

Accounting for motion in fMRI: What part of the spectrum are we characterizing in autism spectrum disorder?

Mary Beth Nebel^{a,b,*}, Daniel E. Lidstone^{a,b}, Liwei Wang^c, David Benkeser^c, Stewart H. Mostofsky^{a,b,d}, Benjamin B. Risk^c

^a*Center for Neurodevelopmental and Imaging Research, Kennedy Krieger Institute*

^b*Department of Neurology, Johns Hopkins University School of Medicine*

^c*Department of Biostatistics and Bioinformatics, Emory University*

^d*Department of Psychiatry and Behavioral Science, Johns Hopkins University School of Medicine*

Abstract

The exclusion of high-motion participants can reduce the impact of motion in functional Magnetic Resonance Imaging (fMRI) data. However, the exclusion of high-motion participants may change the distribution of clinically relevant variables in the study sample, and the resulting sample may not be representative of the population. Our goals are two-fold: 1) to document the biases introduced by common motion exclusion practices in functional connectivity research and 2) to introduce a framework to address these biases by treating excluded scans as a missing data problem. We use a study of autism spectrum disorder to illustrate the problem and the potential solution. We aggregated data from 545 children (8-13 years old) who participated in resting-state fMRI studies at Kennedy Krieger Institute (173 autistic and 372 typically developing) between 2007 and 2020. We found that autistic children were more likely to be excluded than typically developing children, with 29.1% and 16.1% of autistic and typically developing children excluded, respectively, using a lenient criterion and 80.8% and 59.8% with a stricter criterion. The resulting sample of autistic children with usable data tended to be older, have milder social deficits, better motor control, and higher intellectual ability than the original sample. These measures were also related to functional connectivity strength among children with usable data. This suggests that the generalizability of previous studies reporting naïve analyses (i.e., based only on participants with usable data) may be limited by the selection of older children with less severe clinical profiles because these children are better able to remain still during an rs-fMRI scan. We adapt doubly robust targeted minimum loss based estimation with an ensemble of machine learning algorithms to address these data losses and the resulting biases. The proposed approach selects more edges that differ in functional connectivity between autistic and typically developing children than the naïve approach, supporting this as a promising solution to improve the study of heterogeneous populations in which motion is common.

Keywords: causal inference, confounding, functional connectivity, missing data, sampling bias, super learner, targeted minimum loss based estimation

*716 N Broadway, Baltimore, MD 21205

Email address: mb@jhmi.edu (Mary Beth Nebel)

1. Introduction

Resting-state functional magnetic resonance imaging (rs-fMRI) relies on spontaneous, interregional correlations in blood-oxygen-level-dependent signal fluctuations, termed functional connectivity, to characterize brain organization (Biswal et al., 1995). A fundamental challenge in rs-fMRI-based research is to separate the signal reflecting neural activity from a combination of unstructured thermal noise and spatiotemporally structured signals of non-interest. Participant head motion is problematic because even sub-millimeter movements can introduce spatially variable artifacts that are challenging to correct during postprocessing (Power et al., 2012; van Dijk et al., 2012; Satterthwaite et al., 2012). Post-acquisition motion quality control (QC) procedures involve two stages: 1) elimination of scans with gross motion (scan exclusion); and 2) minimization of artifacts due to subtle motion (de-noising). Guidelines for removing motion-corrupted rs-fMRI data have been proposed (Satterthwaite et al., 2013; Parkes et al., 2018; Power, 2017), and many post-acquisition cleaning procedures have been developed (Satterthwaite et al., 2013; Power et al., 2014; Muschelli et al., 2014; Pruim et al., 2015; Mejia et al., 2017; Power et al., 2020). However, this work has focused on maximizing rs-fMRI data quality. The impact of scan exclusion on the study sample composition and selection bias has been largely unexamined.

Motion is particularly common in pediatric and clinical populations (Fassbender et al., 2017; Greene et al., 2018). The focus on maximizing rs-fMRI data quality has been driven by a concern that if motion artifacts are not rigorously cleaned from the data, they may introduce spurious functional connectivity differences between groups of interest. For example, autism spectrum disorder (ASD) is a neurodevelopmental condition affecting approximately 1 in 44 children in the United States that is characterized by impairments in social and communicative abilities as well as restricted interests and repetitive behaviors (Maenner et al., 2021; American Psychiatric Association, 2013). The ‘connectivity hypothesis’ of autism claims that short-range connections are increased at the expense of long-range connections within the brain (for a review, see Vasa et al. (2016)). However, sub-millimeter

28 motion-related artifacts often mimic this pattern. High-motion participants show stronger
29 correlations between nearby brain locations and weaker correlations between distant brain
30 regions compared to low-motion participants, even after controlling for motion in multiple
31 modeling steps [Power et al. \(2012\)](#); [van Dijk et al. \(2012\)](#); [Satterthwaite et al. \(2012\)](#). ASD
32 functional connectivity studies have found conflicting patterns of widespread hypoconnec-
33 tivity, hyperconnectivity, and mixtures of the two ([Di Martino et al., 2011](#); [Supekar et al.,](#)
34 [2013](#); [Keown et al., 2013](#); [Dajani and Uddin, 2016](#); [Lombardo et al., 2019](#)). Moreover, studies
35 using stricter motion QC have reported largely typical patterns of functional connectivity
36 ([Tyszka et al., 2014](#)), suggesting that motion artifacts may have contributed to discrepancies
37 in the literature ([Deen and Pelphrey, 2012](#)).

38 Exclusion of high-motion participants may help alleviate motion artifacts in functional
39 connectivity estimates but may also introduce a new problem by systematically altering the
40 study population. Implementation of scan exclusion guidelines can lead to drastic reductions
41 in sample size. For instance, in a study examining the impact of motion artifact de-noising
42 procedures on predictions of brain maturity from rs-fMRI data, [Nielsen et al. \(2019\)](#) ex-
43 cluded 365 of 487 participants between 6 and 35 years of age due to excessive head motion.
44 Applying similarly stringent scan exclusion criteria to rs-fMRI data from the Adolescent
45 Brain Cognitive Development (ABCD) study, [Marek et al. \(2019\)](#) excluded 40% of partic-
46 ipants despite efforts by the ABCD study to track head motion in real-time ([Dosenbach](#)
47 [et al., 2017](#)) to ensure a sufficient amount of motion-free data would be collected from each
48 participant ([Casey et al., 2018](#)). One strategy for balancing the need to rigorously clean the
49 data with the cost of excluding participants has been to use less stringent scan exclusion
50 criteria and then examine the effect of diagnosis in a linear model controlling for age and
51 summary measures of between-frame head motion (e.g., [Di Martino et al. 2014](#)). However,
52 the possibility of introducing selection bias following scan exclusion remains.

