


RESEARCH ARTICLE OPEN ACCESS

A Site-Wise Reliability Analysis of the ABCD Diffusion Fractional Anisotropy and Cortical Thickness: Impact of Scanner Platforms

Yezhi Pan^{1,2}  | L. Elliot Hong³ | Ashley Acheson⁴ | Paul M. Thompson⁵ | Neda Jahanshad⁵ | Alyssa H. Zhu⁵ | Jiaao Yu⁶ | Chixiang Chen^{2,7} | Tianzhou Ma⁸ | Ho-Ling Liu⁹ | Jelle Veraart¹⁰ | Els Fieremans¹⁰ | Nicole R. Karcher¹¹ | Peter Kochunov³ | Shuo Chen^{1,2} 

¹Maryland Psychiatric Research Center, Department of Psychiatry, School of Medicine, University of Maryland, Baltimore, Maryland, USA | ²Institute for Health Computing, University of Maryland, North Bethesda, Maryland, USA | ³Department of Psychiatry and Behavioral Science, University of Texas Health Science Center, Houston, Texas, USA | ⁴Department of Psychiatry, University of Arkansas for Medical Sciences, Little Rock, Arkansas, USA | ⁵Imaging Genetics Center, Mark & Mary Stevens Institute for Neuroimaging & Informatics, Keck School of Medicine, University of Southern California, Los Angeles, California, USA | ⁶Department of Mathematics, University of Maryland, College Park, Maryland, USA | ⁷Division of Biostatistics and Bioinformatics, Department of Epidemiology and Public Health, School of Medicine, University of Maryland, Baltimore, Maryland, USA | ⁸Department of Epidemiology and Biostatistics, School of Public Health, University of Maryland, College Park, Maryland, USA | ⁹Department of Imaging Physics, Division of Diagnostic Imaging, The University of Texas MD Anderson Cancer Center, Houston, Texas, USA | ¹⁰Center for Biomedical Imaging, Department of Radiology, New York University School of Medicine, New York, New York, USA | ¹¹Department of Psychiatry, Washington University School of Medicine, St. Louis, Missouri, USA

Correspondence: Shuo Chen (shuochen@som.umaryland.edu)

Received: 9 June 2024 | **Revised:** 9 October 2024 | **Accepted:** 25 October 2024

Funding: The project was funded by the National Institute on Drug Abuse of the National Institutes of Health, Award Number 1DP1DA048968-01, and by the following NIH grants: R01 EB015611, R01 MH094520, R01 MH096263, R01 AA012207, and P50 MH103222.

Keywords: brain development | diffusion tensor imaging | longitudinal | quality control | structural MRI | test–retest reliability

ABSTRACT

The Adolescent Brain and Cognitive Development (ABCD) project is the largest study of adolescent brain development. ABCD longitudinally tracks 11,868 participants aged 9–10 years from 21 sites using standardized protocols for multi-site MRI data collection and analysis. While the multi-site and multi-scanner study design enhances the robustness and generalizability of analysis results, it may also introduce nonbiological variances including scanner-related variations, subject motion, and deviations from protocols. ABCD imaging data were collected biennially within a period of ongoing maturation in cortical thickness and integrity of cerebral white matter. These changes can bias the classical test–retest methodologies, such as intraclass correlation coefficients (ICC). We developed a site-wise adaptive ICC (AICC) to evaluate the reliability of imaging-derived phenotypes while accounting for ongoing brain development. AICC iteratively estimates the population-level age-related brain development trajectory using a weighted mixed model and updates age-corrected site-wise reliability until convergence. We evaluated the test–retest reliability of regional fractional anisotropy (FA) measures from diffusion tensor imaging and cortical thickness (CT) from structural MRI data for each site. The mean AICC for 20 FA tracts across sites was 0.61 ± 0.19 , lower than the mean AICC for CT in 34 regions across sites, 0.76 ± 0.12 . Remarkably, sites using Siemens scanners consistently showed significantly higher AICC values compared with those using GE/Philips scanners for both FA (AICC = 0.71 ± 0.12 vs. 0.46 ± 0.17 , $p < 0.001$) and CT (AICC = 0.80 ± 0.10 vs. 0.69 ± 0.11 , $p < 0.001$). These findings

Peter Kochunov and Shuo Chen share the senior authorship.

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial](https://creativecommons.org/licenses/by-nc/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2024 The Author(s). *Human Brain Mapping* published by Wiley Periodicals LLC.

demonstrate site-and-scanner related variations in data quality and underscore the necessity for meticulous data curation in subsequent association analyses.

1 | Introduction

Adolescence is a crucial period for brain development that is associated with myelination of cerebral white matter (WM) tracts and pruning of cortical gray matter that support development of higher cognitive functions (Gogtay et al. 2004; Gogtay and Thompson 2009; Bartzokis et al. 2010; Kochunov et al. 2012; Kochunov et al. 2015). Adolescence is also associated with the onset of symptoms for severe mental illnesses (SMI), such as schizophrenia, bipolar, or major depressive disorders, substance use and others (Rapoport, Addington, and Frangou 2005; Casey, Nigg, and Durston 2007; Kalia 2008). Adolescent Brain and Cognitive Development (ABCD) is the largest longitudinal study of brain development and child health consisting of $N=11,868$ participants aged 9–10 years at baseline, ascertained at 21 sites across the US (Karcher and Barch 2021). The ABCD collection neuroimaging approaches were developed to collect data for quantitative longitudinal analysis of multi-site diffusion, structural and functional maturational changes (Casey et al. 2018; Hagler et al. 2019). This included standardized imaging protocols and preprocessing pipelines designed for multi-site homogenization and phenotype extraction (Casey et al. 2018; Hagler et al. 2019). However, nonbiological variations were reported in ABCD data due to differences in scanners, deviations from the protocol, imaging artifacts, and participant motion (Nielsen et al. 2018). Manual quality control (MQC) of the ABCD T1-weighted images suggested that up to 50% of the scans were affected with nonbiological variance (Elyounssi et al. 2023). Here, we performed quality assessment of longitudinal regional measurements of fractional anisotropy (FA) of water diffusion extracted for major WM tracts using the ABCD recommended pipeline. We specifically evaluated site-related differences in longitudinal fidelity of the FA values, including the differences in data collected using 3T scanners manufactured by Siemens, Philips, and General Electrics (GE). We developed an adaptive intraclass correlation coefficient (AICC) measure to evaluate the test–retest reliability of imaging-derived phenotypes while accounting for brain developmental trends in the population.

We focused on evaluating the impact of nonbiological variance on the longitudinal DTI-FA measurements of cerebral WM. FA is a sensitive biomarker for noninvasive studies of WM development (Basser 1994; Ulug, Barker, and van Zijl 1995; Conturo et al. 1996; Pierpaoli and Basser 1996). Although FA values are sensitive to many parameters (Beaulieu 2002), longitudinal changes in regional FA values during normal maturation are primarily attributed to myelination (Song et al. 2003; Song et al. 2005; Budde et al. 2007; Madler et al. 2008; Ryan et al. 2017; Ryan et al. 2018). Regional changes in cerebral FA values were used to replicate classical findings by Flechsig, who demonstrated that continued myelination of WM during adolescence and early adulthood underpins the development of higher cognitive function (Flechsig 1901; Kochunov et al. 2012). Herein, we evaluated the ability to detect longitudinal changes in FA, compared with longitudinal changes in cortical gray matter

thickness (CT) and speculated on potential causes of nonbiological variance.

