figure 4a**). Within this larger sample, 43 ROIs demonstrated

325      significant (FDR q <.05) relationships between CBCLext and volume (**Figure 4d**, **Table S10**).

326      However, stark differences emerged in comparison to the ground truth (MQC=1) sample,

327      wherein only 3 regions demonstrated significant CBCLext-volume relationships. Stepwise

328      analyses that gradually increased the stringency of QC suggested that this drop in the number of

329      significant ROIs reflected an interplay of QC and power considerations (**Figure 5**). While effect

330      size should not depend on sample size, unlike in the age-thickness analysis, inclusion of lower-

331      quality images resulted in substantial inflation of CBCLext-volume effects, and accordingly

332      Type I errors (i.e., false positives). Numerous ROIs showed statistically significant CBCLext-

333      volume relationships only when MQC=3 and 4 scans were included in the analysis – and, even

14

334    after correction with SHN, these regions demonstrated inflated effect sizes due to inclusion of

335    lower quality scans.  Further, regions with smaller effect sizes in the ground truth sample were

336    more likely to show inflated effect sizes – and, hence Type I error – in the full, non-adjusted

337    sample (**Figure 4d**).  This result was counterintuitive, given that large sample size is often

338    invoked to *reduce* risk of *Type II* error, through improved power to detect small but true effects.

339

340    Further complexity emerged among the 21 ROIs that showed significant CBCLext-volume

341    relationships after including only MQC ≤ 2 images.  For some regions, such as left superior

342    temporal, left precentral, and bilateral postcentral, effect size remained relatively stable as

343    inclusion thresholds loosened (**Figure 5c**).  This pattern suggests that effect sizes were not

344    inflated by artifact, and that failure to reach significance when using only MQC=1 scans

345    reflected a lack of statistical power (i.e., Type II error) – even with a sample size of >4,500.  For

346    others, such as right middle temporal, bilateral insula, and bilateral superior frontal, effect size

347    increased substantially as the MQC inclusion threshold was relaxed to 2 (or higher), likely

348    reflecting artifact.  For these regions, inclusion of even relatively good quality (but not the

349    highest quality) images appeared to result in Type I error (**Figure 5d**).

350

351    ***Effects of manual edits on sMRI indices***

352

353    Image reconstruction errors can influence sMRI measurements and can be exacerbated by

354    head motion and other artifacts [3,4].  These errors include skull strip errors, segmentation errors,

355    intensity normalization errors, pial surface misplacement, and topological defects.  Within

356    FreeSurfer these errors can be corrected through manual editing of voxels in brain and white

15

357     matter masks, watershed thresholds, and the addition of control points [15,16]. Here, we examined

358     effects of manual edits on sMRI indices among scans with relatively higher image quality, to

359     assess whether this intervention might safely be reserved for those with MQC >2.

360

361        A total of 150 Baseline scans with MQC=1 and n=30 Baseline scans with MQC=2 were

362     randomly selected for manual edits by a trained coordinator (see **Methods**). Direction and effect

363     sizes of pre-to-post edit changes across the cortical mantle are displayed in **Figure 6** (MQC=1

364     and 2 combined, n=180) and **Extended Figure 8a** and **b** (MQC=1 and 2 separately), while ROI-

365     level changes across the entire sample of 180 are described in **Table S11a,b,c**. Effects of

366     manual edits were most pronounced for cortical thickness and volume, both of which tended to

367     decrease. These changes reached statistical significance (FDR q <.05) for cortical thickness in

368     40 regions (Cohen's d range 0.16 to 0.92), and for cortical volume in 28 regions (d range 0.18 to

369     0.73). Numerous regions with signal across all scans (MQC=1 and 2) demonstrated stronger

370     effects of editing on cortical thickness and volume in MQC=2 scans than MQC=1 scans (e.g.,

371     bilateral parahippocampal, caudal middle frontal, and superior parietal cortices). Further,

372     cortical volume maps revealed a strong effect of edits in the area of the superior sagittal sinus,

373     particularly impacting superior parietal cortex (**Extended Figure 8c**). In an applied analysis, we

374     then examined the degree to which cortical edits affected effect size for the relationship between

375     cortical thickness and age. Across all 68 cortical ROIs, effect size slightly strengthened (became

376     more negative) for post-edited images compared to pre-edited images (t=2.31, p=0.024, d=-

377     0.10).

378

379      To put these findings in context along with MQC rating effects on sMRI indices, **Extended**

380    **Figure 9** maps all ROIs that showed significant (FDR, q <.05) effects of MQC, surface edits, or

381    both, as well as their direction, among Baseline scans with MQC=1 or 2. Note that even when

382    constrained to the best two scan quality groups, there are diffuse effects of scan quality

383    differences across the cortex, for each of the sMRI indices; and that biases related to poorer

384    overall quality control and to subtle topological defects can induce opposing effects on sMRI

385    measurements.

386

387      Finally, to assess reproducibility and developmental specificity of cortical edit effects, we

388    compared ABCD results to that of a second, non-overlapping MRI cohort of 292 youths, age 8 to

389    18, who received MRI scans that were assessed by radiology reports as free of pathology at

390    Massachusetts General Hospital (MGH; **Table S12a, b, c**). This sample was previously

391    described in an analysis relating prenatal folic acid exposure to cortical development[20]. This

392    sample differed from ABCD by its inclusion of (1) clinical rather than research participants, (2)

393    *all* editable images (not just those of relatively high overall quality), (3) a mix scanner field

394    strengths (1.5 and 3T) as well as manufacturers, and (4) a broader age range. Despite these

395    differences, of the 40 regions demonstrating significant effects of manual edits on thickness in

396    ABCD, 18 again showed nominally significant (15 showed FDR-significant) effects of edits in

397    the same direction within the MGH MRI cohort (Cohen's d range 0.12 to 0.98). Notably, across

398    these 18 regions, differences in pre-to-post edit mean thickness were greater at age 8-10

399    compared to other age groups (11-12, 13-14, and 15-17; omnibus F=8.49, p=0.0001, post hoc

400    comparisons p's≤.0002; **Extended Figure 10a**). Similarly, the standard error of pre-to-post

401    thickness changes across individuals was also greatest at age 8-10 (omnibus F=64.53, p=2.25E-

402    17, post hoc comparisons of age 8-10 vs. other groups, p's≤6.53E-10).  Finally, the effect of edits

403    on the relationship between age and cortical thickness differed among age groups (F=21.54,

404    p=3.88E-12); specifically, the effect of edits on the age-thickness relationship was stronger at

405    age 8-10 (d=-1.18) than for any other age group (p's≤7.73E-09; **Extended Figure 10b**). These

406    results indicate that manual edits result in replicable, diffuse changes in cortical thickness in

407    early adolescence that can influence effect sizes in clinical-MRI associations, but also suggest

408    that effects of edits become less pronounced later in adolescence.

409

**Discussion**

*Implications for brain-wide association studies in youth*

410
411

412

413

414 These present findings identify nuances related to scan quality in large pediatric brain MRI

415 cohorts that are pervasive and complex, and that likely require multi-pronged intervention to

416 avoid error in MRI-based analyses. Leveraging one of the largest collections of uniformly

417 collected sMRI data from children and adolescents, we used manual quality control (MQC) to

418 separate high quality scans and contrast them to those with various degrees of observable

419 artifact. While inclusion of lower-quality scans diminished variance in estimates of widely used

420 sMRI metrics, such as cortical thickness and surface area, they also introduced substantial bias.

421 These effects were partially mitigated by inclusion of surface hole number (SHN), an automated

422 measure of topological complexity that accounted for quality-related variance in sMRI measures

423 akin to MQC. However, inclusion of SHN failed to safeguard against most Type I and II errors

424 when poorer quality scans were included in applied analyses that associated sMRI measures with

425 clinical data. Further, even among the highest quality scans, manual editing associated with

426 significant changes in cortical thickness and surface area -- changes that in some regions were

427 oppositely signed to those observed when controlling for SHN or MQC, and that replicated in a

428 non-overlapping clinical cohort. As a whole, these results challenge assumptions that large

429 sample size alone improves sensitivity to detect valid brain-behavior relationships, or mitigates

430 the effects of variable image quality on error risk.

431

432 Implications of these findings extend not only to studies that map trajectories of healthy and

433 aberrant brain development, but also to applied analyses that relate structural indices to clinical

19

434    measures.  Comparison of effect sizes for sMRI-clinical relationships (**Figure 5**) to those of bias

435    related to poor scan quality (d=0.14-2.84) or manual edits (d=0.15-0.92) – which are generally

436    higher by an order of magnitude – demonstrates the susceptibility of these relationships to

437    artifact.  Recent analyses illustrate the need to include thousands of individuals in brain-wide

438    association studies [28,29], reflecting the small effect sizes intrinsic to these relationships.  Here,

439    inclusion of the best quality (MQC=1, n=4,617) scans was inadequate to detect relationships

440    between cortical volume and externalizing psychopathology in several regions, effects that

441    became statistically significant when scans of marginally lower quality (MQC=2, n=4,057) were

442    included.  However, further inclusion of even lower quality (MQC=3 and 4, n=1,585) scans

443    resulted in statistically significant but errant associations between volume and externalizing

444    psychopathology in dozens of brain regions, as demonstrated by markedly inflated effect sizes

445    compared to analyses that included exclusively higher quality scans.  These results indicate

446    complex trade-offs between sample size and scan quality that warrant careful consideration in

447    large MRI studies, especially in the setting of small effect sizes.