53 In studies with missing data, an estimate of an association may be biased if the data are
54 not missing at random in the sense that the difference in the mean outcome between the

55 groups of interest in the observed data differs from the difference if all data were observed
56 (Rubin, 1976). In rs-fMRI studies, if scan exclusion changes the distribution of participant
57 characteristics related to functional connectivity, naïve estimators of group-level functional
58 connectivity based only on participants with usable rs-fMRI data may be biased. In the case
59 of ASD, studies excluding high-motion participants have reported functional connectivity
60 differences between autistic and typically developing children, as well as associations between
61 functional connectivity strength and the severity of motor and social skill deficits (Uddin
62 et al., 2013; Lake et al., 2019; Wymbs et al., 2021; D’Souza et al., 2021), but these studies did
63 not examine the impact of scan exclusion on the composition of the study sample with usable
64 data. The graph in Figure 1 illustrates how excluding high-motion participants could obscure
65 the relationship between a diagnosis of ASD (A) and functional connectivity (Y) by changing
66 the joint distribution of diagnosis and a covariate related to symptom severity (W). If autistic
67 children with usable rs-fMRI data are phenotypically more similar to typically developing
68 children than those that were excluded, observed group differences may be reduced relative
69 to group differences if we were able to collect usable rs-fMRI data from all participants.

70 In this study, we first describe our motivating dataset, an aggregation of phenotypic and
71 rs-fMRI data from 173 autistic children and 372 typically developing children who partici-
72 pated in one of several neuroimaging studies at Kennedy Krieger Institute (KKI) between
73 2007 and 2020. We then explore the impact of commonly used head motion exclusion crite-
74 ria on the composition of the sample of participants with usable rs-fMRI data, so that we
75 can better understand what part of the spectrum we are characterizing after accounting for
76 motion. Next we introduce a method for estimating functional connectivity adjusting for the
77 observed sampling bias following participant exclusion due to motion QC, which we call the
78 deconfounded group difference (Figure 1, grey panel). We propose to treat the excluded rs-
79 fMRI scans as a missing data problem. We use an ensemble of machine learning algorithms to
80 estimate the relationship between behavioral phenotypes and rs-fMRI data usability, which
81 is called the propensity model, and between behavioral phenotypes and functional connectiv-