There is no single established metric for evaluating test–retest reproducibility of neuroimaging measurements in longitudinal developmental studies. Many studies used intraclass correlation coefficients (ICC) to demonstrate reproducibility for metrics such as FA, CT, and other measurements (Shrout and Fleiss 1979; Wijtenburg et al. 2013; Zuo and Xing 2014; Acheson et al. 2017; Xue et al. 2021). However, ICC and other approaches may not be applicable to the ABCD data because the neuroimaging measures are collected biennially at the time when participants are undergoing rapid development (Kochunov, Glahn, Nichols et al. 2011; Kochunov, Glahn, Lancaster et al. 2011; Kochunov et al. 2012). ICC is performed with the assumption of repeated measures performed under similar conditions, and the dynamic alterations in the developing adolescent brain will bias the reliability measures (Barnea-Goraly et al. 2005; Barnhart, Haber, and Lin 2007; Casey, Jones, and Hare 2008; Konrad, Firk, and Uhlhaas 2013). In the present work, we describe an AICC measure to evaluate the reliability of imaging-derived phenotypes acquired from the ABCD study while accounting for the normative age-related changes (Figure 1). We first estimate developmental trajectory using the complete dataset, for it is more robust and accurate than site-wise estimation. The iterative AICC estimation process also factors site-wise data reliability into the calculation of age effects. The resulting site-wise AICC can be integrated into subsequent statistical analyses to reduce bias and enhance inference efficiency. A simulation study was also conducted to assess the accuracy and robustness of our method.

2 | Materials and Methods

2.1 | Study Samples

This study used baseline, two-year, and four-year follow up data from the NIMH Data Archive ABCD Curated Data Release 5.0 (<https://abcdstudy.org/>). The cohort and study protocols can be found in Garavan et al. (2018). Overall, the ABCD release 5.0 included early longitudinal data on 11,868 demographically diverse subjects, including neuroimaging data and other phenotypic data. For inclusion in the analyses, subjects were required to have both imaging data and relevant imaging acquisition information available. Additionally, subjects meeting any of the following exclusion criteria were excluded for the evaluation of site-wise data reliability: (1) attendance at different study sites during follow-up visits; (2) absence of longitudinal information.

2.2 | ABCD Image Acquisition

In the ABCD study, imaging data were acquired using Siemens (Prisma VE11B-C), Philips (Achieva dStream, Ingenia), and GE (MR750, DV25-26) 3-Tesla MRI scanners. Siemens scanners were equipped with either 32 or 64 channel head coils. Philips scanners used 32 channel head coil. GE protocol required the use of

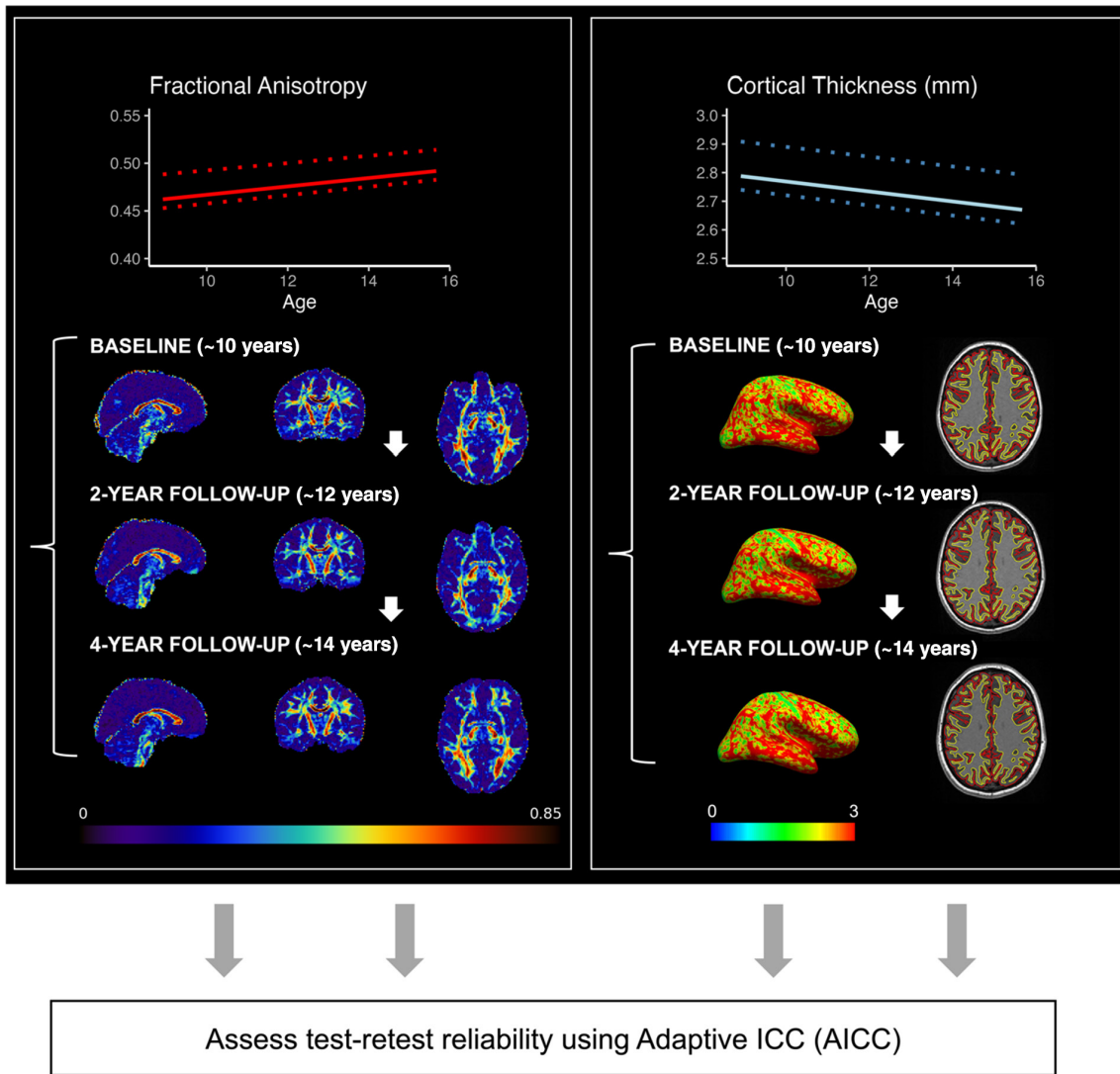


FIGURE 1 | Growth patterns of fractional anisotropy and cortical thickness in the developing human brain between ages 9 and 16 years. Normative trajectories of FA and CT at ages 9–16 years were estimated using cross-sectional linear regression, with the 25th and 95th percentiles indicated by dotted lines derived from quantile regression. The voxel-level FA map from a sample in the ABCD study across three biennial visits is presented in the bottom left. The columns on the right display the CT map of the same subject on an inflated surface, along with T1-weighted images overlaid with the pial (red) and white matter (yellow) surfaces from FreeSurfer. Longitudinal changes associated with growth can undermine conventional test–retest reliability assessments, highlighting the need for a novel method that accounts for these growth-related changes.

Nova Medical 32 channel coil. The scanner-specific sequences and sequence parameters were standardized across different scanner platform with exception of minor discrepancies due to the hardware and software constraints (Casey et al. 2018; Hagler et al. 2019).

2.2.1 | Siemens Scanners

For the T1-weighted session, a matrix of 256×256 , 176 slices, 1.0 mm isotropic resolution, a repetition time (TR) of 2500 ms, echo time (TE) of 2.88 ms, and a field of view (FOV) of 256×256 were used. T2-weighted employed the same matrix, slices, and resolution, with a TR of 3200 ms and a TE of 565 ms. The dMRI session had a matrix of 140×140 , 81 slices, 1.7 mm isotropic resolution, TR of 4100 ms, TE of 88 ms, multi-band

acceleration of 3 (Moeller et al. 2010), and a FOV of 240×240 . The protocol used 6/8 partial Fourier acquisition in phase direction.

2.2.2 | Philips Scanners

Matrix, FOV, and resolution parameters were identical across Philips and Siemens scanners in sMRI sessions. The T1-weighted session employed 225 slices, TR of 6.31 ms, and TE of 2.9 ms; the T2-weighted session used 256 slices, TR of 2500 ms, and TE of 251.6 ms. dMRI sessions shared similar slices, FOV, resolution, and multi-band acceleration factors, with variations in matrix (140×141). Philips protocol had longer TR (5300 ms) and TE (89 ms). The protocol used 0.6 partial Fourier acquisition in phase direction.