448

449    ***Quality control in the era of "big data" sMRI studies***

450

451     Large and diverse study samples offer clear advantages such as statistical power and

452    improved generalizability, and in the case of psychology and neuropsychiatry research, such

453    designs help to mitigate well-described problems of publication bias and reproducibility failures

454    [28,30].  However, several pitfalls within "big data" science have also been described, including

455    inadequate control for multiple comparisons, sampling bias, measurement error, and

456    discrepancies between statistical and clinical significance.  These issues have hampered other

457     areas of clinical-translational research, such as electronic health record, epidemiology, and health

458     services studies [31]. With regard to brain imaging research, a recent study [32] used theoretical data

459     to model trade-offs of increasing sample size well into the thousands, demonstrating the risk of

460     latent bias to outweigh the benefit of reduced variance. This concern bears out in the present

461     real-world analysis, which cautions against equating data quantity and quality in youth sMRI

462     studies [33]. The findings also have important implications for large-scale MRI studies of other

463     populations where head motion occurs more frequently, including those with psychiatric and

464     neurological disorders [34], and those at the extremes of age [35,36].

465

466     Beyond best practices to minimize participant motion[23], the present findings suggest that

467     relatively labor-intensive approaches – visual QC and manual editing – conducted in concert

468     with automated measures such as SHN – provide the best protection against errant sMRI findings

469     in youth cohorts. However, manual edits pose their own challenges with regard to feasibility in

470     studies with tens of thousands of participants, as editing each poorer quality scan can require

471     hours of personnel time. The present analyses offer several QC alternatives that may be weighed

472     in the context of available resources, the nature of particular findings, and the characteristics of

473     the study population. For example, the SHN benchmarks identified and validated in this report

474     provide an alternative QC approach that is imperfect but less time-consuming. Investigators who

475     find associations of sMRI indices with behavioral, genetic, or environmental factors described in

476     ABCD and other neurodevelopmental cohorts may wish to consider both the effect sizes and

477     anatomical distribution of these associations in deciding which QC approach is appropriately

478     conservative. The present analyses of older adolescents (ABCD Year 2, MGH) are encouraging

479     and suggest that less intervention may be needed with advancing participant age. Further, as

480     automated methods continue to gain sophistication [37,38] they may continue to improve the

481     efficiency of QC and further strengthen causal inference in neurodevelopmental MRI research.

482

483

484    **Methods**

485    *Sample from ABCD*

486    The ABCD Study has collected data from 11,875 children from 22 sites across the United

487    States. Primary analyses used baseline data from children aged 9-10 years old.  Institutional

488    Review Board (IRB) approval for the ABCD study is described in Auchter et al. [39] All parents

489    provided written informed consent and all youth provided assent. We excluded subjects whose

490    baseline MRI scans were flagged for clinical consultation (N=451), and those without available

491    T1 data (N=160) from all analyses.

492    *MRI acquisition*

493    All MRI images were obtained using harmonized parameters with 3T MRI scanners

494    manufactured by Siemens, Philips, or GE.  We used T1 weighted images (256x256 matrix,

495    slices=176-225, TR=6.31-2500, TE=2-2.9, 1x1x1 mm resolution) for our analysis.  Images

496    acquired from Siemens and GE scanners included real-time motion detection using volumetric

497    navigators that automatically triggered re-scans [11,12].  Additional details of MRI sequences are

498    described elsewhere [23].

499    *Image processing*

500    Minimally processed baseline T1 images from 11,264 participants, and year 2 follow-up T 1

501    images from 6,941 of these participants, were downloaded from the ABCD Data Archive

502    (release 4.0).  Scans underwent N4 field bias correction to correct low frequency intensity non-

503    uniformities or field bias [40]. Subsequently whole brain processing and analyses were conducted

23

504     using FreeSurfer version 7.1 (http://surfer.nmr.mgh.harvard.edu/). One baseline scan failed

505     Freesurfer processing. Using automated segmentation (Desikan-Killiany atlas), cortical

506     thickness, surface area, and volume of 68 regions of interest (ROI) were extracted, as were 20

507     subcortical volumes.

508     ***Manual quality control (MQC)***

509     A single, trained rater (S.E.) conducted visual assessment of all processed Baseline scans.

510     The rater was blinded to any potential identifying, clinical, or demographic information

511     regarding participants.  This rater had been trained by the PI (J.L.R.) and a clinical research

512     coordinator (K.F.D.) who had experience conducting manual edits of >300 MRI scans acquired

513     from children and adolescents aged 8 to 18 [20]. The system of rating was developed by using a

514     randomly selected set of 200 baseline T1 scans.  The final manual quality control (MQC) ratings

515     scheme, developed in consensus with S.E., K.F.D., and J.L.R., included 5 categories:  A rating of

516     "1" was given to scans of minimal artifact, only needing about ½ hour to complete edits. A rating

517     of "2" was given to scans with moderate artifact, requiring 1-2 hours for manual edits. A rating

518     of "3" was given to scans with substantial artifact, requiring several hours of edits. A rating of

519     "4" was given to scans with severe artifact, such as motion artifact, and would not be possible to

520     fix with manual edits. Lastly, a rating of "5" was given to scans with a processing defect which

521     resulted in segmentation errors and apparent loss of tissue.  Scans that included cysts that were

522     greater than 1 cm$^3$ were not rated and excluded from subsequent analysis.   The order in which

523     scans were evaluated for MQC was semi-random.  Scans originating from N=5,105 participants

524     of European ancestry were prioritized and randomly sequenced to facilitate a genomic analysis

525     [27].  However, this initial group also contained 373 randomly interspersed scans from randomly

24

526   selected non-European participants.  Following assessment of this initial set of 5,105 scans, the

527   remainder were evaluated in random order.  Of the evaluated scans, 368 were coded within the

528   ABCD NIMH Data Archive as "inclusion not recommended" based on an automated overall QC

529   measure in the FreeSurfer preprocessing stream and/or corrupted raw data at the time of scan

530   acquisition (imgincl_t1w_include=0); the remainder received the "inclusion recommended"

531   code.

532   ***Characterizing apparent tissue loss due to segmentation errors***

533   MQC=5 scans (N=228) were re-rated as MQC from "1" through "4" to assess the quality of

534   the remaining volume that was unaffected by segmentation errors. Ratings were performed by

535   the same trained rater who assigned ratings to all baseline scans. The sagittal, coronal, and axial

536   extents of the drop out region were measured in Freeview

537   (https://surfer.nmr.mgh.harvard.edu/fswiki/FreeviewGuide). Approximate volumes of

538   segmentation error-related tissue loss were calculated assuming an ellipsoid shape and measured

539   x, y, and z dimensions. For purposes of displaying location and overlap of drop-out across scans,

540   rectangular cuboids were constructed in MarsBar in SPM 12 (http://www.fil.ion.ucl.ac.uk/spm/)

541   using measured dimensions and coordinates. Rectangular cuboids were combined across subjects

542   and were thresholded by a whole brain mask in MarsBar. Areas of dropout were thresholded at

543   n>10 subjects with dropout in that region and drop-out was displayed on an exemplar structural

544   image in xjView (https://www.alivelearn.net/xjview).

545   ***Surface hole number***

546     We used surface hole number (SHN) as an automated quality control measure extracted from

547     FreeSurfer aparc tabulated data. SHN is a topological measurement referring to geometrical

548     holes (imperfections) in the tessellated brain surface. SHN is related to the Euler number by the

549     formula, Euler number = 2 - 2 × SHN. Previous, smaller studies have suggested that SHN can

550     serve as a proxy for overall T1 scan quality [5,16]. Here, SHN from baseline scans were used to

551     determine optimal proxies for MQC, through creating of 4 tiers (A, B, C, D) that approximated

552     the 4 levels of MQC ratings (1, 2, 3, 4).

553     *Psychopathology measurement*

554     We used the parent-reported Child Behavior Checklist (CBCL) as a measure of dimensional

555     psychopathology. The CBCL is a frequently used scale comprising eight subscales

556     (anxious/depressed, withdrawn/depressed, somatic, social, thought, attention, rulebreaking, and

557     aggressive symptoms) that can be summarized by total, internalizing, and externalizing scores.

558     Raw scores are converted to t-scores which are normed for age and gender [41].

559     *Year 2 T1 replication*

560     We examined all available Year 2 T1 weighted images to assess the reliability of SHN tiers

561     derived from Baseline scans. The most recent ABCD data release (4.0) contains Year 2 scans

562     from 7,829 participants. Using the same method as for Baseline scans, we used FreeSurfer to

563     process images from 6,941 individuals whose baseline image passed the inclusion criteria and

564     received MQC ratings of 1-5. SHN were calculated by FreeSurfer for each of these scans. In

565     addition, 1,000 Year 2 scans were semi-randomly selected for MQC ratings, such that they

566     contained (1) a range of scan quality, operationalized by selecting for an approximately

26

567   equivalent number of scans that fell into tiers A, B, C, and D; and (2) a distribution of magnet

568   types (Siemens, Philips, GE) that was equivalent to the analyzed Baseline sample.  These scans

569   were then rated for MQC in random sequence by two raters (E.L, K.A.K.) who had previously

570   been trained by the rater of all Baseline scans (S.E), such that the three raters achieved an

571   intraclass coefficient of >0.75 (two-way mixed effects model for absolute agreement) across a

572   training set of 1,000 Baseline scans.