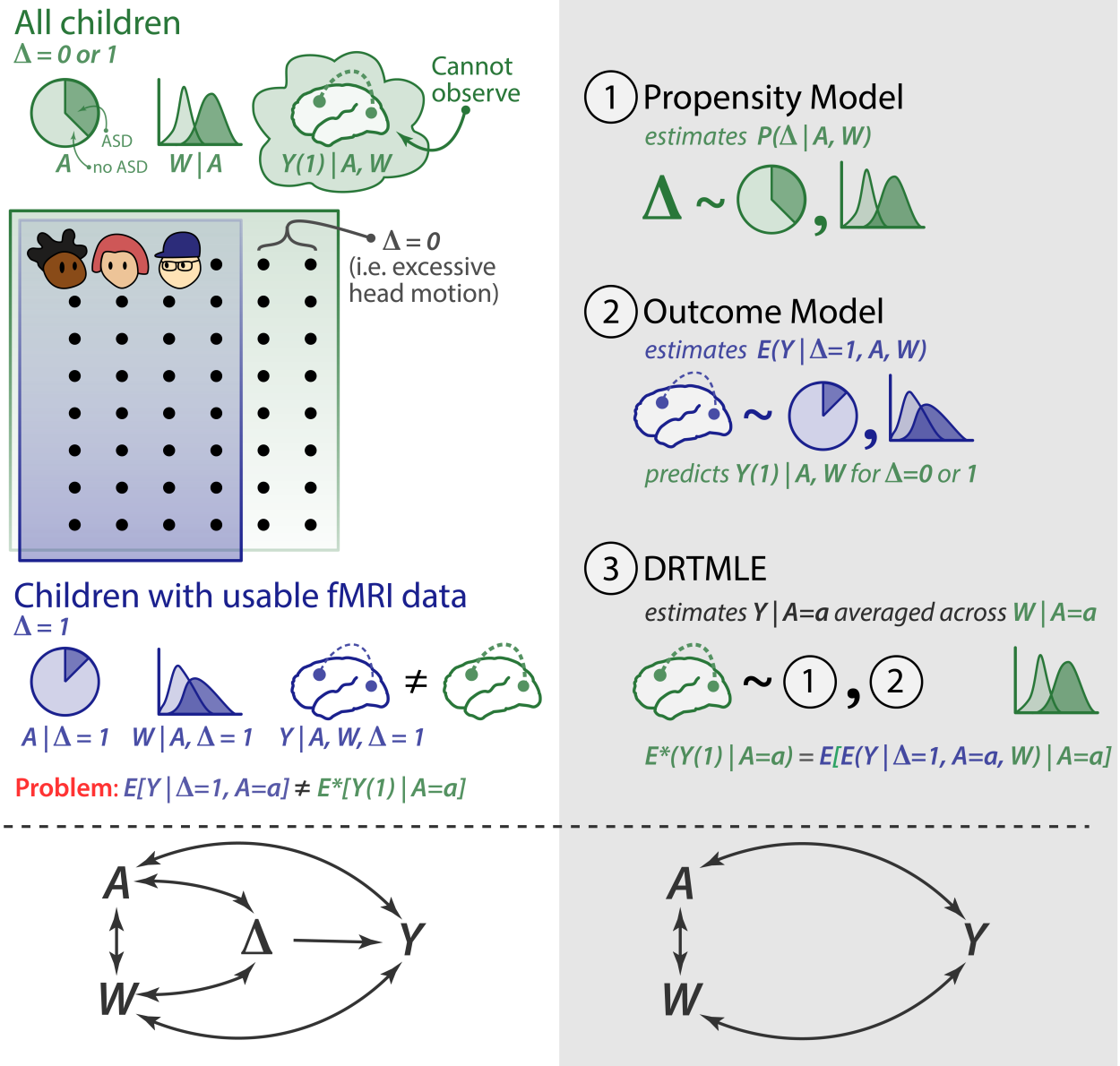


Figure 1: **Scan exclusion may induce confounding.** A indicates diagnosis, where $A = 0$ (lighter shading) represents the typically developing group and $A = 1$ (darker shading) represents the autism spectrum disorder (ASD) group. W represents a covariate that reflects symptom severity; Δ indicates resting-state fMRI usability, where $\Delta=1$ is usable and $\Delta=0$ is unusable. Y is the functional connectivity between two brain regions. Additional details are in Section 2.3.1. (Left panel) Children with usable resting-state fMRI data (in purple) may systematically differ from all enrolled children (in green). If the distribution of W differs between children with usable and unusable fMRI data ($W \leftrightarrow \Delta$) and W is related to functional connectivity ($W \leftrightarrow Y$), then naïve estimators of group-level functional connectivity based only on participants with usable data may be biased. (Right panel) We propose to address this confounding using doubly robust targeted minimum loss based estimation (DRTMLE), which involves three steps. 1. Fit the propensity model. 2. Fit the outcome model, which predicts functional connectivity from the covariates for participants with usable rs-fMRI data. Then use this model to predict functional connectivity for both usable and unusable participants. 3. Apply the DRTMLE algorithm, which uses the inverse probability of usability from step 1 and predictions of functional connectivity for all subjects (usable and unusable) from step 2.

82 ity, which is called the outcome model. The propensity and outcome models are then used in
83 the doubly robust targeted minimum loss based estimation (DRTMLE) of the deconfounded
84 group difference (Benkeser et al., 2017; van der Laan and Rose, 2011; van der Laan et al.,
85 2007). We apply this approach to estimate the deconfounded group difference between autis-
86 tic and typically developing children in the KKI dataset and compare our findings to the
87 naïve approach. Finally, we discuss the costs and benefits of motion quality control and our
88 proposed solution.

89 2. Methods

90 2.1. Dataset

91 2.1.1. Study Population

92 Our initial cohort is an aggregate of 545 children between 8- and 13-years old who partic-
93 ipated in one of several neuroimaging studies at Kennedy Krieger Institute (KKI) between
94 2007 and 2020. Participants included 173 autistic children (148 boys) and 372 typically de-
95 veloping children (258 boys); rs-fMRI scans and a limited set of phenotypic data from 266 of
96 these children (78 with ASD) were previously shared with the Autism Brain Imaging Data
97 Exchange (ABIDE) (Di Martino et al., 2014, 2017). Participants were recruited through
98 local schools, community-wide advertisement, volunteer organizations, medical institutions,
99 and word of mouth. The data collecting studies were all approved by the Johns Hopkins
100 University School of Medicine Institutional Review Board. After providing a complete study
101 description, informed consent was obtained from a parent/guardian prior to the initial phone
102 screening; written informed consent and assent were obtained from the parent/guardian and
103 the child, respectively, upon arrival at the initial laboratory visit.

104 Children were ineligible to participate if their full scale intelligence quotient (FSIQ) from
105 the Wechsler Intelligence Scale for Children, Fourth or Fifth Edition (WISC-IV or WISC-V;
106 (Wechsler, 2003)) was less than 80 and they scored below 65 on 1) the Verbal Comprehen-
107 sion Index and 2) the Perceptual Reasoning Index (WISC-IV) or the Visual Spatial Index

108 and the Fluid Reasoning Index (WISC-V), depending on which version of the WISC was
109 administered. Children were also excluded if they had a) a history of a definitive neurological
110 disorder, including seizures (except for uncomplicated brief febrile seizures), tumor, lesion,
111 severe head injury, or stroke, based on parent responses during an initial phone screening;
112 b) a major visual impairment; or c) conditions that contraindicate or make it challenging
113 to obtain MRI data (e.g., cardiac pacemaker, surgical clips in the brain or blood vessels, or
114 dental braces).

115 A diagnosis of ASD was determined using the Autism Diagnostic Observation Schedule-
116 Generic (Lord et al., 2000) or the Autism Diagnostic Observation Schedule, Second Edition
117 (ADOS-2) (Lord and Jones, 2012), depending on the date of enrollment. Diagnosis was
118 verified by a board-certified child neurologist (SHM) with more than 30 years of experience
119 in the clinical assessment of autistic children. Autistic children were excluded if they had
120 identifiable causes of autism (e.g., fragile X syndrome, Tuberous Sclerosis, phenylketonuria,
121 congenital rubella), documented history of prenatal/perinatal insult, or showed evidence of
122 meeting criteria for major depression, bipolar disorder, conduct disorder, or adjustment dis-
123 order based on parent responses during an initial phone screening. Within the ASD group,
124 a secondary diagnosis of attention deficit hyperactivity disorder (ADHD) was determined
125 using the DSM-IV or DSM-5 (APA, 2000; APA, 2013) criteria and confirmed using a struc-
126 tured parent interview, either the Diagnostic Interview for Children and Adolescents-IV
127 (DICA-IV; Reich (2000)) or the Kiddie Schedule for Affective Disorders and Schizophrenia
128 for School-Age Children (K-SADS; Kaufman et al. (2013)), as well as parent and teachers
129 versions of the Conners-Revised (Conners, 1999) or the Conners-3 Rating Scale (Conners,
130 2008), and parent and teacher versions of the DuPaul ADHD Rating Scale (DuPaul et al.,
131 1998). To be classified as having comorbid ASD and ADHD (ASD+ADHD), a child with
132 ASD had to receive one of the following: 1) a t-score of 60 or higher on the inattentive or
133 hyperactive subscales of the Conners' Parent or Teacher Rating Scale, or 2) a score of 2 or 3
134 on at least 6 of 9 items on the Inattentive or Hyperactivity/Impulsivity scales of the ADHD

135 Rating Scale-IV (DuPaul et al., 1998). Diagnosis was verified by a board-certified child neu-
136 rologist (SHM) or clinical psychologist with extensive experience in the clinical assessment
137 of children with ADHD. Children taking stimulant medications were asked to withhold their
138 medications the day prior to and the day of their study visit to avoid the effects of stimulants
139 on cognitive, behavioral, and motor measures.