2.2.3 | GE Scanners

GE sMRI sessions mirrored Siemens sessions in matrix, FOV, resolution, and TR, with 208 slices and TE of 2ms for T1-weighted and TE of 60ms for T2-weighted. GE dMRI sessions slightly differed from Siemens dMRI sessions in TE (81.9ms). The protocol used 5.5/8 partial Fourier acquisition in phase direction.

2.3 | Imaging Data Preprocessing

2.3.1 | ABCD sMRI and dMRI Preprocessing

The details of the preprocessing pipelines from ABCD Data Analysis and Informatics Core are described elsewhere (Hagler et al. 2019). Briefly, for sMRI scans, this pipeline performs gradient distortion correction (Wald, Schmitt, and Dale 2001; Jovicich et al. 2006) using scanner-specific and nonlinear transformations provided by each scanner manufacturer, and intensity inhomogeneity correction (Sled, Zijdenbos, and Evans 1998; Ashburner and Friston 2000). T2-weighted images are then registered to T1-weighted images by maximizing mutual information (Wells et al. 1996), and resampled with a 1.0mm isotropic resolution reference brain in standard space. Subsequently, cortical reconstruction was performed using FreeSurfer 5.3.0, and reconstructed cortical surfaces were then registered to the Desikan atlas (Desikan et al. 2006). Average cortical thickness (CT) within each fuzzy-cluster parcels was calculated using smoothed surface maps (Chen et al. 2012).

dMRI scans underwent eddy current correction using a diffusion gradient and model-based approach (Zhuang et al. 2006). Dark slices affected by abrupt head motion were identified through robust tensor fitting, and frames exceeding a normalized residual error threshold were censored. Head motion correction was performed by rigid-body-registering each frame to the corresponding volume synthesized from the post-ECC censored tensor fit. Gradient nonlinearity distortions were also corrected for each frame. The dMRI images were registered to T1-weighted structural images using mutual information (Wells et al. 1996), resampled with 1.7mm isotropic resolution, and adjusted for head rotation to achieve consistent diffusion orientations across participants.

2.3.2 | ABCD Regional FA Analysis

Regional FA values were extracted for major WM fiber tracts using AtlasTrack approach that included fitting of the diffusion tensor model to the pre-processed diffusion images (Hagler et al. 2008). Diffusion tensor imaging (DTI) measures, including FA were obtained using a standard linear estimation approach, with two different tensor model fits: one excluding frames with $b > 1000\text{s/mm}^2$ (DTI inner shell) and another including all gradient strengths/shells (Alexander et al. 2007). In this study, we analyzed ABCD FA data extracted the full shell WM FA modality.

2.3.3 | Final Sample Composition

The ABCD study release 5.0 includes the baseline scans from 11,868 subjects at 21 sites. The follow up scans included: $N = 3360$ participants underwent one assessment only; 5744 participants underwent two assessments; 2619 participants completed all three assessments. After implementing the QC and preprocessing pipelines, the ABCD study provided brain WM FA data for 11,542 subjects (mean baseline age 9.9 years [SD 0.6]; 47.9% female) and morphometric measures (i.e., CT) for 11,802 subjects (mean baseline age 9.9 years [SD 0.6]; 47.8% female). Subjects with only one assessment were excluded, as longitudinal information is imperative for the evaluation of imaging data quality in our proposed method. Ninety-seven participants who went to different study sites during the baseline visit and follow-up visits were also excluded from the study. The distribution of subjects with multiple assessments across the 21 study sites is summarized in Table S1. There are 7889 subjects with FA data (mean baseline age 9.9 years [SD 0.6]; 46.6% female) and 8326 with CT data (mean baseline age 9.9 years [SD 0.6]; 46.5% female).

2.3.4 | Other dMRI and sMRI Measures

We performed a secondary reliability analysis on all DTI and sMRI traits provided by the ABCD study. For DTI, this included longitudinal diffusivity, mean diffusivity, and transverse diffusivity, all derived using the AtlasTrack approach. For sMRI, we assessed surface area, sulcal depth, T1 intensity, T2 intensity, and cortical volume, using the Desikan atlas as the reference.

2.4 | Reliability Analysis

In reliability analysis of neuroimaging data, the commonly used test-retest reliability scores are often built on the contrast of intra-subject versus inter-subject variances. A commonly used rationale is that the reliability is higher when the proportion of intra-subject variance is lower. The underlying assumption for this rationale is that each subject is measured repeatedly in the same condition (or with very high similarity). However, this assumption may not be applied to neuroimaging data that are subject to the rapid brain development of the adolescents, where intra-subject changes are assumed and the relatively lower or lack of intra-subject changes in specific tracts or individuals may even infer neurodevelopment abnormality. Therefore, the contribution of intra-subject variability to reliability estimate under typical ICC is inflated, and age correction is required for unbiased estimation of testing-retesting measures. Due to the similar enrollment age across sites, the population-level age-related brain (imaging measures) developmental trajectories is relatively invariant across the 21 sites, allowing us to integrate data from all sites to estimate the age trajectories. Specifically, we calculate the AICC as follows:

- Step 1: Fit a weighted mixed model across all-sites to estimate the age trajectory parameters.

- Step 2: For each site, perform age correction using the estimated age trajectory parameters and calculate the ICC. Update the weights based on the site-wise ICC values.
- Step 3: Repeat Steps 1 and 2 until convergence, and report the updated AICC values for all sites.

In Step 1, the weighted mixed model is expressed as follows:

$$y_{i,t} = \alpha_0 + f(\text{age}_i | \mathbf{X}_i) + \beta^T \mathbf{X}_i + b_i + \epsilon_{i,t} \quad (1)$$

where $y_{i,t}$ is the outcome (e.g., FA value on a WM fiber tract), i denotes the subject ($i = 1, \dots, n$), and t denotes the time point; the function $f(\text{age}_i | \mathbf{X}_i)$ models the developmental trajectory by age, which can be specific to sex and racial/ethnic groups; \mathbf{X}_i is the vector of demographic variables; b_i is the random intercept assumed to follow a normal distribution. Random slopes are generally not used because they tend to give rise to convergence issues due to singular fits (i.e., overfitting) (Pinheiro and Bates 2000). The age-related growth trajectory $f(\text{age}_i | \mathbf{X}_i)$ is assumed to be a population-level function, invariant across site, and is estimated using the weighted mixed model in Equation (1) across all sites. The weight is incorporated into the covariance matrix $\mathbf{V} = \mathbf{Z}^T \mathbf{G} \mathbf{Z} + \mathbf{W}^{-\frac{1}{2}} \mathbf{R} \mathbf{W}^{-\frac{1}{2}}$, where \mathbf{Z} is the design matrix giving the values of random effects to each observation, \mathbf{G} is the covariance for the random effects, and \mathbf{R} denotes the dependence between repeated measures. The weights $\{w_i\}$ are calculated in Step 2. For the first iteration, $w_i = 1$ for all i .

In Step 2, we first calculate the age-adjusted ICC for each site, denoted as ICC_s , where s represents the index for the study site ($s = 1, \dots, 21$). Let $y'_{i,s,t}$ denote the outcome for each participant i in site s at time t after correcting for the developmental trajectory $f(\text{age}_i | \mathbf{X}_i)$. The adjusted outcome $y'_{i,s,t}$ can be described as follows:

$$y'_{i,s,t} = \mu_s + b_{i,s} + \epsilon_{i,s,t} \quad (2)$$

where μ_s is the mean of all observations in site s across all visits, $b_{i,s}$ represents the random effects, and $\epsilon_{i,s,t}$ is the residual error. The random effects $b_{i,s}$ and the residuals $\epsilon_{i,s,t}$ are independent and are both assumed to be zero-mean and identically distributed (Donner and Koval 1980). Therefore, the ICC for site s can be calculated as follows:

$$\text{ICC}_s = \frac{\sigma_{b,s}^2}{\sigma_{b,s}^2 + \sigma_{\epsilon,s}^2} \quad (3)$$

The weights for each participant i in site s are then calculated based on the ICC_s using min-max scaling: $w_{i,s} = \left(\frac{\text{ICC}_s - \text{ICC}_{\min}}{\text{ICC}_{\max} - \text{ICC}_{\min}} \right)^2$.