573   *Manual cortical edits of ABCD scans*

574       A subset of the rated Baseline scans was randomly selected for manual editing (N=150 with

575   MQC=1, N=30  with MQC=2). Each structural scan was loaded into Freeview version 7.1.1 with

576   the following volumes: brainmask, wm, brain.finalsurfs.manedit, T1, and the following surfaces:

577   rh.pial, rh.white, lh.pial, lh.white. The scans were primarily displayed in the coronal view,

578   although sagittal and horizontal views were used as needed. Criteria for editing were primarily

579   based off overestimation and underestimation of the pial and white matter boundaries. Edits to

580   the white matter boundary were made directly on the wm volume using control points and the

581   erasing tool. Edits to the pial surface were made on the brainmask volume. Errors between the

582   pial surface and cerebellum were corrected using the brain.finalsurfs.manedit volume. Edits were

583   considered to be complete when, after post-edit re-processing in FreeSurfer, there appeared only

584   minimal errors remaining, meaning the generated pial and white matter boundaries more closely

585   matched the actual boundaries on the T1 image.

586   *Manual edits of Massachusetts General Hospital (MGH) scans*

587    The MGH sample was included as a replication set for effects of manual editing on cortical

588    MRI indices and to assess whether such effects change later in adolescence.  Study sample,

589    scanner characteristics, and editing methods were previously described by Eryilmaz and

590    colleagues [20].  Study procedures were approved by Partners Human Research Committee, which

591    granted a waiver of informed consent, since this retrospective study of the medical record

592    involved only deidentified data.  Briefly, clinical brain MRI scans from 292 individuals aged 8 to

593    17, conducted at MGH between 2005 and 2015, were selected based on date of birth, adequate

594    scan quality on visual inspection (i.e., artifacts could reasonably be addressed with manual edits),

595    and absence of pathology as indicated on radiology reports.  Scans were edited by a trained

596    research coordinator (KFD) as described above.  Pre-to-post edit changes in cortical thickness,

597    volume, and surface area were measured across 68 ROIs using FreeSurfer 5.0.

598    ***Statistical analysis***

599    *Stability of MQC ratings over time*

600    MQC ratings of baseline scans that did not show signal dropout or cysts were divided into 10

601    equally sized time groups, reflecting the sequence in which scans were evaluated.  Initial

602    analyses were conducted to assess whether factors known to affect scan quality, including age,

603    gender, scanner manufacturer, and psychopathology (CBCL) differed over time, using time

604    period as either a categorical or continuous variable.  Then, ANOVA was used to assess

605    significant linear or quadratic changes in mean MQC rating across time groups, controlling for

606    variation in these other factors and in their interactions with time and time-squared.

607    *Surface-based sMRI analyses*

608    Surface maps for group-based and within-subject analyses were generated using Freesurfer

609    7.0. Images from each subject were smoothed by 22mm full width-half maximum. For between-

610    group analyses we fit general linear models with following covariates: age, gender, estimated

611    intracranial volume, study site, and scanner.  Continuous predictor variables were z-transformed

612    prior to analysis.  Models assessed both linear effects of MQC ratings (i.e., 1 to 4) as well as

613    pairwise contrasts (1 vs. 2, 1 vs. 3, 1 vs. 4) on cortical thickness, surface area, and volume.

614    Sensitivity analyses assessed linear effects of SHN on these indices, as well as effects of MQC

615    after controlling for SHN and vice versa.  Results were visualized using uncorrected significance

616    maps (log p-value) and effect size maps (Cohen's d) as appropriate.

617    *ROI-based sMRI analyses*

618    Following extraction of ROI-based data from Freesurfer, analyses involving cortical

619    thickness, cortical surface area, and cortical and subcortical volumes were conducted with R

620    version 4.1.2 (https://www.R-project.org/). Mixed-effects linear regression was run with "lme4"

621    package (https://github.com/lme4/lme4/), unless specifically mentioned. The covariates included

622    in the analysis were age, gender, estimated intracranial volume (fixed effect), site, scanner, and

623    family ID (random effects), the latter accounting for inclusion of sibling groups.  Analyses were

624    corrected for multiple comparisons using FDR (q<.05), based on the number of included ROIs.

625    *SHN tiers*

626    We conducted receiver operating characteristic (ROC) analyses to evaluate the sensitivity of

627    SHN to detect poorer quality scans.  Analyses were conducted in R using the "pROC" package.

628    Using Baseline scan data, we contrasted SHN for three breakpoints: MQC=1 vs. 2, 3, and 4;

629   MQC=1 and 2 vs. 3 and 4; and MQC=1, 2, and 3 vs. 4. We used the Youden Index to select an

630   optimal threshold to discriminate higher versus lower quality scans for each of the three

631   breakpoints.  These three thresholds were used to define SHN tiers A, B, C, and D, respectively –

632   such that scans in the A tier best represented MQC=1, those in the B tier best represented

633   MQC=2, etc.  As a sensitivity analysis, we also included MQC and SHN values for scans with

634   segmentation-related tissue loss into the analysis, and examined whether thresholds were altered

635   by inclusion of these scans.  Then, to test reliability, we grouped all available Year 2 scans

636   according to SHN tiers, and conducted MQC on 1,000 of these scans (described above).

637   Sensitivity, specificity, and accuracy of the SHN tiers to distinguish MQC levels were assessed.

638   These metrics could then be compared to those from the Baseline analysis, as well as to those

639   from a new set of ROC analyses that determined optimal thresholds for SHN tiers in the 1,000

640   Year 2 scans.

641   *Applied analyses relating quality control to MRI-clinical associations*

642      Linear mixed models examined associations between cortical thickness and age, and between

643   cortical thickness and externalizing psychopathology, conditioned on the degree to which lower-

644   quality scans were included in the analyses (e.g., inclusion of MQC=1 only, versus MQC 1 and

645   2; 1, 2, and 3; and 1, 2, 3, and 4).  Overall surface-based and ROI methods were similar to those

646   described above, but now using age or CBCL externalizing score rather than MQC as the

647   predictors of interest.  Sensitivity analyses examined effects of including SHN as an additional

648   predictor in the models.

649   *Effects of manual edits on sMRI indices*

650     For Baseline ABCD scans, within-subject analyses that contrasted cortical thickness, surface

651     area, and volume before vs. after manual edits were conducted using general linear models in

652     Freesurfer (for surface maps of effect size) or paired t-tests in R (for ROI analyses).  These

653     analyses were conducted without covariates, following upon sensitivity analyses that

654     demonstrated no significant effects of age, gender, scanner, or CBCL externalizing symptoms on

655     pre-to-post edit changes in sMRI measures.  ROI analyses were corrected for multiple

656     comparisons using FDR (q<.05), based on the number of included ROIs.  Analyses of MGH

657     scans focused on cortical regions that replicated significant effects of manual edits on cortical

658     thickness that were seen in the ABCD cohort.  Potential changes in magnitude and variance of

659     pre-to-post edit changes across these regions were assessed as a function of age group (8-10, 11-

660     12, 13-14, 15-17 years) using ANOVA.

661     ***Data availability***

662     Data from all ABCD-related analyses were downloaded from the NIMH Data Archive

663     (NDA), version 4.0.  Derived variables, including MQC ratings and SHN, as well as region-of-

664     interest level data for cortical thickness, surface area, and volume processed in FreeSurfer 7.0,

665     have been uploaded to the NDA (Study ID #1944, doi 10.15154/1528507).  Data from MGH

666     analyses contain sensitive patient information that was obtained following a waiver of informed

667     consent, and as such has not been uploaded to a publicly available repository.  Please contact the

668     corresponding author for additional information.

669

## References

1. Becht, A. I. & Mills, K. L. Modeling Individual Differences in Brain Development. *Biol Psychiatry* **88**, 63–69 (2020).

2. Dick, A. S. *et al.* Meaningful Associations in the Adolescent Brain Cognitive Development Study. *Neuroimage* **239**, 118262 (2021).

3. Alexander-Bloch, A. *et al.* Subtle in-scanner motion biases automated measurement of brain anatomy from in vivo MRI. *Hum Brain Mapp* **37**, 2385–2397 (2016).

4. Blumenthal, J. D., Zijdenbos, A., Molloy, E. & Giedd, J. N. Motion artifact in magnetic resonance imaging: implications for automated analysis. *Neuroimage* **16**, 89–92 (2002).

5. Rosen, A. F. G. *et al.* Quantitative Assessment of Structural Image Quality. *Neuroimage* **169**, 407–418 (2018).

6. Karcher, N. R. & Barch, D. M. The ABCD study: understanding the development of risk for mental and physical health outcomes. *Neuropsychopharmacology* **46**, 131–142 (2021).

7. Thompson, P. M. *et al.* ENIGMA and global neuroscience: A decade of large-scale studies of the brain in health and disease across more than 40 countries. *Transl Psychiatry* **10**, 100 (2020).

8. Mills, K. L. & Tamnes, C. K. Methods and considerations for longitudinal structural brain imaging analysis across development. *Dev Cogn Neurosci* **9**, 172–190 (2014).

9. Marquand, A. F. *et al.* Conceptualizing mental disorders as deviations from normative functioning. *Mol Psychiatry* **24**, 1415–1424 (2019).