140 Children were excluded from the typically developing group if they had a first-degree
141 relative with ASD, if parent responses to either the DICA-IV or for more recent participants,
142 the K-SADS, revealed a history of a developmental or psychiatric disorder, except for simple
143 phobias, or if they scored above clinical cut-offs on the parent and teacher versions of the
144 Conners' and ADHD Rating Scales.

145 The Hollingshead Four-Factor Index was used to generate a composite score of family so-
146 cioeconomic status (SES) for each participant based on each parent's education, occupation,
147 and marital status (Hollingshead, 1975). Higher scores reflect higher SES.

148 *2.1.2. Phenotypic Assessment*

149 Available phenotypic data varied according to the study in which participants enrolled.
150 The severity of core ASD symptoms was quantified within the ASD group using scores
151 from the ADOS or the ADOS-2 calibrated to be comparable across instrument versions
152 (Hus et al., 2014). Higher total scores indicate more severe ASD symptoms. These semi-
153 structured ASD observation schedules are rarely administered to control participants; they
154 were not designed to characterize meaningful variability in unaffected individuals, and scores
155 are usually equal or close to zero in typically developing children. However, ASD-like traits
156 vary among non-clinical individuals, with those meeting criteria for a diagnosis of ASD
157 falling at one extreme of a spectrum encompassing the population at large. To supplement
158 ADOS information, parent and teacher responses to the Social Responsiveness Scale (SRS)
159 questionnaire (Constantino and Todd, 2003) or the SRS-2 (Constantino and Gruber, 2012)
160 were also used. The SRS asks a respondent to rate a child's motivation to engage in social
161 interactions and their ability to recognize, interpret, and respond appropriately to emotional

162 and interpersonal cues. The SRS yields a total score ranging between 0 and 195, with a
163 higher total score indicating more severe social deficits. Total raw scores were averaged
164 across respondents.

165 We also quantified the severity of ADHD symptoms using parent responses to the Du-
166 Paul ADHD Rating Scale (DuPaul et al., 1998) due to the high comorbidity of ASD and
167 ADHD (Simonoff et al., 2008) and previous reports associating in-scanner movement with
168 ADHD-like traits (Kong et al., 2014). The DuPaul ADHD Rating Scale asks a caregiver
169 to rate the severity of inattention and hyperactivity/impulsivity symptoms over the last six
170 months and yields a total raw score as well as two domain scores: inattention and hyper-
171 activity/impulsivity. Our analyses focus on the two domain scores; higher DuPaul scores
172 indicate more severe symptoms.

173 In addition to ASD and ADHD trait severity, basic motor control was examined using
174 the Physical and Neurological Exam for Subtle Signs (PANESS), as the children were, in
175 effect, asked to complete a motor task by remaining as still as possible during the scan.
176 The PANESS assesses basic motor control through a detailed examination of subtle motor
177 deficits, including overflow movements, involuntary movements, and dysrhythmia (Denckla,
178 1985), which also allows for the observation of handedness. We focused on total motor over-
179 flow as our primary measure of motor control derived from the PANESS. Motor overflow is
180 a developmental phenomenon defined as unintentional movements that mimic the execution
181 of intentional movements. Motor overflow is common in early childhood and typically de-
182 creases as children age into adolescence. Excessive degree and abnormal persistence of motor
183 overflow is thought to reflect an impaired capacity to inhibit unintentional movements and
184 has been associated with a number of developmental and clinical conditions, in particular
185 ADHD (Mostofsky et al., 2003). Higher total motor overflow scores indicate poorer basic
186 motor control.

187 Intellectual ability was quantified using the General Ability Index (GAI) derived from
188 the WISC-IV or WISC-V (Wechsler, 2003). We used GAI because we wanted a measure

189 of intellectual ability that was independent of motor control. GAI discounts the impact of
190 tasks involving working memory and processing speed, the latter of which is abnormal in
191 ASD and associated with poor motor control (Mayes and Calhoun, 2008). Higher GAI scores
192 indicated greater intellectual ability.

193 *2.1.3. Study Sample*

194 The available phenotypic data varied according to the study in which participants en-
195 rolled. The study sample for our application of the deconfounded group difference is defined
196 as the subset of participants with a complete set of demographic information (sex, socioe-
197 conomic status, and race) and the selected predictors along with nineteen children in which
198 motor overflow is imputed as described in Section 2.3.2. The missingness of the data is de-
199 picted in the Web Supplement Figure S.1. This subset contains 151 autistic and 353 typically
200 developing children from the original 173 autistic and 373 typically developing children, and
201 we refer to these 504 participants as the complete predictor cases. The socio-demographic
202 characteristics of the complete predictor cases are summarized in Table 1. The impacts of
203 motion exclusion criteria on this subset are discussed in Section 3.1.1.

204 *2.1.4. rs-fMRI Acquisition and Preprocessing*

205 All participants completed at least one mock scan training session to habituate to the MRI
206 environment during a study visit prior to their MRI session. Rs-fMRI scans were acquired on
207 a Phillips 3T scanner using an 8-channel or a 32-channel head coil and a single-shot, partially
208 parallel, gradient-recalled echo planar sequence with sensitivity encoding (repetition time
209 [TR]/echo time = 2500/30 ms, flip angle = 70°, sensitivity encoding acceleration factor of 2,
210 3-mm axial slices with no slice gap, in-plane resolution of 3.05 × 3.15 mm [84 × 81 acquisition
211 matrix]). An ascending slice order was used, and the first 10 seconds were discarded at the
212 time of acquisition to allow for magnetization stabilization. The duration of rs-fMRI scans
213 varied between 5 min 20 seconds (128 timepoints) and 6.75 min (162 timepoints), depending
214 on the date of enrollment.

Table 1: **Socio-demographic characteristics of complete predictor cases.** For continuous variables, mean and standard deviation (SD) are indicated; Kruskal-Wallis rank-sum tests were used to assess diagnosis group differences. For binary and categorical variables, frequencies and percentages are summarized, and differences between diagnosis groups were assessed using either the Chi-square test or Fisher’s exact test. Despite aggregating data from several studies, age and handedness were balanced between diagnosis groups. In contrast, sex, race, and socioeconomic status were imbalanced. ASD=autism spectrum disorder. TD=typically developing. SD=standard deviation.