By iterating Steps 1 and 2, the function $f(\text{age}_i | \mathbf{X}_i)$ can be estimated based on all participants across sites while minimizing the influence of measurement errors on the age-trajectory estimation. In practice, AICC converges quickly. Since AICC is built upon the mixed model, it is robust to drop-outs. Extensive simulation analysis was performed, and the results demonstrated that AICC can estimate the underlying ICC more accurately than classical testing-testing measures. Detailed

results from the simulation study are provided in Supporting Information: [SI.1](#).

3 | Results

3.1 | Reliability of CT and FA

The AICC was higher for the whole-brain gray matter CT measurements ($\text{AICC} = 0.76 \pm 0.12$), followed by whole-brain FA (0.61 ± 0.19). The AICC measurements per site are presented in Figure 2a and in Table S2. Sites that used Siemens scanners showed higher AICC on average than GE and Philips for both FA (0.71 ± 0.12 vs. 0.46 ± 0.17) and CT (0.80 ± 0.10 vs. 0.69 ± 0.11). Both differences were significant ($p < 0.001$) (Figure 2b).

The site-wise AICC measurements for CT and FA showed significant and positive correlation, that is, sites with higher AICC for CT also had higher AICC for FA ($r = 0.86$, $p < 0.001$) (Figure 2c). This correlation was driven by differences in scanners—partial correlation adjusting for scanner manufacturers dropped to 0.62 with $p = 0.005$.

3.2 | Regional Differences

Regional differences in AICC for FA and CT are shown in Figure 3 and Tables S3 and S4. For FA values, the highest AICC values were observed for the Superior Cortico Striatum (~0.75 (SD 0.16)). The lowest AICC were observed in Forceps Minor and Fornix (excluding fimbria) with a mean value of 0.49 (SD 0.17) and 0.53 (SD 0.28), respectively. The regional AICC for CT were lowest for Cingulate and Parahippocampal gyri and highest for temporal pole and insula.

3.3 | Age Trajectories

For Siemens scanners, the whole-brain FAs showed significant correlation with age ($r = 0.38$, $p < 0.001$) (Figure 4). Data collected from GE and Philips scanners showed lower correlations (FA: $r = -0.05$, $p < 0.001$). CT demonstrated higher reliability across all sites and show smaller differences between the trajectories obtained from Siemens ($r = -0.35$, $p < 0.001$) and non-Siemens sites ($r = -0.31$, $p < 0.001$), but the effect size of the correlation was still attenuated in non-Siemens sites.

3.4 | Reliability of all DTI and sMRI Traits

The test-retest reliability analysis of additional DTI measures (longitudinal, mean, and transverse diffusivity) and sMRI measures (sulcal depth, T1 intensity, T2 intensity, and cortical volume) revealed similar patterns to our primary findings. Consistent with FA and CT, data collected from Siemens scanner sites demonstrated higher reliability compared with non-Siemens sites. The average AICC for DTI traits from Siemens sites was 0.75 ± 0.10 , significantly higher (Wilcoxon test, $p < 0.001$) than the AICC for non-Siemens sites, which was 0.42 ± 0.11 . Similarly, the mean AICC for sMRI traits from Siemens sites

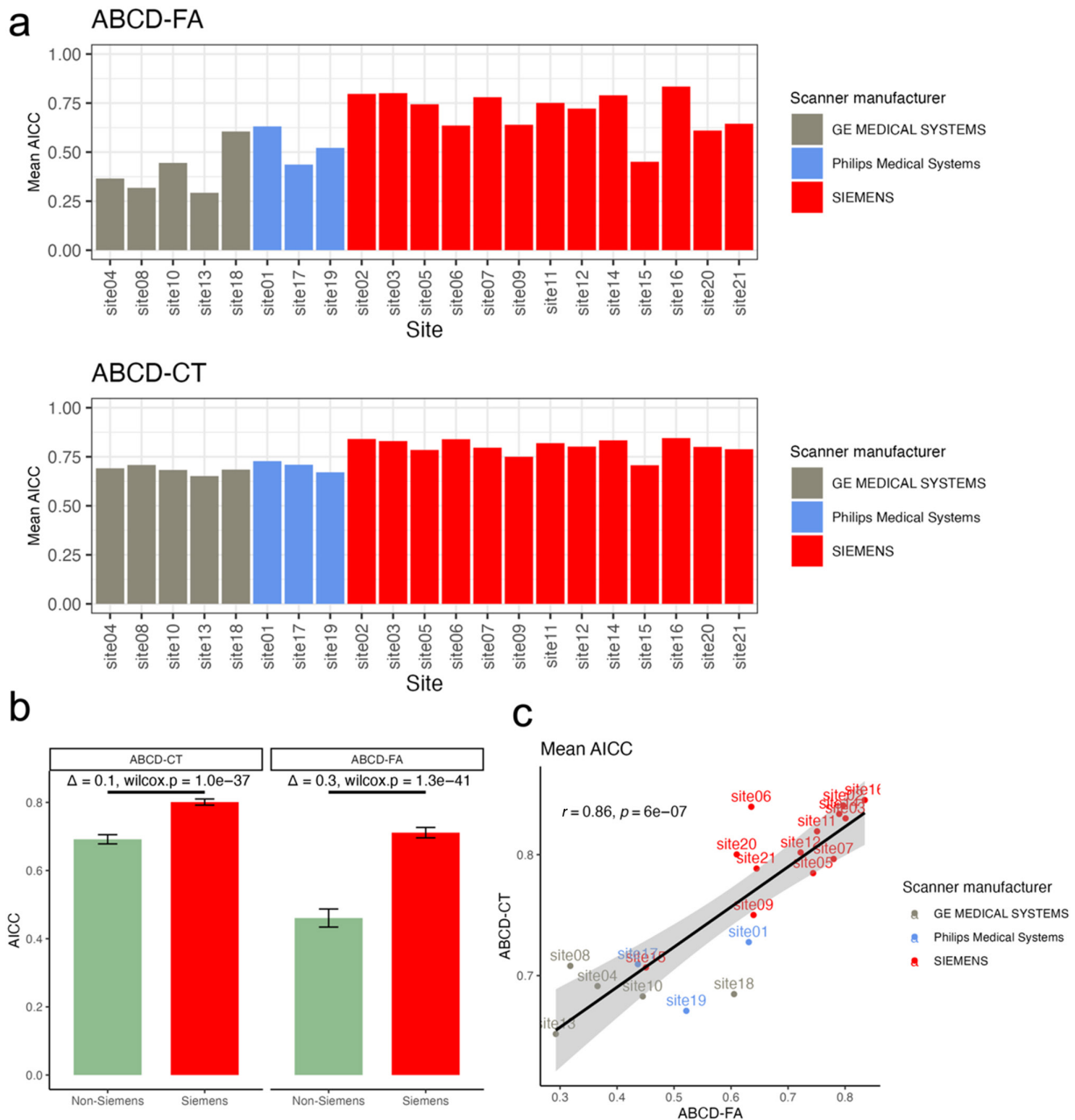


FIGURE 2 | Mean AICC across brain regions for each site. The bar plots demonstrate the mean data reliability (i.e., AICC) of each site in the ABCD study. *Red* bars indicate sites using Siemens scanners, *blue* Philips, and *gray* GE. Panel (a) shows the mean AICC across all brain regions for FA and CT, respectively. Overall, data reliability was better in Siemens sites for both FA (AICC = 0.71 ± 0.12 vs. 0.46 ± 0.17 , $p < 0.001$) and CT (AICC = 0.80 ± 0.10 vs. 0.69 ± 0.11 , $p < 0.001$). Morphometry measures (i.e., CT) has higher reliability (AICC = 0.76 ± 0.12) than FA measures (AICC = 0.61 ± 0.19) and less variations across sites-and-scanners. Panel (b) displays the significantly higher AICC of FA and CT from Siemens sites compared with non-Siemens sites. The error bars indicate standard errors across sites and regions. Panel (c) shows the correlation of mean AICC between the two measures. Linear regression analysis showed a significant and positive correlation ($r = 0.86$, $p < 0.001$). This correlation was primarily driven by differences in scanners.

was also significantly higher than from non-Siemens sites (0.78 ± 0.04 vs. 0.72 ± 0.02 , $p < 0.001$). Overall, sMRI measures showed greater reliability with less variability across sites and scanners compared with DTI measures. Individually, all measures demonstrated significantly higher reliability at Siemens sites compared with non-Siemens sites at the $\alpha = 0.001$ level, except for sulcal depth ($p = 0.01$) and T2 intensity ($p = 0.14$). Detailed results are provided in Table S5 and Figure S2.