10. Reuter, M. *et al.* Head motion during MRI acquisition reduces gray matter volume and thickness estimates. *NeuroImage* **107**, 107–115 (2015).
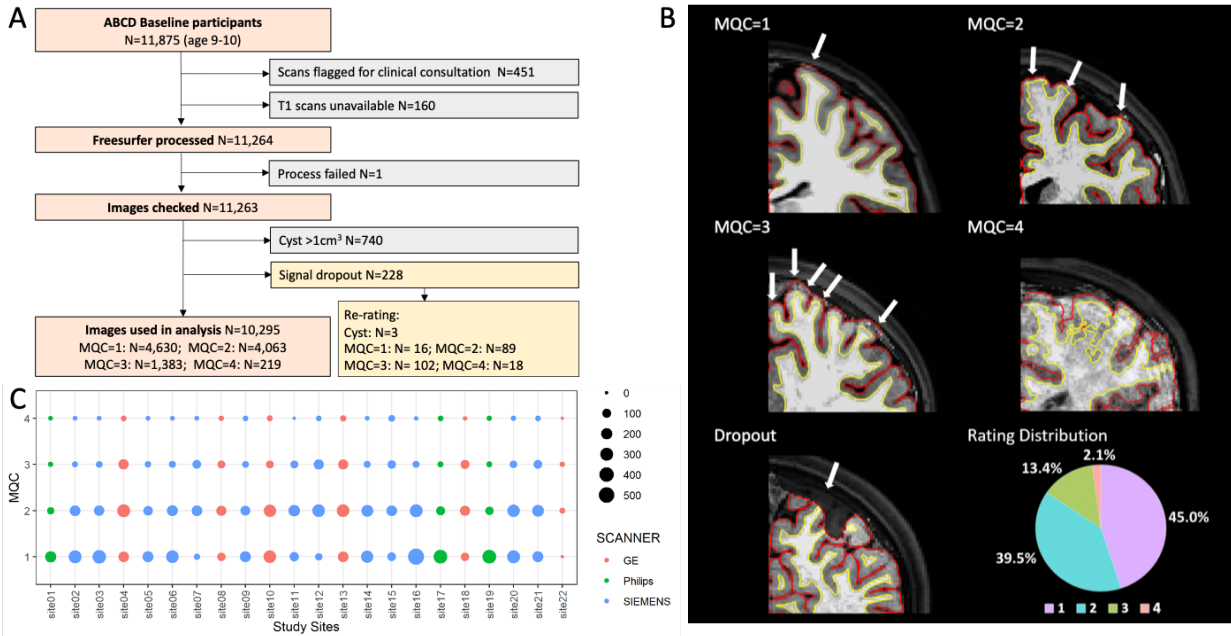
694    11.    White, N. *et al.* PROMO: Real-time prospective motion correction in MRI using image-

695        based tracking. *Magn Reson Med* **63**, 91–105 (2010).

696    12.    Tisdall, M. D. *et al.* Prospective motion correction with volumetric navigators (vNavs)

697        reduces the bias and variance in brain morphometry induced by subject motion. *Neuroimage*

698        **127**, 11–22 (2016).

699    13.    Dale, A. M., Fischl, B. & Sereno, M. I. Cortical surface-based analysis. I. Segmentation

700        and surface reconstruction. *Neuroimage* **9**, 179–194 (1999).

701    14.    White, T. *et al.* Automated quality assessment of structural magnetic resonance images in

702        children: Comparison with visual inspection and surface-based reconstruction. *Hum Brain*

703        *Mapp* **39**, 1218–1231 (2018).

704    15.    Waters, A. B., Mace, R. A., Sawyer, K. S. & Gansler, D. A. Identifying errors in

705        Freesurfer automated skull stripping and the incremental utility of manual intervention. *Brain*

706        *Imaging Behav* **13**, 1281–1291 (2019).

707    16.    Monereo-Sánchez, J. *et al.* Quality control strategies for brain MRI segmentation and

708        parcellation: Practical approaches and recommendations - insights from the Maastricht study.

709        *Neuroimage* **237**, 118174 (2021).

710    17.    Ross, M. C. *et al.* Gray matter volume correlates of adolescent posttraumatic stress

711        disorder: A comparison of manual intervention and automated segmentation in FreeSurfer.

712        *Psychiatry Res Neuroimaging* **313**, 111297 (2021).

713    18.    McCarthy, C. S. *et al.* A comparison of FreeSurfer-generated data with and without

714        manual intervention. *Front Neurosci* **9**, 379 (2015).

715    19.    Beelen, C., Phan, T. V., Wouters, J., Ghesquière, P. & Vandermosten, M. Investigating

716         the Added Value of FreeSurfer's Manual Editing Procedure for the Study of the Reading

717         Network in a Pediatric Population. *Front Hum Neurosci* **14**, 143 (2020).

718    20.    Eryilmaz, H. *et al.* Association of Prenatal Exposure to Population-Wide Folic Acid

719         Fortification With Altered Cerebral Cortex Maturation in Youths. *JAMA Psychiatry* **75**, 918–

720         928 (2018).

721    21.    Pulli, E. P. *et al.* Feasibility of FreeSurfer Processing for T1-Weighted Brain Images of 5-

722         Year-Olds: Semiautomated Protocol of FinnBrain Neuroimaging Lab. *Front Neurosci* **16**,

723         874062 (2022).

724    22.    Garavan, H. *et al.* Recruiting the ABCD sample: Design considerations and procedures.

725         *Dev Cogn Neurosci* **32**, 16–22 (2018).

726    23.    Casey, B. J. *et al.* The Adolescent Brain Cognitive Development (ABCD) study: Imaging

727         acquisition across 21 sites. *Dev Cogn Neurosci* **32**, 43–54 (2018).

728    24.    Ducharme, S. *et al.* Trajectories of cortical thickness maturation in normal brain

729         development – The importance of quality control procedures. *Neuroimage* **125**, 267–279

730         (2016).

731    25.    Wainberg, M., Jacobs, G. R., Voineskos, A. N. & Tripathy, S. J. Neurobiological,

732         familial and genetic risk factors for dimensional psychopathology in the Adolescent Brain

733         Cognitive Development study. *Mol Psychiatry* **27**, 2731–2741 (2022).

734    26.    Wu, X. *et al.* Symptom-Based Profiling and Multimodal Neuroimaging of a Large

735         Preteenage Population Identifies Distinct Obsessive-Compulsive Disorder-like Subtypes With

736         Neurocognitive Differences. *Biol Psychiatry Cogn Neurosci Neuroimaging* **7**, 1078–1089

737         (2022).

738    27.    Hughes, D. *et al.* Genetic Patterning for Child Psychopathology is Distinct from Adults

739        and Implicates Fetal Cerebellar Development. *Nature Neuroscience*. In Press.

740     28.    Marek, S. *et al.* Reproducible brain-wide association studies require thousands of

741        individuals. *Nature* **603**, 654–660 (2022).

742    29.    Szucs, D. & Ioannidis, J. P. Sample size evolution in neuroimaging research: An

743        evaluation of highly-cited studies (1990-2012) and of latest practices (2017-2018) in high-

744        impact journals. *Neuroimage* **221**, 117164 (2020).

745    30.    Open Science Collaboration. PSYCHOLOGY. Estimating the reproducibility of

746        psychological science. *Science* **349**, aac4716 (2015).

747    31.    Kaplan, R. M., Chambers, D. A. & Glasgow, R. E. Big Data and Large Sample Size: A

748        Cautionary Note on the Potential for Bias. *Clin Transl Sci* **7**, 342–346 (2014).

749    32.    Smaczny, S. *et al.* Disconnection in a left-hemispheric temporo-parietal network impairs

750        multiplication fact retrieval. *Neuroimage* **268**, 119840 (2023).

751    33.    Sonuga-Barke, E. J. S. Editorial: 'Safety in numbers'? Big data discovery strategies in

752        neuro-developmental science - contributions and caveats. *J Child Psychol Psychiatry* **64**, 1–3

753        (2023).

754    34.    Pardoe, H. R., Kucharsky Hiess, R. & Kuzniecky, R. Motion and morphometry in clinical

755        and nonclinical populations. *NeuroImage* **135**, 177–185 (2016).

756    35.    Smith, J. *et al.* Can this data be saved? Techniques for high motion in resting state scans

757        of first grade children. *Dev Cogn Neurosci* **58**, 101178 (2022).

758    36.    Saccà, V. *et al.* Aging effect on head motion: A Machine Learning study on resting state

759        fMRI data. *J Neurosci Methods* **352**, 109084 (2021).

760    37.    Backhausen, L. L., Herting, M. M., Tamnes, C. K. & Vetter, N. C. Best Practices in

761          Structural Neuroimaging of Neurodevelopmental Disorders. *Neuropsychol Rev* **32**, 400–418

762          (2022).

763    38.    Duffy, B. A. *et al.* Retrospective motion artifact correction of structural MRI images

764          using deep learning improves the quality of cortical surface reconstructions. *Neuroimage* **230**,

765          117756 (2021).

766    39.    Auchter, A. M. *et al.* A description of the ABCD organizational structure and

767          communication framework. *Dev Cogn Neurosci* **32**, 8–15 (2018).

768    40.    Tustison, N. J. *et al.* N4ITK: Improved N3 Bias Correction. *IEEE Trans Med Imaging* **29**,

769          1310–1320 (2010).

770    41.    Achenbach, T. M. The Achenbach System of Empirically Based Assessment (ASEBA):

771          Development. *Findings, Theory, and Applications* (2009).

772
773

**Figure 1. Manual quality control (MQC) protocol.** (A) Among 11,875 total participants at Baseline, we excluded participants with clinical findings (see Methods), broken or blank T1 images, or repeated failed FreeSurfer preprocessing. After excluding additional images found to have cysts or signal dropout, we rated 10,295 images on MQC=1-4 scale (1: best, 4: worst). (B) Distribution of MQC ratings, stratified by site and scanner. (C) Representative examples of MQC=1-4 scans and a scan with apparent tissue loss due to segmentation error. Arrows indicate areas where manual edits are needed to correct for errant automated segmentation of the pial surface from the underlying cortex. Distribution of MQC ratings among all scans is displayed at lower right.