	TD (N=353)	ASD (N=151)	p value
Sex			<0.001 ¹
Male	245 (69.4%)	127 (84.1%)	
Female	108 (30.6%)	24 (15.9%)	
Age			0.826 ²
Mean (SD)	10.363 (1.248)	10.324 (1.363)	
Range	8.020 - 12.980	8.010 - 12.990	
Race			0.004 ³
African American	36 (10.2%)	9 (6.0%)	
Asian	27 (7.6%)	3 (2.0%)	
Biracial	45 (12.7%)	12 (7.9%)	
Caucasian	245 (69.4%)	127 (84.1%)	
Socioeconomic Status			0.006 ²
Mean (SD)	54.135 (9.390)	51.964 (9.379)	
Range	18.500 - 66.000	27.000 - 66.000	
Handedness			0.364 ¹
Right	317 (89.8%)	128 (84.8%)	
Left	17 (4.8%)	12 (7.9%)	
Mixed	19 (5.4%)	11 (7.3%)	
Currently On Stimulants			
No	353 (100.0%)	97 (64.2%)	
Yes	0 (0.0%)	54 (35.8%)	

¹ Pearson Chi-Square test

² Kruskal-Wallis rank sum test

³ Fisher’s exact test

215 Rs-fMRI scans were either aborted or not attempted for seven participants in the complete
216 predictor case set (3 ASD) due to noncompliance. Rs-fMRI scans for the remaining 497 par-
217 ticipants in the complete predictor case set were visually inspected for artifacts and prepro-
218 cessed using SPM12 (Wellcome Trust Centre for Neuroimaging, London, United Kingdom)
219 and custom code written in MATLAB (The Mathworks, Inc., Natick Massachusetts), which
220 is publicly available (https://github.com/KKI-CNIR/CNIR-fmri_preproc_toolbox). Rs-
221 fMRI scans were slice-time adjusted using the slice acquired at the middle of the TR as a
222 reference, and head motion was estimated using rigid body realignment. Framewise displace-
223 ment was calculated from these realignment parameters (Power et al., 2012). The volume
224 collected in the middle of the scan was spatially normalized using the Montreal Neurological
225 Institute (MNI) EPI template with 2-mm isotropic resolution (Calhoun et al., 2017). The
226 estimated rigid body and nonlinear spatial transformations were applied to the functional
227 data in one step. Each rs-fMRI scan was linearly detrended on a voxel-wise basis to remove
228 gradual trends in the data. Rs-fMRI data were spatially smoothed using a 6-mm FWHM
229 Gaussian kernel.

230 *2.1.5. Motion QC*

231 We considered two levels of gross motion exclusion:

232 1. In the lenient case, scans were excluded/deemed unusable if the participant had less
233 than 5 minutes of continuous data after removing frames in which the participant
234 moved more than the nominal size of a voxel between any two frames (3 mm) or their
235 head rotated 3° , where a 3° rotation corresponds to an arc length equal to 2.6 mm
236 assuming a brain radius of 50 mm (Power et al., 2012) or 4.2 mm assuming 80 mm
237 (Jenkinson et al., 2002; Yan et al., 2013). This procedure was modeled after common
238 head motion exclusion criteria for task fMRI data, which rely on voxel size to determine
239 thresholds for unacceptable motion (Johnstone et al., 2006; Fassbender et al., 2017).

240 2. In the strict case, scans were excluded if mean FD exceeded .2 mm or they included

241 less than five minutes of data free from frames with FD exceeding .25 mm (Ciric et al.,
242 2017).

243 Eighty-five participants in the complete predictor case set (19 ASD) completed more than
244 one rs-fMRI scan. For these participants, if more than one scan passed the lenient level of
245 motion QC, we selected the scan with the lowest mean FD to include in our analyses.

246 *2.1.6. Group ICA and Partial Correlations*

247 Thirty components were estimated using group independent component analysis (Group
248 ICA) with 85 principal components retained in the initial subject-level dimension reduction
249 step from the scans that passed lenient motion QC (GIFT v3.0b: <https://trendscenter.org/software/gift/>; Medical Image Analysis Lab, Albuquerque, New Mexico) (Calhoun
250 et al., 2001; Erhardt et al., 2011). Detailed methods for Group ICA can be found in Allen
251 et al. (2011). We used the back-reconstructed subject-level timecourses for each independent
252 component to construct subject-specific partial correlation matrices (30x30) using ridge re-
253 gression ($\rho = 1$) (Lombardo et al., 2019; Mejia et al., 2018). After Fisher z-transforming the
254 partial correlation matrices, we extracted the lower triangle for statistical analysis. Following
255 the taxonomy for macro-scale functional brain networks in Uddin et al. (2019), we identified
256 18 signal components from the 30 group components. A partial correlation is equal to zero
257 if the two components are conditionally independent given the other components. By using
258 the partial correlations, we control for correlations due to the twelve non-signal components,
259 which include some motion artifacts, as well as components mainly composed of white mat-
260 ter or cerebrospinal fluid, which capture other signals of non-interest that impact the brain
261 globally (Bijsterbosch et al., 2020).

263 *2.2. Impact of motion QC on the sample size and composition*

264 *2.2.1. Impact of motion QC on group sample size*

265 For each level of motion exclusion, Pearson's chi-squared tests were used to assess whether
266 the proportion of excluded children differed between the ASD and typically developing

267 groups.

268 *2.2.2. rs-fMRI exclusion probability as a function of phenotypes*

269 We used univariate generalized additive models (GAMs) to examine the relationship be-
270 tween the log odds of exclusion and seven covariates: ADOS (ASD group), SRS, inattention,
271 hyperactivity/impulsivity, motor overflow, age, and GAI. We used the subset of children
272 included in the final study sample (Section 2.1.3) for the strict and lenient motion exclu-
273 sion criteria. We used automatic smoothing determined using random effects with restricted
274 maximum likelihood estimation (REML) (Wood, 2017). We used univariate models rather
275 than a model with all covariates simultaneously because some of the variables are correlated,
276 such that the impact of each variable on rs-fMRI usability may be difficult to estimate.
277 These models are related to the propensity models that will be used in the estimation of the
278 deconfounded group difference (Section 2.3.1). While the propensity models use an ensemble
279 of machine learning models to predict usability from multiple predictors, our focus for this
280 analysis is on interpretable models. We controlled for multiple comparisons using the false
281 discovery rate (FDR) for the seven univariate models, in which FDR is applied separately
282 to the lenient and strict criteria models (Benjamini and Hochberg, 1995). Although FDR
283 correction was popularized by high-throughput studies conducted in computational biology,
284 Benjamini and Hochberg (1995) originally illustrated the utility of their approach for con-
285 trolling the expected number of falsely rejected null hypotheses using a study in which a
286 moderate number of tests (15) were performed, which is comparable to our analysis.