4 | Discussion

ABCD is the largest longitudinal cohort tracking adolescent brain development using standardized protocols for multi-site data collection. We examined reproducibility of FA values extracted using the standard ABCD DTI workflows and that of CT derived from the ABCD structural MRI pipelines. We analyzed neuroimaging data collected over a 5-year period and therefore

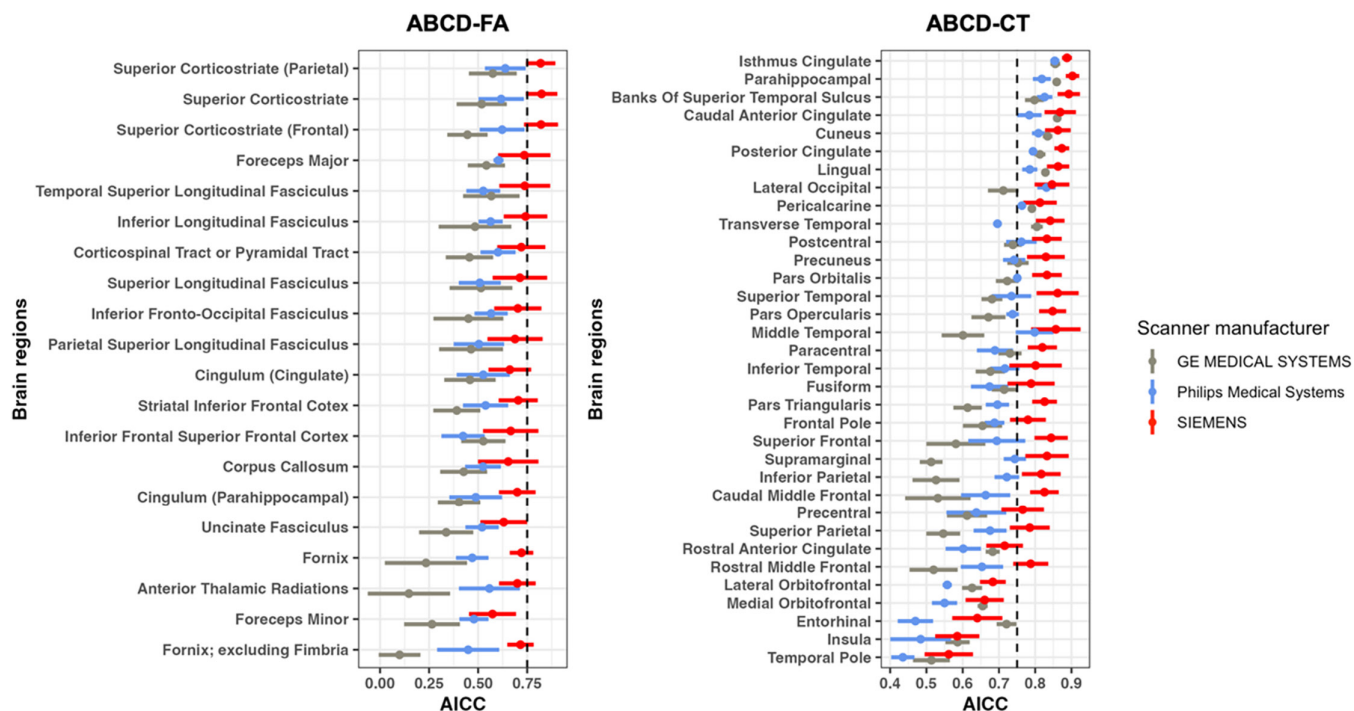


FIGURE 3 | Regional differences of the mean AICC across sites, grouped by the type of scanners used. The forest plots represent the mean regional AICC across sites using scanners from the same manufacturer, with error bars indicating standard deviation. *Red* color shows the measures obtained from Siemens scanners, *blue* Philips, and *gray* GE. The dotted line is at the conventionally good ICC score, 0.75. Overall, for FA, Siemens sites have the highest AICC with small variations across all the individual tracts among the three types of scanners. For CT, differences in AICC related to scanners are notably smaller compared with both FA measures. Nevertheless, the temporal pole, insula, and entorhinal cortex consistently exhibited low mean AICC across sites for CT, regardless of the scanner type employed.

expected that maturational change in FA and CT. We developed a novel iterative ICC approach that considers population-level developmental trends. We showed that despite the great efforts in ensuring consistent data collection across Siemens, Philips, and GE scanners, the scanner-manufacturer-related variations were large, with data collected on Philips and GE platforms showing significantly poorer test-retest reliability. We also showed that this manufacturer-specific reproducibility difference was reflected in both diffusion and CT data. The scanner-related data differences in ABCD structural and functional data have been reported before (Nielson et al. 2018; Sinha and Raamana 2023), although they focused on baseline scans only. Here we provided quantification of reliability using longitudinal data. We speculate on two likely causes for these scanner-related issues across these analyses.

The first likely source of nonbiological variance that lowered the reproducibility metrics for GE and Philips scanners may include possible protocol deviation and differences in machine hardware and software technology that cannot be homogenized across the manufacturers. The deviations from prescribed protocol parameters are more pervasive for data collected on Philips and GE versus Siemens scanners (Sinha and Raamana 2023). The deviation can also be caused by operator errors. For example, the GE MR750 software may not display all slices for DTI sequence and to overcome this, operator must change the prescribed TR parameter and then change it back to the protocol mandated value. This can introduce human error, unlike the Siemens platform that does not require this step. The deviation of the protocol can also be caused by the differences in the automated

scanner optimization algorithms built in by manufacturers that lead to changes in critical sequence parameters such as echo and TR, unbeknown to the operator (Sinha and Raamana 2023). The second source of the scanner related difference may arise from the overall stability of the magnet and gradient systems, differences in the k-space sampling trajectories and others. For example, Philips and GE scanners used by ABCD have slower gradient slew rates and therefore longer echo spacing. This led the ABCD protocol to use more aggressive partial Fourier imaging. Siemens DTI protocol collected 75% of the k-space versus 69% for GE and only 55% for Philips. The slower gradient performance forced the DTI protocol for Philips scanners to be split into two parts that are concatenated at the analysis stage. This was not the case for Siemens and GE scanners where all data were collected in a single step. The default direction in which k-space is transversed is also different between Siemens and GE and Philips scanners. These differences are hard to harmonize and likely cause complex interaction with the multi-slice excitation, leading to differences in signal, noise, distortions, and Nyquist ghost locations that reduce the overall reproducibility.