786



**Figure 2. Association between MQC ratings and sMRI indices, n=10,261**. Maps at left show linear associations of MQC rating (1 to 4) with cortical thickness (A), surface area (B), and volume (C). Maps at right contrast thickness, surface area, and volume highest quality images (MQC=1) with those assigned to lower quality ratings. Covariates included age, gender, estimated intracranial volume (fixed effects), site, and scanner manufacturer (random effects). Of the initial 10,295 scans with MQC, 34 were excluded due to missing covariates or FreeSurfer processing errors.

**Figure 3. Effects of surface hole number (SHN) on sMRI indices, and derivation of SHN tiers in conjunction with MQC ratings, n=10,261.** Linear associations of SHN (non-transformed) with cortical thickness (A), surface area (B), and volume (C) closely resembled those of MQC ratings (compare to Figure 2). Covariates included age, gender, estimated intracranial volume (fixed effects), site, and scanner manufacturer (random effects). Additional adjustment for SHN diminished the effect sizes of pairwise MQC contrasts for thickness (D), surface area (E), and volume (F). Markers represent effect sizes for pairwise MQC contrasts in each of 68 cortical regions-of-interest, and solid lines reflect best-fit across all 68 regions for a given pairwise MQC contrast. Note reduced slopes compared to dashed unity line. (G) Density plot of SHN values, stratified by MQC ratings. Panel (H) illustrates the overall approach for deriving SHN tiers from MQC ratings. The SHN tiers were developed to provide quality control estimates in the absence of manual ratings, and are based on optimized SHN thresholds for parsing higher versus lower manual quality scan groupings. Receiver-operating characteristic (ROC) analyses for various thresholds are shown in panels (I), (J), and (K) along with related specificity, sensitivity, and accuracy indices. For example, with an optimized SHN threshold of 29.5, 81.3% of scans with MQC=2 and higher are eliminated. This threshold was used as a breakpoint for SHN tiers A and B. Blue shaded regions indicate 95% confidence intervals. AUC: area under the curve.
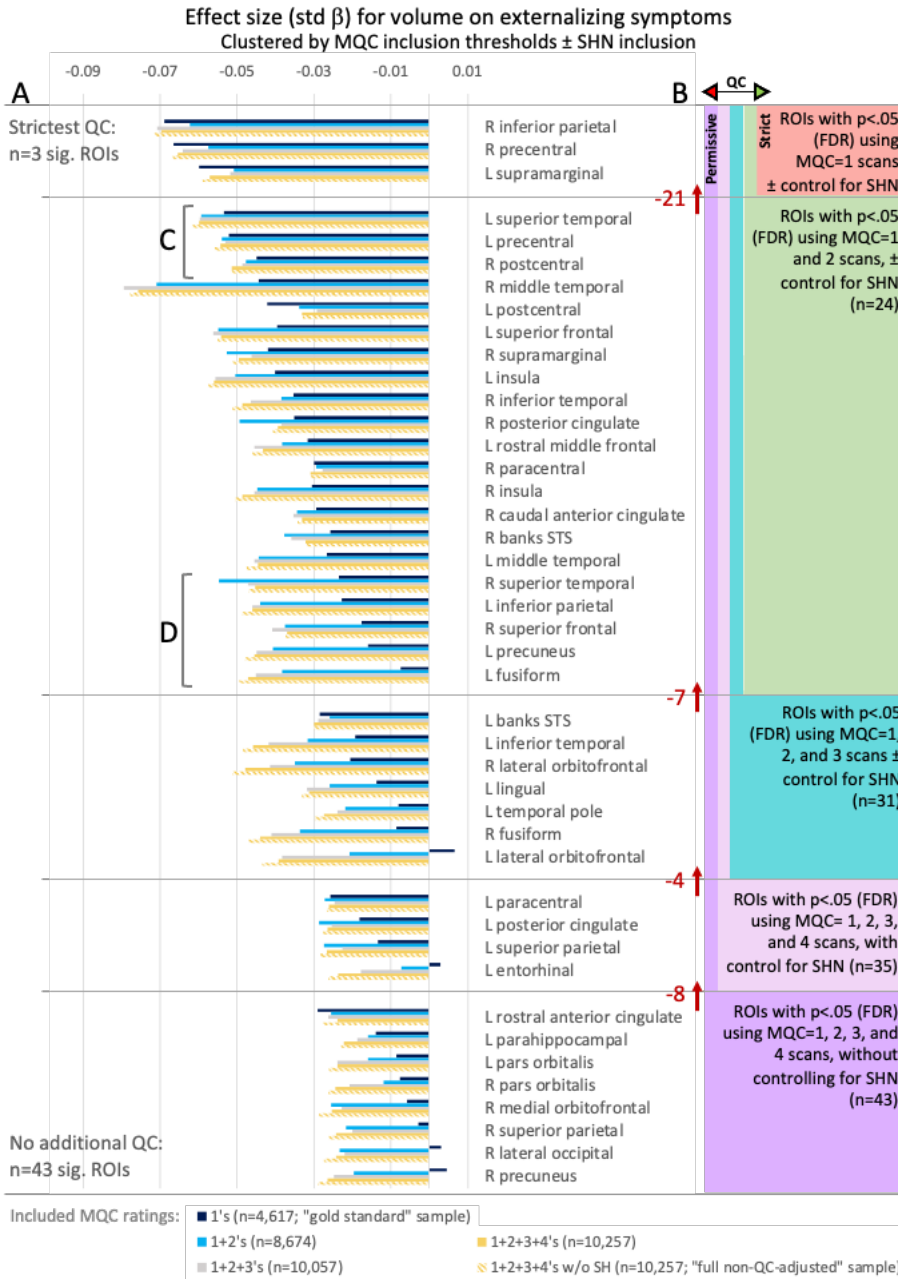
39

**Figure 4. Effects of variable quality control on applied analyses of sMRI data.** (A) Association of age with cortical thickness, without adjusting for manual quality control (MQC) rating or surface hole number (SHN). Note the substantially smaller effect size scale compared to Figure 2, which shows effects of quality control variance on sMRI measurements. (B) Association of externalizing symptoms (CBCL externalizing subscale) with cortical volume, without adjusting for MQC or SHN. Note the even smaller effect size compared to effects of age on thickness. (C) Age-thickness effects stratified by region of interest (ROI) and MQC inclusion threshold. Broken lines indicate best-fit lines across all ROIs for each inclusion threshold group. Note tendency toward *diminished* effect size (and increased risk for false negatives) with broader inclusion thresholds. (D) Externalizing symptoms-volume effects stratified by ROI and MQC inclusion threshold. Broken lines indicate best-fit lines for each inclusion threshold group. Note tendency toward *inflated* effect size (and increased risk for false positives) with broader inclusion thresholds. All analyses covaried for age, gender, estimated intracranial volume (fixed effects), site, and scanner manufacturer (random effects); ROI-based analyses also included family ID as a random effect.
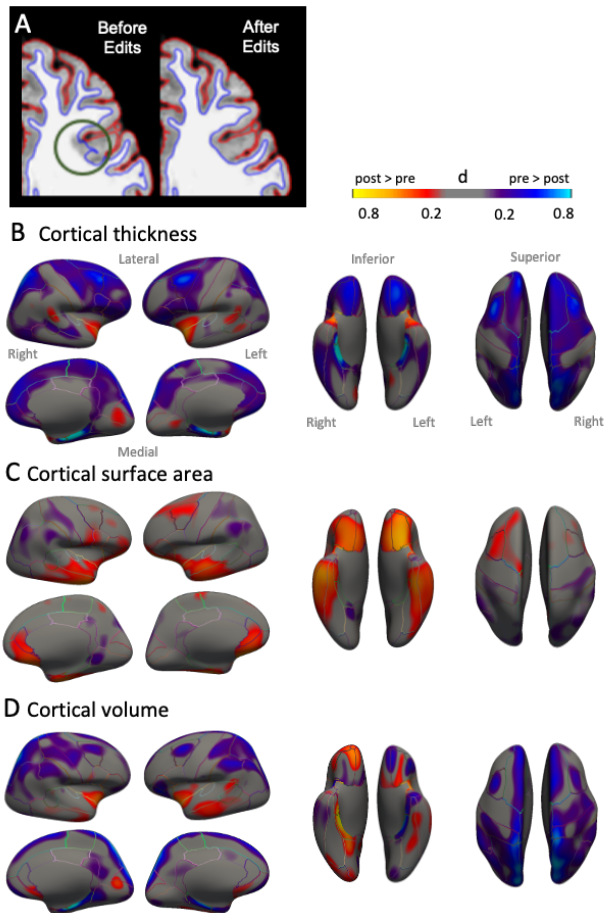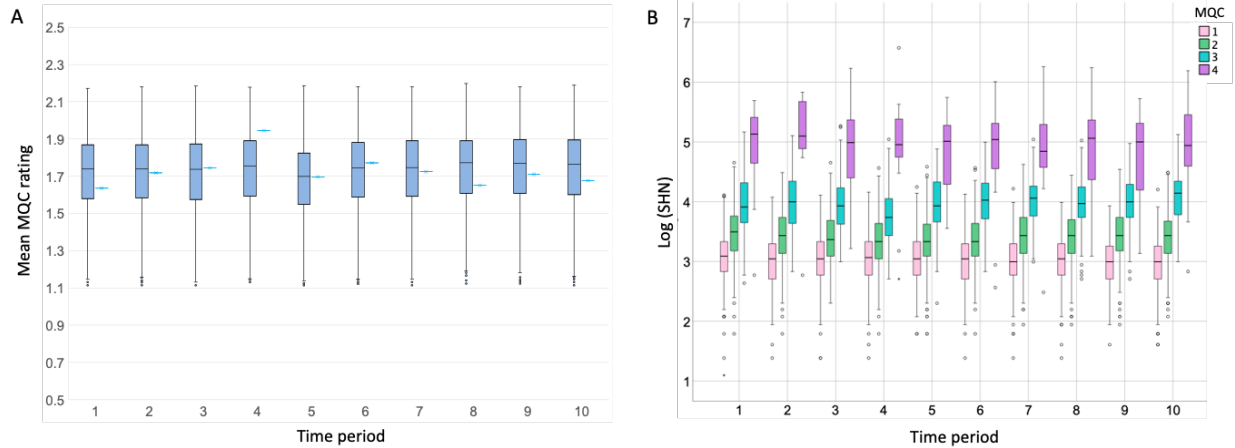
40