287 *2.2.3. Impact of motion QC on distributions of phenotypes among children with usable data*

288 We examined how the distribution of ADOS (ASD group), SRS, inattention, hyperac-
289 tivity/impulsivity, motor overflow, age, and GAI differed between included and excluded
290 participants. For additional insight into how scan exclusion may differentially affect autistic
291 versus typically developing children, we stratified this analysis by diagnosis. We visualized
292 the densities using kernel density estimation with default bandwidths in ggplot2 (Wickham,

293 2016). We then used one-sided Mann-Whitney U tests to test for differences between in-
294 cluded and excluded participants for each measure stratified by diagnosis. We hypothesized
295 that 1) included children would have less severe social, inattentive, hyperactive/impulsive,
296 and motor deficits than excluded children, and 2) included children would be older and have
297 higher GAI. We controlled for multiple comparisons by applying the FDR separately to the
298 thirteen tests (7 for the ASD group and 6 for the typically developing group) performed for
299 the lenient and strict motion QC cases.

300 2.2.4. Functional connectivity as a function of phenotypes

301 We also characterized the relationship between phenotypes and functional connectivity.
302 For each level of motion exclusion, we used univariate GAMs to examine the relationship
303 between each phenotypic measure and the adjusted residuals for each edge of signal-to-
304 signal components in the partial correlation matrix. The adjusted residuals are the same
305 data inputted to the deconfounded group difference and are calculated from the residuals of
306 a linear model with mean FD, max FD, the number of frames with $FD < 0.25$ mm, sex, race,
307 socioeconomic status, and diagnosis with the effect of diagnosis added back in as described
308 in Section 2.3.2. Smoothing was determined using the random effects formulation of spline
309 coefficients with restricted maximum likelihood estimation (REML) (Wood, 2017).

310 2.3. Addressing data loss and reducing sampling bias using the deconfounded group difference

311 2.3.1. Theory: Deconfounded group difference

312 Our goal is to estimate the difference in average functional connectivity between autistic
313 and typically developing children. Let Y be a random variable denoting the functional
314 connectivity between two locations (or nodes defined using independent component analysis)
315 in the brain. In practice, these will be indexed by v and v' , but we suppress this notation for
316 conciseness. Let A denote the diagnosis indicator variable equal to one if the participant has
317 ASD and zero otherwise. We first consider the hypothetical case in which all participants
318 have usable rs-fMRI data. We use the potential outcomes notation and let $Y(1)$ denote the

319 functional connectivity in this hypothetical world (Hernan and Robins, 2020). Let W denote
320 the covariates, which include measures that may be related to functional connectivity and
321 ASD severity. Our *parameter of interest* is the difference in functional connectivity between
322 autistic and typically developing children: $\psi^* = E^*(Y(1)|A = 1) - E^*(Y(1)|A = 0)$, where
323 $E^*(\cdot)$ denotes an expectation with respect to the probability measure of $\{Y(1), A, W\}$. It is
324 important to observe that W is not independent of A , as the distribution of some behavioral
325 variables differ by diagnosis group. We rewrite ψ^* using the law of iterated expectations to
326 gain insight into our parameter of interest:

$$\begin{aligned}\psi^* &= E^*(Y(1)|A = 1) - E^*(Y(1)|A = 0) \\ &= E^* \{E^*(Y(1)|A = 1, W) |A = 1\} - E^* \{E^*(Y(1)|A = 0, W) |A = 0\}.\end{aligned}$$

327 Here, the outer expectation integrates across the conditional distribution of the variables
328 given diagnosis. This estimand differs from an average treatment effect (ATE) commonly
329 considered in causal inference (Hernan and Robins, 2020), which integrates across the distri-
330 bution of the covariates for the pooled population (autistic and typically developing children).

331 In contrast to the hypothetical world, many children in the observed world move too
332 much during their rs-fMRI scan for their data to be usable, but we are still able to collect
333 important behavioral and socio-demographic covariates from them. We regard data that fail
334 motion quality control as “missing data.” Let Δ denote a binary random variable capturing
335 the missing data mechanism that is equal to one if the data are usable and zero otherwise.
336 Then data are realizations of the random vector $\{Y, A, W, \Delta\}$. Expectation with respect to
337 their probability measure is denoted $E(\cdot)$ (no asterisk). Additionally, $Y|(\Delta = 0)$ is missing.
338 Then the naïve difference is $\psi_{naive} = E(Y|\Delta = 1, A = 1) - E(Y|\Delta = 1, A = 0)$. We define
339 confounding as $\psi^* \neq \psi_{naive}$ (Greenland et al., 1999). Confounding can occur when a covariate
340 is related to data usability/missingness, $W \leftrightarrow \Delta$, and also related to functional connectivity,
341 $W \leftrightarrow Y$. Then if the covariate is related to diagnosis, i.e., $W \leftrightarrow A$, we have $\psi^* \neq \psi_{naive}$.

342 If there are interactions between W and A , then we can also have $\psi^* \neq \psi_{naive}$. These
343 relationships are summarized in the graph in Figure 1. We now define our *target parameter*
344 as a function of usable data. We call this quantity the *deconfounded group difference*:

$$\psi = E\{E(Y|A = 1, \Delta = 1, W) | A = 1\} - E\{E(Y|A = 0, \Delta = 1, W) | A = 0\}. \quad (1)$$

345 The mathematical distinction between this and the naïve estimator is that in the naïve
346 estimator, $E(Y|\Delta = 1, A = 1) = E\{E(Y|\Delta = 1, A = 1, W)|\Delta = 1, A = 1\}$, which differs
347 from $E\{E(Y|\Delta = 1, A = 1, W) | A = 1\}$, with a similar distinction for $E(Y|\Delta = 1, A = 0)$.
348 In the deconfounded group difference, we integrate across the conditional distribution of
349 phenotypic variables given diagnosis versus the naïve approach that integrates across the
350 conditional distribution of phenotypic variables given diagnosis and data usability. We will
351 show in Section 3.1.3 that the distribution of phenotypic variables given diagnosis differs
352 from the distribution given diagnosis and data usability.