We observed significant scanner effects in the by-region AICC analysis with data collected using Siemens scanners showing uniformly higher AICC across all sites. The regional differences in AICC were also in agreement with reproducibility study of FA values derived from the Enhancing Neuro Imaging Genetic Meta Analyses (ENIGMA) DTI pipeline (Acheson et al. 2017). Specifically, lower reproducibility of FA values in fornix tracts was reported in two independent cohorts: adolescents and adults. The Fornix is a long and narrow bundle that

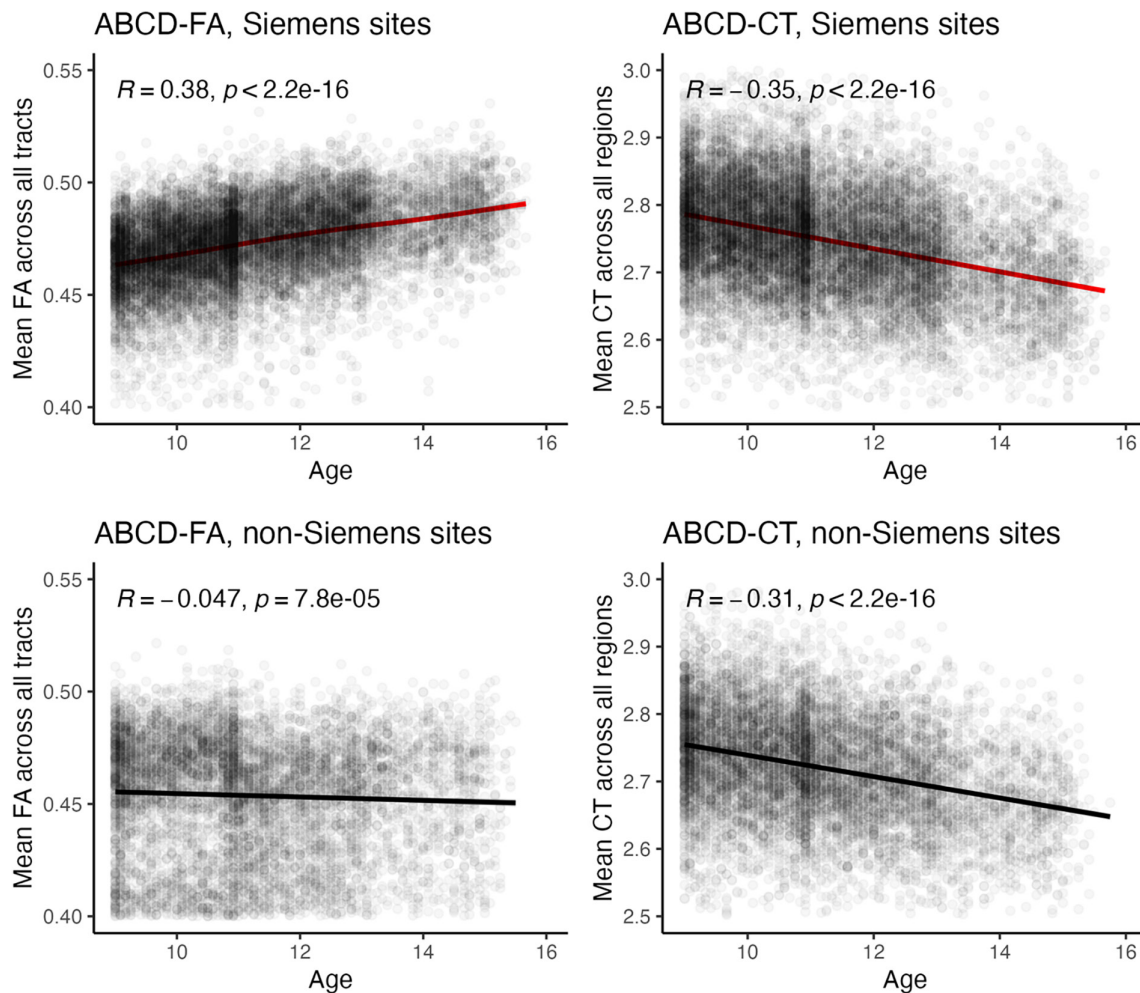


FIGURE 4 | Correlation between age and imaging outcomes. The correlations between age and FA are positive in Siemens sites, indicating a myelination of WM in development as expected ($r=0.38$, $p<0.001$), while the age effect is reversed in non-Siemens sites ($r=-0.05$, $p<0.001$). The differences in the correlations between CT values obtained from different scanners are smaller (Siemens sites: $r=-0.35$, $p<0.001$ vs. non-Siemens sites: $r=-0.31$, $p<0.001$).

is hard to correct for individual anatomical variability (Acheson et al. 2017). Additionally, another reproducibility study on multimodal MRI brain phenotypes in youth subjects reported consistent results: the ICC of FA in Forceps Minor was the lowest (0.76, 95%CI [0.61, 0.86]), followed by left Uncinate Fasciculus (ICC = 0.78, 95%CI [0.65, 0.87]). This implies that the low reproducibility observed in these tracts is likely to be attributable to technical aspects of FA measurements including alignment-related methodological confounds (e.g., misalignments in registration or partial voluming effects) rather than a failure to account for the developmental trajectory (Bach et al. 2014). Furthermore, CT measurements demonstrated consistently low AICC regions within the temporal lobe (i.e., temporal pole, insula, entorhinal cortex) and in the medial/lateral orbitofrontal lobe. Seiger et al.'s (2018) work showed that these regions have the largest difference in region-wise CT estimations using two different methods, suggesting these regions are susceptible to methodological confounds when estimating CT values (Seiger et al. 2018). The orbitofrontal region and temporal poles are commonly vulnerable tracts that may suffer from low reproducibility if the signal-to-noise ratio is low (Farrell et al. 2007; Mac Donald et al. 2011; Shahim et al. 2017).

Scanner platform had a large impact on the measured effects of biological age on longitudinal FA values in ABCD. During adolescent development, longitudinal rise in FA is hypothesized to indicate ongoing myelination of cerebral WM (Kochunov, Glahn, Nichols et al. 2011; Kochunov, Glahn, Lancaster et al. 2011). This trend was readily observed for FA collected using Siemens scanners. However, data collected on GE or Philips scanners showed a weak and negative association between FA and age. MRI-based CT is also expected to show developmental change. The ongoing myelination of WM adjacent to cortical ribbon and pruning of cortical neurons leads to reduction of measured CT during adolescence (Kochunov, Glahn, Nichols et al. 2011; Kochunov, Glahn, Lancaster et al. 2011). The negative trend in CT values was again readily detected in the data collected across Siemens, and the age effect observed in data collected from non-Siemens sites was also attenuated. The smaller age effect difference of CT between Siemens and non-Siemens sites can be explained by lower across-site variation in AICC and reduced disparity in measurement errors between Siemens and non-Siemens sites. These results underscore the substantial influence of imaging data reliability on association analyses and the credibility of neurobiological findings. The AICC of FA values from GE or

Philips scanners are significantly lower compared with FA from Siemens scanners and CT from all scanners, indicating greater measurement errors. These elevated measurement errors could obscure the association between FA values and age. Therefore, a thorough quality assessment should be conducted prior to the primary analysis to ensure the reliability and replicability of the results.

The presence of measurement errors in neuroimaging data introduces biases in association analyses. Measurement error in imaging predictors can attenuate or weaken observed brain-behavior effect sizes, while measurement error in imaging outcomes can inflate estimates of standard errors (Kenny 1979). In the context of longitudinal data analysis, nonnegligible measurement errors can also cause the ordinary maximum likelihood estimators to be inconsistent (Fuller 2009). To address measurement errors, potential strategies involve careful consideration of data quality based on both data-driven methods such as AICC and MQC measures. Quality assessment using the MQC method on the ABCD T1-weighted images revealed that 55.1% of the scans were identified as lower quality images, contributing to bias in the statistical analysis outcomes (Elyounssi et al. 2023). Therefore, the evaluation of data reliability becomes crucial for the robustness and accuracy of longitudinal multimodal imaging data analysis. While MQC is a valuable method for flexible quality assessment, its inherent labor-intensive nature poses practical challenges, particularly in the context of large-scale studies. In such scenarios, quantitative data consistency metrics like AICC can offer a more efficient and scalable alternative, allowing for data quality assessments with enhanced speed and consistency.