**Figure 5. Effects of increasingly stringent quality control on effect size and statistical significance of externalizing symptoms-volume findings.** (A) At left, bars indicate effect sizes for the relationships between externalizing symptoms and cortical volume for each ROI, stratified by the stringency of quality control of included scans (see legend). At one extreme, dark blue bars indicate effect sizes generated by using the most conservative approach, i.e., only MQC=1 scans were included in the analysis, which also corrected for SHN ("gold standard" sample, n=4,617). At the other extreme, thatched yellow bars indicate effect sizes generated by using the most permissive approach, i.e., scans with all MQC levels were included in the analysis, and no SHN correction was applied ("full non-QC-adjusted" sample, n=10,257). Note that for most regions, more permissive quality control was associated with inflated effect sizes.

850    (B) As seen at right, ROIs were grouped based on whether they continued to show statistically

851    significant (q<.05, FDR) relationships between externalizing symptoms and cortical volumes as

852    lower quality scans were iteratively removed.  Red numbers and arrows indicate the number of

853    ROIs that dropped out of significance with each level of tightened QC.  The purple group

854    contains 43 regions that were significant after using the most permissive quality control (no

855    removed scans).  In contrast, the red group (gold standard) contains only 3 regions that were

856    significant after using the most conservative quality control (MQC=2, 3, and 4 removed).  Note

857    that when including the next best of scans (MQC=2), several regions that become significant in

858    this larger sample (n=8,674), e.g., those near (C), did not show inflated effect sizes when lower

859    quality scans were included, and thus appeared robust to poor scan quality.  That these regions

860    are significant when MQC 1+2 scans are included in the analysis – but *not* significant when only

861    MQC=1 scans are included (n=4,617) – indicates that using only the highest quality scans

862    potentially results in false negatives (type II error) due to lack of statistical power.  However,

863    other regions that were significant in the MQC 1+2 group, e.g., those near (D), showed

864    substantial effect size inflations when scans rated as MQC=2 or higher were included.  For these

865    regions, statistically significant findings likely reflected false positives (type I error) – even when

866    all included scans were of relatively good quality.  All analyses covaried for age, gender,

867    estimated intracranial volume (fixed effects), and site, scanner manufacturer, and family ID
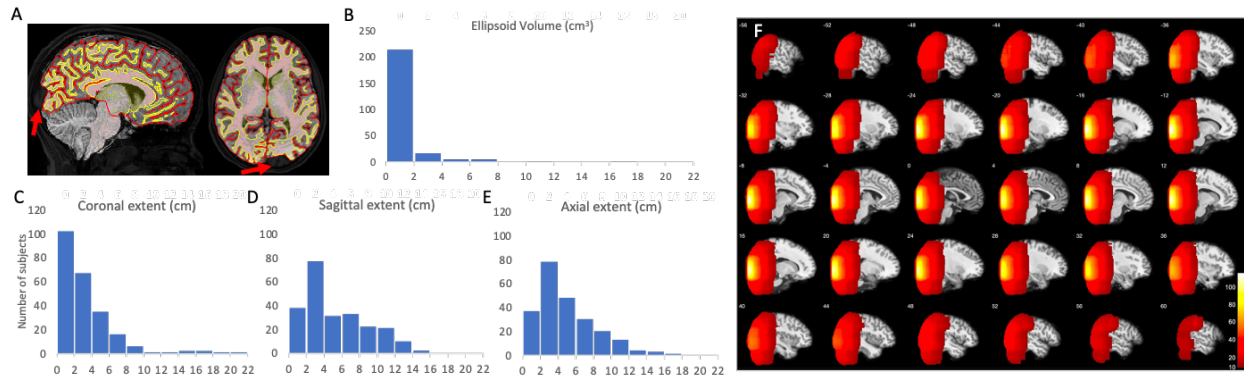
868    (random effects).

869

870
871
872 **Figure 6. Effects of manual edits on sMRI indices (n=180).** Manual edits (e.g., A, which
873 corrects a gray-white matter boundary segmentation error) were conducted on 150 scans with
874 MQC=1 and 30 scans with MQC=2.  Maps reflect effect sizes (Cohen's d) of pre-to-post edit
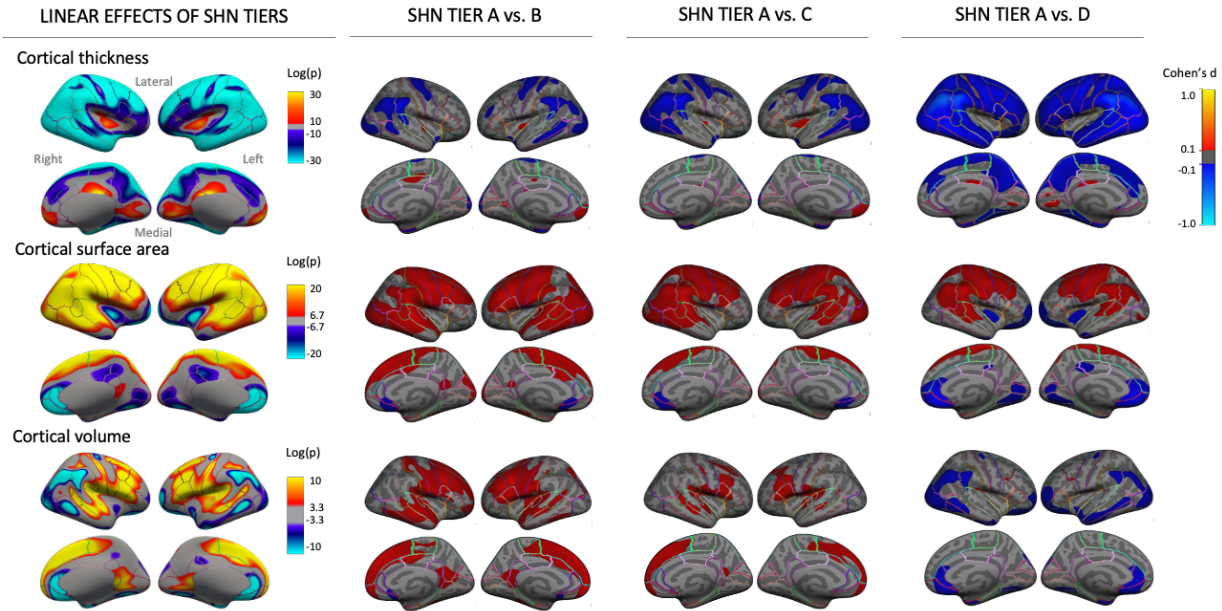875 changes in (B) cortical thickness, (C) cortical surface area, and (D) cortical volume.
876