353 Identifying the parameter of interest ψ^* from the target parameter ψ requires three
354 assumptions:

355 (A1.1) *Mean exchangeability*:

$$356 \quad \text{for } a = 0, 1, \quad E^*\{Y(1) | A = a, W\} = E^*\{Y(1) | \Delta = 1, A = a, W\}.$$

357 (A1.2) *Positivity*: for $a = 0, 1$ and all possible w , $P(\Delta = 1 | A = a, W = w) > 0$.

358 (A1.3) *Causal Consistency*: for all i such that $\Delta_i = 1$, $Y_i(1) = Y_i$.

359 Assumption (A1.1) implies that W is sufficiently rich as to contain all variables simultane-
360 ously associated with mean functional connectivity and exclusion due to failed motion QC.

361 This assumption is also called ignorability or the assumption of no unmeasured confounders.

362 In the missing data literature, this is closely related to the assumption that data are missing
363 at random: $P(\Delta = 1|Y, A = a, W) = P(\Delta = 1|A = a, W)$ (van der Laan and Robins, 2003).

364 Assumption (A1.2) implies that there are no phenotypes in the population who uniformly

365 fail motion QC. Assumption (A1.3) stipulates that Y from children with usable fMRI data
366 is the same as the outcome that would have been observed under a hypothetical intervention
367 that allows the child to pass motion control (Vander Weele, 2009). Under A1.1 and A1.3,
368 we have $E^*\{Y(1) | A = a, W\} = E\{Y | \Delta = 1, A = a, W\}$, which allows us to identify the
369 potential outcomes from the observable data.

370 We estimate our target using doubly robust targeted minimum loss based estimation
371 [DRTMLE, (Benkeser et al., 2017; van der Laan and Rose, 2011)], which involves three steps
372 enumerated below and illustrated in Figure 1 .

- 373 1. Fit the propensity model: $P(\Delta|A, W)$. This model characterizes the probability that
374 the rs-fMRI data pass motion quality control. It uses all data to fit the model. Then
375 the usable functional connectivity will be weighted by their inverse probabilities of
376 usability (propensities) during step three.
- 377 2. Fit the outcome model: $E(Y|\Delta = 1, A, W)$. This step estimates functional connec-
378 tivity from the covariates for participants with usable rs-fMRI data. It then predicts
379 functional connectivity for both usable and unusable participants.
- 380 3. Use DRTMLE to combine functional connectivity from the usable subjects weighted
381 by the inverse probability of usability from step 1 with predictions of functional con-
382 nectivity for all subjects (usable and unusable) from step 2. Here, DRTMLE is applied
383 separately to each diagnosis group, which calculates mean functional connectivity by
384 integrating across the diagnosis-specific distribution of the covariates from usable and
385 non-usable participants.

386 Steps 1 and 2 use super learner, an ensemble machine learning technique. The super learner
387 fits multiple pre-specified regression models and selects a weight for each model by minimizing
388 cross-validated risk (Polley et al., 2019). Step 3 combines the propensity and outcome models
389 using DRTMLE. An appealing property of DRTMLE is that the estimate of the deconfounded
390 group difference and its variance are statistically consistent even if either the propensity

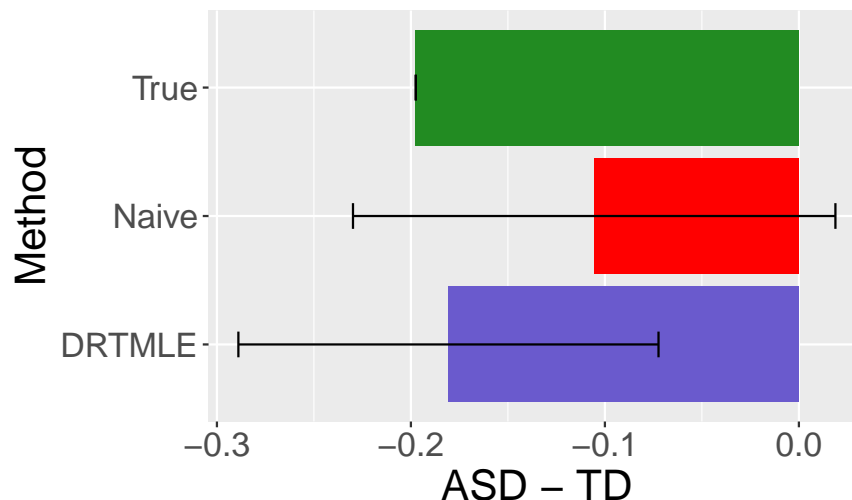


Figure 2: **An illustration of the improvement in functional connectivity from DRTMLE compared to the naïve approach from a single simulated dataset.** The true mean ASD-TD difference in functional connectivity is negative (green bar), with the true mean in the ASD group being negative and the the true mean in the TD group being slightly positive. The estimate of the mean ASD-TD difference from the naïve approach (red bar) is also negative but closer to zero due to confounding. Additionally, the 95% confidence interval includes zero. Using DRTMLE, the deconfounded group difference (purple bar) is closer to the truth and the 95% confidence interval does not include zero. Code to reproduce this example is available at [thebrisklab github](https://github.com/thebrisklab).

391 model or the outcome model is inconsistently estimated. See [Benkeser et al. \(2017\)](#). By
392 statistical consistency, we mean that our estimate converges to the true difference as the
393 sample size goes to infinity, which is different from the causal consistency assumption in
394 A1.3. Here, we know that the missingness mechanism is deterministic based on motion,
395 but we are replacing it with a stochastic model that estimates missingness based on the
396 behavioral phenotypes. Details of our implementation are in Section 2.3.2. We simulate
397 a dataset with confounding and estimate the deconfounded group difference in a tutorial
398 available at <https://github.com/thebrisklab/DeconfoundedFMRI>, and the results from
399 this simulation are included in Fig. 2.