The reliability analysis of imaging measures derived from DTI and sMRI suggests significantly and consistently higher reliability of Siemens sites compared with non-Siemens sites. Additionally, sMRI measures show overall higher reliability with lower variability across sites compared with DTI. Future research should also consider reliability analyses that account for the growth-related trends in functional MRI (fMRI) data. To address site and scanner effects in ABCD-related analyses, one viable approach is to begin by analyzing participants from sites using Siemens scanners to establish initial associations between imaging-derived phenotypes (e.g., FA) and clinical variables. Afterward, data from sites with non-Siemens scanners can be incorporated, followed by sensitivity analyses to evaluate potential biases and assess the robustness of the findings. Additionally, advanced statistical data integration techniques can also be applied to mitigate site-specific effects (Johnson, Li, and Rabinovic 2007; Fortin et al. 2018). For example, methods such as weighting or Bayesian approaches can account for variability across different scanners and sites, allowing flexible inferences both for Siemens-specific sites and for data from all sites combined.

Ensuring high data quality is critical for large-scale studies seeking meaningful insights from imaging data analysis. Both correction methods and rigorous quality control filtering are essential for reducing measurement errors. Excluding subjects or even sites with lower data quality decreases the sample size, while maintaining an optimal sample size could result in latent bias. Therefore, the challenge lies in striking a balance between sample size and noise. Toward these goals, we have devised a

novel longitudinal reliability statistical method. This approach takes into account the expected normal age-related intra-subject variance, a crucial consideration to maximize the value of neurodevelopmental data. We identified critical scanner-type-related confounds for imaging data, particularly longitudinal diffusion imaging data collected during rapid neurodevelopment, and have recommended at least partial solutions and attention to the use of data from different sites for different types of data use.

Conflicts of Interest

The authors declare no conflicts of interest.

Data Availability Statement

Data used in this study are publicly accessible from the Adolescent Brain and Cognitive Development (ABCD) study (<https://abcdstudy.org>) in the NIMH Data Archive (NDA).

References

- Acheson, A., S. Wijtenburg, L. Rowland, et al. 2017. "Reproducibility of Tract-Based White Matter Microstructural Measures Using the ENIGMA-DTI Protocol." *Genes, Brain, and Behavior* 7, no. 2: e00615. <https://doi.org/10.1002/brb1003.1615>.
- Alexander, A. L., J. E. Lee, M. Lazar, and A. S. Field. 2007. "Diffusion Tensor Imaging of the Brain." *Neurotherapeutics* 4, no. 3: 316–329.
- Ashburner, J., and K. J. Friston. 2000. "Voxel-Based Morphometry—The Methods." *NeuroImage* 11, no. 6 Pt 1: 805–821.
- Bach, M., F. B. Laun, A. Leemans, et al. 2014. "Methodological Considerations on Tract-Based Spatial Statistics (TBSS)." *NeuroImage* 100: 358–369.
- Barnea-Goraly, N., V. Menon, M. Eckert, et al. 2005. "White Matter Development During Childhood and Adolescence: A Cross-Sectional Diffusion Tensor Imaging Study." *Cerebral Cortex* 15, no. 12: 1848–1854.
- Barnhart, H. X., M. J. Haber, and L. I. Lin. 2007. "An Overview on Assessing Agreement With Continuous Measurements." *Journal of Biopharmaceutical Statistics* 17, no. 4: 529–569.
- Bartzokis, G., P. H. Lu, K. Tingus, et al. 2010. "Lifespan Trajectory of Myelin Integrity and Maximum Motor Speed." *Neurobiology of Aging* 31, no. 9: 1554–1562.
- Basser, P. J. 1994. "Focal Magnetic Stimulation of an Axon." *IEEE Transactions on Biomedical Engineering* 41, no. 6: 601–606.
- Beaulieu, C. 2002. "The Basis of Anisotropic Water Diffusion in the Nervous System—A Technical Review." *NMR in Biomedicine* 15, no. 7–8: 435–455.
- Budde, M. D., J. H. Kim, H. F. Liang, et al. 2007. "Toward Accurate Diagnosis of White Matter Pathology Using Diffusion Tensor Imaging." *Magnetic Resonance in Medicine* 57, no. 4: 688–695.
- Casey, B. J., J. T. Nigg, and S. Durston. 2007. "New Potential Leads in the Biology and Treatment of Attention Deficit-Hyperactivity Disorder." *Current Opinion in Neurology* 20, no. 2: 119–124.
- Casey, B. J., R. M. Jones, and T. A. Hare. 2008. "The Adolescent Brain." *Annals of the New York Academy of Sciences* 1124: 111–126.
- Casey, B. J., T. Cannonier, M. I. Conley, et al. 2018. "The Adolescent Brain Cognitive Development (ABCD) Study: Imaging Acquisition Across 21 Sites." *Developmental Cognitive Neuroscience* 32: 43–54.
- Chen, C.-H., E. D. Gutierrez, W. Thompson, et al. 2012. "Hierarchical Genetic Organization of Human Cortical Surface Area." *Science (New York, N.Y.)* 335, no. 6076: 1634–1636.