**Extended Data Figure 1. Stability of manual quality control (MQC) ratings over time (n=10,295).** Scans were assigned to deciles based on the sequence in which they received MQC ratings by a single trained rater. (A) Box and whisker plots show distribution of MQC ratings for each time period, after adjusting for age, gender, scanner manufacturer, and externalizing psychopathology. Adjacent marks show unadjusted mean ratings for the same period. (B) Box and whisker plots show distribution of the log of surface hole numbers (SHN), stratified by decile and MQC rating.
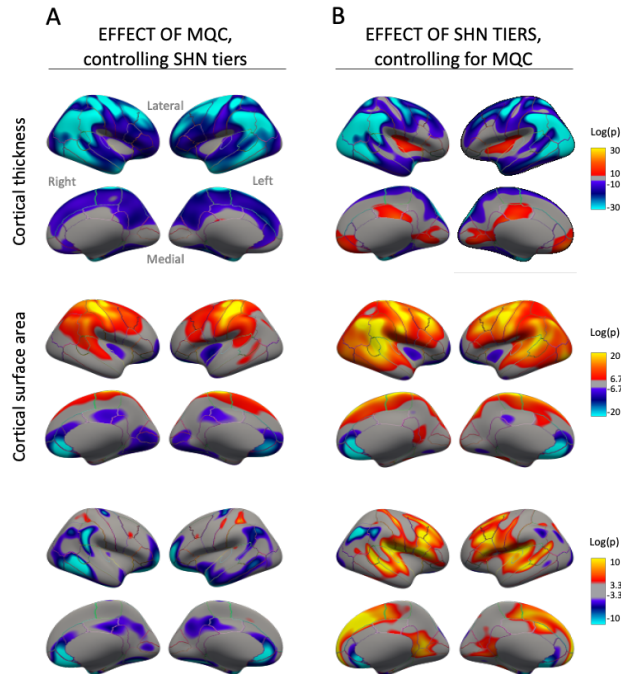
**Extended Data Figure 2. Signal dropout in sMRI processing (n=228).** (A) Examples of dropout regions where FreeSurfer segmentation failed and did not include a substantial portion of cortex. (B) Distribution of approximate volume of dropout area estimated by ellipsoid volume calculated and distribution of (C) sagittal, (D) coronal, and (E) axial extent. (F) Distributions of drop-out regions overlaid on exemplar brain thresholded at n=10 subjects. Heat map represents number of overlapping subjects.
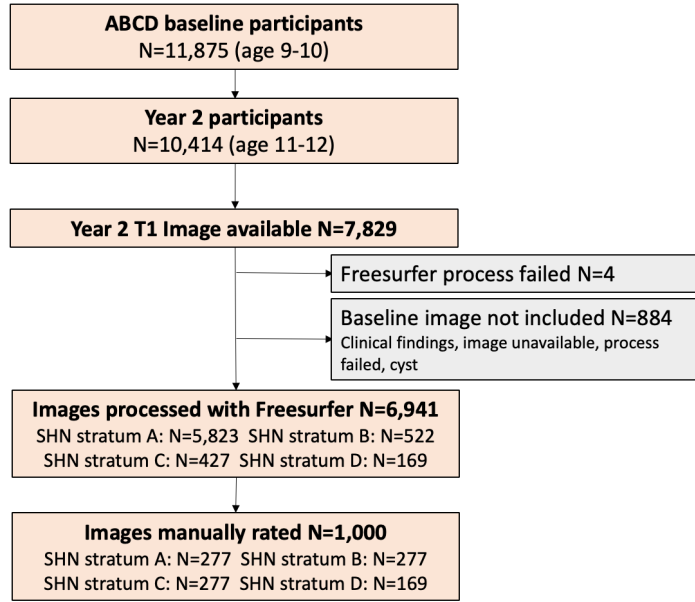
**Extended Figure 3. Comparison of SHN tier effects on sMRI indices at Baseline, n=10,295; compare to Figure 2.** Maps at left show linear associations of SHN tier (A to D) with cortical thickness, surface area, and volume. Maps at right contrast thickness, surface area, and volume highest quality images (SHN=A) with those assigned to lower quality ratings. Covariates included age, gender, estimated intracranial volume (fixed effects), site, and scanner manufacturer (random effects).

**Extended Data Figure 4. Unique contributions of SHN tiers versus MQC to variance in sMRI indices, n=10,295.** (A) Linear association of MQC on cortical indices after controlling for SHN tiers. (B) Linear association of SHN tiers on cortical indices after controlling for MQC. Covariates included age, gender, estimated intracranial volume (fixed effects), site, and scanner manufacturer (random effects).
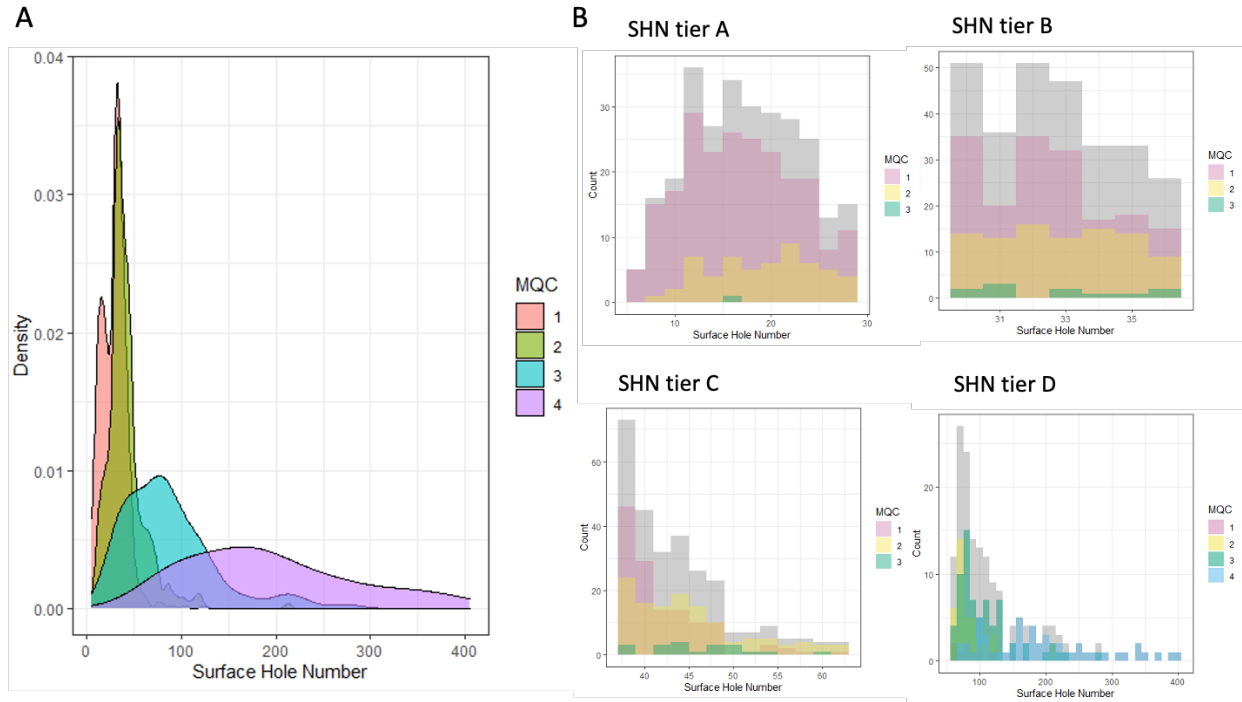
```
┌─────────────────────────────────────┐
│      ABCD baseline participants      │
│         N=11,875 (age 9-10)          │
└─────────────────────────────────────┘
                 ↓
┌─────────────────────────────────────┐
│          Year 2 participants         │
│         N=10,414 (age 11-12)         │
└─────────────────────────────────────┘
                 ↓
┌─────────────────────────────────────┐
│   Year 2 T1 Image available N=7,829  │
└─────────────────────────────────────┘
         │      ┌──────────────────────────────────────┐
         │─────▶│    Freesurfer process failed N=4     │
         │      └──────────────────────────────────────┘
         │      ┌──────────────────────────────────────┐
         │      │   Baseline image not included N=884  │
         │─────▶│ Clinical findings, image unavailable,│
         │      │ process failed, cyst                 │
         │      └──────────────────────────────────────┘
         ↓
┌─────────────────────────────────────────────────┐
│   Images processed with Freesurfer N=6,941      │
│  SHN stratum A: N=5,823  SHN stratum B: N=522   │
│  SHN stratum C: N=427  SHN stratum D: N=169     │
└─────────────────────────────────────────────────┘
                 ↓
┌─────────────────────────────────────────────────┐
│        Images manually rated N=1,000            │
│  SHN stratum A: N=277  SHN stratum B: N=277     │
│  SHN stratum C: N=277  SHN stratum D: N=169     │
└─────────────────────────────────────────────────┘
```
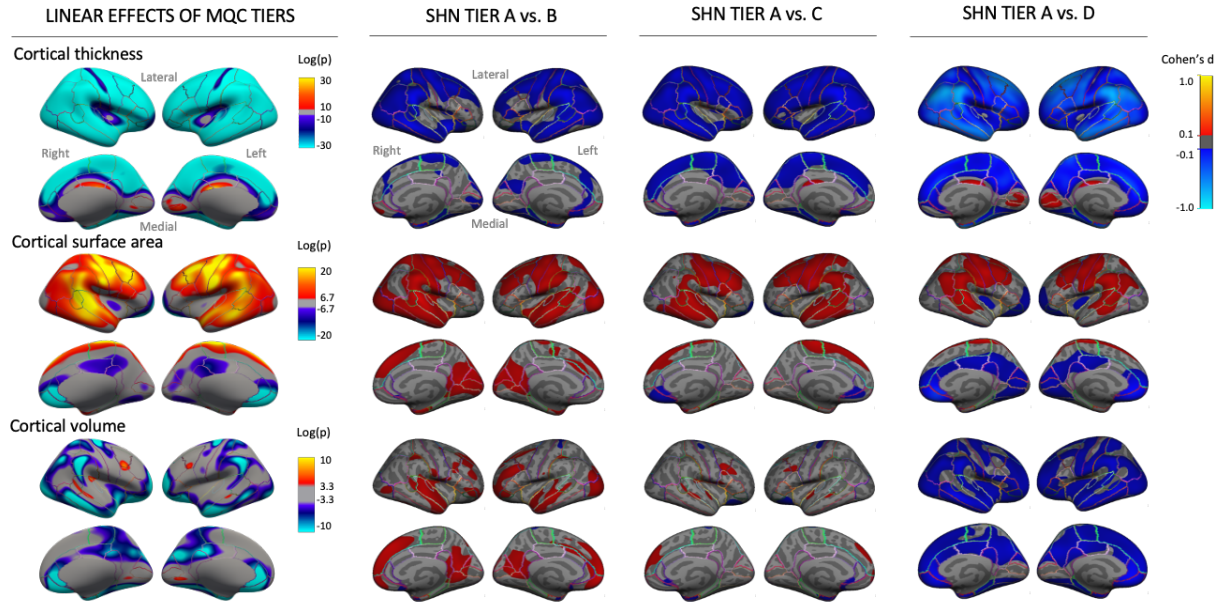
**Extended Data Figure 5. Included Year 2 follow-up scans.** Among 11,875 total participants at baseline, Year 2 T1 scans were available from 7,829; of these, 6,941 were eligible for processing with FreeSurfer, and 1,000 were semi-randomly selected for MQC ratings (see Methods for additional details).