- Conturo, T. E., R. C. McKinstry, E. Akbudak, and B. H. Robinson. 1996. "Encoding of Anisotropic Diffusion With Tetrahedral Gradients: A General Mathematical Diffusion Formalism and Experimental Results." *Magnetic Resonance in Medicine* 35, no. 3: 399–412.
- Desikan, R. S., F. Ségonne, B. Fischl, et al. 2006. "An Automated Labeling System for Subdividing the Human Cerebral Cortex on MRI Scans Into Gyral Based Regions of Interest." *NeuroImage* 31, no. 3: 968–980.
- Donner, A., and J. J. Koval. 1980. "The Estimation of Intraclass Correlation in the Analysis of Family Data." *Biometrics* 36, no. 1: 19.
- Elyounssi, S., K. Kunitoki, J. A. Clauss, et al. 2023. "Uncovering and Mitigating Bias in Large, Automated MRI Analyses of Brain Development." bioRxiv: 2023.2002.2028.530498.
- Farrell, J. A. D., B. A. Landman, C. K. Jones, et al. 2007. "Effects of Signal-to-Noise Ratio on the Accuracy and Reproducibility of Diffusion Tensor Imaging-Derived Fractional Anisotropy, Mean Diffusivity, and Principal Eigenvector Measurements at 1.5T." *Journal of Magnetic Resonance Imaging* 26, no. 3: 756–767.
- Flechsig, P. 1901. "Developmental (Myelogenetic) Localisation of the Cerebral Cortex in the Human." *Lancet* 158: 1027–1030.
- Fortin, J.-P., N. Cullen, Y. I. Sheline, et al. 2018. "Harmonization of Cortical Thickness Measurements Across Scanners and Sites." *NeuroImage* 167: 104–120.
- Fuller, W. A. 2009. *Measurement Error Models*. Hoboken, New Jersey: John Wiley & Sons.
- Garavan, H., H. Bartsch, K. Conway, et al. 2018. "Recruiting the ABCD Sample: Design Considerations and Procedures." *Developmental Cognitive Neuroscience* 32: 16–22.
- Gogtay, N., and P. M. Thompson. 2009. "Mapping Gray Matter Development: Implications for Typical Development and Vulnerability to Psychopathology." *Brain and Cognition* 72, no. 1: 6–15.
- Gogtay, N., J. N. Giedd, L. Lusk, et al. 2004. "Dynamic Mapping of Human Cortical Development During Childhood Through Early Adulthood." *Proceedings of the National Academy of Sciences of the United States of America* 101, no. 21: 8174–8179.
- Hagler, D. J., M. E. Ahmadi, J. Kuperman, et al. 2008. "Automated White-Matter Tractography Using a Probabilistic Diffusion Tensor Atlas: Application to Temporal Lobe Epilepsy." *Human Brain Mapping* 30, no. 5: 1535–1547.
- Hagler, D. J., S. N. Hatton, M. D. Cornejo, et al. 2019. "Image Processing and Analysis Methods for the Adolescent Brain Cognitive Development Study." *NeuroImage* 202: 116091.
- Johnson, W. E., C. Li, and A. Rabinovic. 2007. "Adjusting Batch Effects in Microarray Expression Data Using Empirical Bayes Methods." *Biostatistics* 8, no. 1: 118–127.
- Jovicich, J., S. Czanner, D. Greve, et al. 2006. "Reliability in Multi-Site Structural MRI Studies: Effects of Gradient Non-Linearity Correction on Phantom and Human Data." *NeuroImage* 30, no. 2: 436–443.
- Kalia, M. 2008. "Brain Development: Anatomy, Connectivity, Adaptive Plasticity, and Toxicity." *Metabolism* 57: S2–S5.
- Karcher, N. R., and D. M. Barch. 2021. "The ABCD Study: Understanding the Development of Risk for Mental and Physical Health Outcomes." *Neuropsychopharmacology* 46, no. 1: 131–142.
- Kenny, D. A. 1979. *Correlation and Causality*. Hoboken, New Jersey: John Wiley & Sons.
- Kochunov, P., D. C. Glahn, J. Lancaster, et al. 2011. "Fractional Anisotropy of Cerebral White Matter and Thickness of Cortical Gray Matter Across the Lifespan." *NeuroImage* 58, no. 1: 41–49.
- Kochunov, P., D. E. Williamson, J. Lancaster, et al. 2012. "Fractional Anisotropy of Water Diffusion in Cerebral White Matter Across the Lifespan." *Neurobiology of Aging* 33, no. 1: 9–20.
- Kochunov, P., D. Glahn, T. Nichols, et al. 2011. "Genetic Analysis of Cortical Thickness and Fractional Anisotropy of Water Diffusion in the Brain." *Frontiers in Neuroscience* 5, no. 120: 1–15.
- Kochunov, P., P. M. Thompson, A. Winkler, et al. 2015. "The Common Genetic Influence Over Processing Speed and White Matter Microstructure: Evidence From the Old Order Amish and Human Connectome Projects." *NeuroImage* 125: 189–197.
- Konrad, K., C. Firk, and P. J. Uhlhaas. 2013. "Brain Development During Adolescence." *Deutsches Ärzteblatt International* 110, no. 25: 425–431.
- Mac Donald, C. L., A. M. Johnson, D. Cooper, et al. 2011. "Detection of Blast-Related Traumatic Brain Injury in U.S. Military Personnel." *New England Journal of Medicine* 364, no. 22: 2091–2100.
- Madler, B., S. A. Drabycz, S. H. Kolind, K. P. Whittall, and A. L. MacKay. 2008. "Is Diffusion Anisotropy an Accurate Monitor of Myelination? Correlation of Multicomponent T2 Relaxation and Diffusion Tensor Anisotropy in Human Brain." *Magnetic Resonance Imaging* 26, no. 7: 874–888.
- Moeller, S., E. Yacoub, C. A. Olman, et al. 2010. "Multiband Multislice GE-EPI at 7 Tesla, With 16-Fold Acceleration Using Partial Parallel Imaging With Application to High Spatial and Temporal Whole-Brain fMRI." *Magnetic Resonance in Medicine* 63, no. 5: 1144–1153. <https://doi.org/10.1002/mrm.22361>.
- Nielson, D. M., F. Pereira, C. Y. Zheng, et al. 2018. "Detecting and Harmonizing Scanner Differences in the ABCD Study—Annual Release 1.0." bioRxiv: 309260.
- Pierpaoli, C., and P. J. Basser. 1996. "Toward a Quantitative Assessment of Diffusion Anisotropy." *Magnetic Resonance in Medicine* 36, no. 6: 893–906.
- Pinheiro, J. C., and D. M. Bates. 2000. *Mixed-Effects Models in S and S-PLUS*. New York: Springer.
- Rapoport, J. L., A. Addington, and S. Frangou. 2005. "The Neurodevelopmental Model of Schizophrenia: What Can Very Early Onset Cases Tell Us?" *Current Psychiatry Reports* 7, no. 2: 81–82.
- Ryan, M. C., P. Sherman, L. M. Rowland, et al. 2017. "Miniature Pig Model of Human Adolescent Brain White Matter Development." *Journal of Neuroscience Methods* 296: 99–108.
- Ryan, M. C., P. Sherman, L. M. Rowland, et al. 2018. "Miniature Pig Magnetic Resonance Spectroscopy Model of Normal Adolescent Brain Development." *Journal of Neuroscience Methods* 296: 99–108.
- Seiger, R., S. Ganger, G. S. Kranz, A. Hahn, and R. Lanzenberger. 2018. "Cortical Thickness Estimations of FreeSurfer and the CAT12 Toolbox in Patients With Alzheimer's Disease and Healthy Controls." *Journal of Neuroimaging* 28, no. 5: 515–523.
- Shahim, P., L. Holleran, J. H. Kim, and D. L. Brody. 2017. "Test-Retest Reliability of High Spatial Resolution Diffusion Tensor and Diffusion Kurtosis Imaging." *Scientific Reports* 7, no. 1: 11141.
- Shrout, P. E., and J. L. Fleiss. 1979. "Intraclass Correlations: Uses in Assessing Rater Reliability." *Psychological Bulletin* 86, no. 2: 420–428.
- Sinha, H., and P. R. Raamana. 2023. "Solving the Pervasive Problem of Protocol Non-Compliance in MRI Using an Open-Source Tool mrQA." bioRxiv: 2023.2007.2017.548591.
- Sled, J. G., A. P. Zijdenbos, and A. C. Evans. 1998. "A Nonparametric Method for Automatic Correction of Intensity Nonuniformity in MRI Data." *IEEE Transactions on Medical Imaging* 17, no. 1: 87–97.
- Song, S. K., J. Yoshino, T. Q. Le, et al. 2005. "Demyelination Increases Radial Diffusivity in Corpus Callosum of Mouse Brain." *NeuroImage* 26, no. 1: 132–140.
- Song, S. K., S. W. Sun, W. K. Ju, S. J. Lin, A. H. Cross, and A. H. Neufeld. 2003. "Diffusion Tensor Imaging Detects and Differentiates Axon and Myelin Degeneration in Mouse Optic Nerve After Retinal Ischemia." *NeuroImage* 20, no. 3: 1714–1722.

- Ulug, A. M., P. B. Barker, and P. C. van Zijl. 1995. "Correction of Motional Artifacts in Diffusion-Weighted Images Using a Reference Phase Map." *Magnetic Resonance in Medicine* 34, no. 3: 476–480.
- Wald, L., F. Schmitt, and A. Dale. 2001. "Systematic Spatial Distortion in MRI due to Gradient Non-Linearities." *NeuroImage* 6, no. 13: 50.
- Wells, W. M., P. Viola, H. Atsumi, S. Nakajima, and R. Kikinis. 1996. "Multi-Modal Volume Registration by Maximization of Mutual Information." *Medical Image Analysis* 1, no. 1: 35–51.
- Wijtenburg, S. A., F. E. Gaston, E. A. Spieker, et al. 2013. "Reproducibility of Phase Rotation STEAM at 3T: Focus on Glutathione." *Magnetic Resonance in Medicine* 72, no. 3: 603–609.
- Xue, C., J. Yuan, G. G. Lo, et al. 2021. "Radiomics Feature Reliability Assessed by Intraclass Correlation Coefficient: A Systematic Review." *Quantitative Imaging in Medicine and Surgery* 11, no. 10: 4431–4460.
- Zhuang, J., J. Hrabe, A. Kangarlu, et al. 2006. "Correction of Eddy-Current Distortions in Diffusion Tensor Images Using the Known Directions and Strengths of Diffusion Gradients." *Journal of Magnetic Resonance Imaging* 24, no. 5: 1188–1193.
- Zuo, X.-N., and X.-X. Xing. 2014. "Test-Retest Reliabilities of Resting-State fMRI Measurements in Human Brain Functional Connectomics: A Systems Neuroscience Perspective." *Neuroscience & Biobehavioral Reviews* 45: 100–118.

Supporting Information

Additional supporting information can be found online in the Supporting Information section.