48

**Extended Data Figure 6. Relationship of surface hole number (SHN) to manual quality control (MQC) in selected Year 2 follow-up scans (n=999).** (A) Density plot of SHN values, stratified by MQC ratings. (B) Distribution of MQC ratings as related to SHN for each SHN tier.

**Extended Figure 7. SHN tier effects on sMRI indices at Year 2, n=6,941 (compare to Extended Data Figure 3).** Maps at left show linear associations of SHN tier (A to D) with cortical thickness, surface area, and volume. Maps at right contrast thickness, surface area, and volume highest quality images (SHN=A) with those assigned to lower quality ratings. Covariates included age, gender, estimated intracranial volume (fixed effects), site, and scanner manufacturer (random effects).
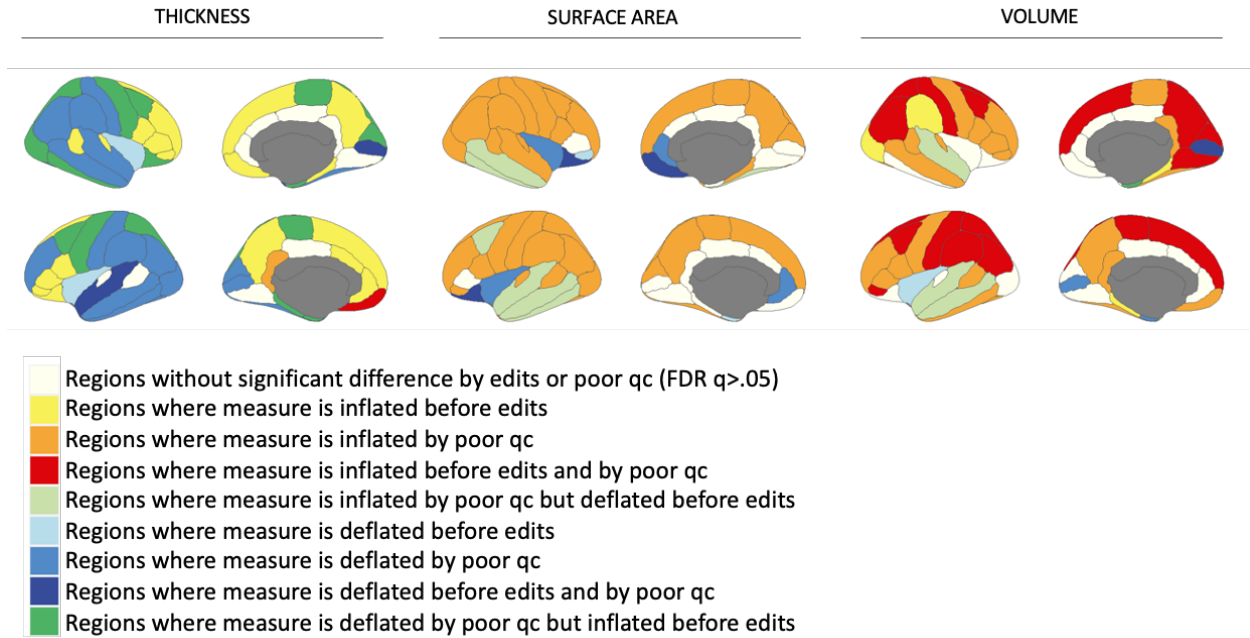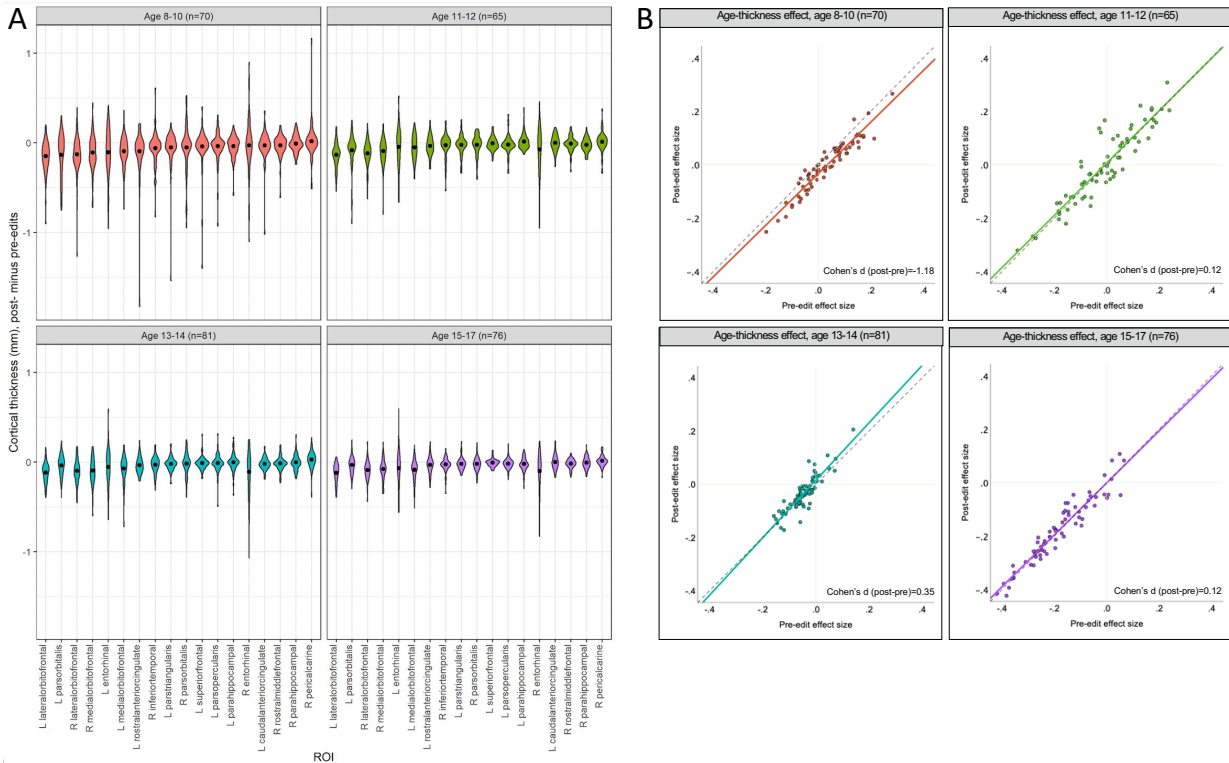
936



937
938
939 **Extended Data Figure 8.   Effects of manual edits on sMRI indices, stratified by MQC**
940 **rating**.  Edits were conducted on 150 scans with MQC=1 and 30 scans with MQC=2.  Maps
941 reflect effect sizes of pre-to-post edit changes in (A) cortical thickness, (B) cortical surface area,
942 and (C) cortical volume.  Note increased effects of edits in MQC=2 relative to MQC=1. (D)
943 Post-edit thickness reduction along the superior sagittal sinus, which is frequently misattributed
944 to pial surface during preprocessing.
945

**Extended Data Figure 9.  Composite maps showing location and direction of sMRI measurement errors detected by manual quality control and cortical edits, among MQC=1 and 2 scans only.**  Highlighted regions show either significant differences in sMRI indices between MQC=1 and MQC=2 scans, significant effects of cortical edits, or both.  Note that, when co-occurring within the same region, errors due to poor scan quality (assessed by MQC) do not necessarily occur in the same direction as errors requiring manual edits.

**Extended Data Figure 10. Effects of manual edits on cortical thickness and age-thickness relationships MGH sample, stratified by age group (n=292).** (A) Violin plots show effect size and related variance of manual edits on cortical thickness in the MGH sample, stratified by age group. The 18 included ROIs are those that also showed significant effects of edits on cortical thickness in the ABCD cohort, in the same direction. Regions are ordered by effect size in the 8- to 10-year-old group. Means are represented by black circles. Note that effect sizes and variance diminished with age. (B) Effects of edits on the magnitude of age-thickness relationships within the MGH sample across 68 cortical ROIs, stratified by age group. Each marker shows the age-thickness effect size for a given ROI. Edits strengthened age-thickness effects (i.e., effect sizes became more negative, indicated by lower intercept of the best-fit line compared to the dashed unity line) at age 8-10, but not in other age groups.

970 **<u>Acknowledgments</u>**

971      The authors are grateful to Drs. Randy L. Buckner and Erin C. Dunn for helpful comments on

972    the manuscript.

973

974 **Author contributions**

975 **Conception and experimental design:** Kunitoki, Clauss, Doyle, Lee, Tervo-Clemmens,

976 Eryilmaz, Satterthwaite, Roffman.

977 **Data acquisition:** Hopkinson, Eryilmaz, Gollub, Dowling, Roffman.

978 **Data analysis:** Elyounssi, Kunitoki, Clauss, Laurent, Kane, Hughes, Bazer, Sussman, Lee,

979 Dowling, Roffman.

980 **Data interpretation:** Elyounssi, Kunitoki, Class, Laurent, Kane, Hughes, Bazer, Doyle, Lee,

981 Tervo-Clemmens, Gollub, Barch, Satterthwaite, Dowling, Roffman.

982 **Drafting and revision of manuscript:** All authors.

983 All authors have approved the submitted version of the manuscript and have agreed to be

984 personally accountable to their own contributions.

985