400 2.3.2. Application: Deconfounded group difference in the KKI Dataset

401 Recall Step 1 involves fitting a propensity model and Step 2 involves fitting an outcome
402 model (Section 2.3.1). We use the same predictors in the propensity and outcome models: age
403 at scan, handedness (left, right, mixed), primary diagnosis, secondary diagnosis of ADHD,

404 indicator variable for a current prescription for stimulants (all participants were asked to
405 withhold the use of stimulants the day prior to and on the day of the scan), motor overflow,
406 GAI, DuPaul inattention, DuPaul hyperactivity/impulsivity, and ADOS. The ADOS is only
407 administered to children in the ASD group, since it is usually equal or close to zero in typically
408 developing children. We set ADOS equal to zero for all typically developing children. Social
409 responsiveness score was not included due to missing values in 19.5% of observations.

410 Since motor overflow was missing in 5.5% of children, we imputed its value using super
411 learner from all variables above plus sex, SES, and race (details of the learners described
412 below). This resulted in the imputation of motor overflow scores for nineteen children. Thus
413 the study sample for our application of the deconfounded group difference is the subset of
414 participants with a complete set of predictors after the imputation of motor overflow for
415 these nineteen children (Section 2.1.3) as depicted in Web Supplement Figure S.1.

416 We focus on the lenient motion QC case because too few participants have usable data
417 following strict motion QC to accurately estimate the outcome model. Functional con-
418 nectivity metrics based on partial correlations were recently shown to be less sensitive to
419 motion artifacts than those based on full correlations (Mahadevan et al., 2021), but to guard
420 against lingering impacts of motion on functional connectivity and to account for possible
421 confounders due to sampling design, we adjust the partial correlations as follows. For each
422 edge, we fit a linear model with mean FD, max FD, number of frames with $FD < 0.25$ mm,
423 sex (reference: female), race (reference: African American), socioeconomic status, and pri-
424 mary diagnosis (reference: Autism) as predictors. We include sex, race, and socioeconomic
425 status in this model because they differed between autistic and typically developing children
426 (see Section 3.1.1). We then extracted the residuals and added the estimated intercept and
427 effect of primary diagnosis. Then the “naïve” approach is comparable to the approach used
428 in Di Martino et al. (2014), who included diagnosis, sex, age, and mean FD in a linear model.
429 See Section 4.4 for additional discussion.

430 Steps 1 and 2 (see Section 2.3.1): We use the following learners and R packages when

431 using super learner: multivariate adaptive regression splines in the R package earth, lasso in
432 glmnet, generalized additive models in gam, generalized linear models in glm, random forests
433 with ranger, step-wise regression in step, step-wise regression with interactions, xgboost,
434 and the intercept only (mean) model; for the outcome model (continuous response), we
435 additionally used ridge from MASS and support vector machines in e1071. Parameters were
436 set to their defaults except for the following: the family was equal to binomial (logistic
437 link) in the propensity model with method set to minimize the negative log likelihood; in
438 the motor overflow and outcome models, the method was set to minimize the squared error
439 loss. Note the outcome model is fit separately for each of the 153 edges, whereas the same
440 propensities are used for all edges. The propensity model is fit using the complete predictor
441 cases. The outcome model is fit using the complete usable cases.

442 Step 3: DRTMLE is applied to ASD for each edge, then to TD. This step uses both
443 the propensities and the predicted outcomes to result in an estimate of the deconfounded
444 mean for the ASD group, the deconfounded mean for the typically developing group, and
445 their variances. We use the non-parametric regression option for both the reduced-dimension
446 propensity and reduced-dimension outcome regression. A z-statistic is formed from their dif-
447 ference under the assumption of independent groups, which is used to test the null hypothesis
448 that functional connectivity is equal in autistic and typically developing children.

449 Since super learner uses cross validation, its results differ for different random seeds. We
450 ran the entire procedure (motor overflow imputation, propensity model, and 153 outcome
451 models) for two hundred different seeds, calculated the DRTMLE-based z-statistic for the
452 difference in functional connectivity, and averaged the z-statistics at each edge from the
453 two hundred seeds. We calculated adjusted p-values using FDR=0.2, which means that we
454 expect 20% of the rejected null hypotheses to be falsely rejected. This threshold has been
455 used in recent papers on FDR ([Barber and Candès, 2015](#)). We also report edges that survive
456 the more stringent FDR=0.05. We repeated this entire procedure a second time with a
457 different set of 200 seeds. The correlation between the average z-statistics across the 153

458 edges was greater than 0.999. Eleven edges were selected at false discovery rate FDR=0.20
459 in the first set of seeds and nine of these eleven edges in the second set. The same two edges
460 were selected in both sets for FDR=0.05. For the final input to the figures, we pooled both
461 sets of seeds and averaged their z-statistics, which resulted in eleven edges at FDR=0.20.

462 We examined the stability of the propensity scores across the first five random seeds.
463 Propensities near zero can increase the bias and variance of causal effects (Petersen et al.,
464 2010) and indicate a possible violation of the positivity assumption (A1.2). The smallest
465 propensity ranged from 0.30-0.36. This indicates that there is a reasonable probability of data
466 inclusion across the range of $\{W, A\}$ and that Assumption (A1.2) is likely to be adequately
467 satisfied. The AUCs for predicting usability across the five seeds ranged from 0.75 to 0.92,
468 whereas the AUC was 0.68 using logistic regression and 0.69 using a logistic additive model,
469 which indicates that the super learner often improves the accuracy of the propensity model.

470 For the naïve approach, we calculated the z-statistic of the average group differences
471 between autistic and typically developing children from the complete usable cases for each of
472 the 153 edges. This test statistic is nearly equivalent to the t-statistic from the linear model
473 with motion variables, sex, socioeconomic status, and diagnosis.

474 *2.4. Data and code availability*

475 All data used for this study can be made available by written request through the study's
476 corresponding author under the guidance of a formal data-sharing agreement between in-
477 stitutions that includes the following: 1) using the data only for research purposes and not
478 attempting to identify any participant; 2) limiting analyses to those described in both insti-
479 tutions IRB-approved protocols; and 3) no redistribution of any shared data without a data
480 sharing agreement.

481 The code for recreating all analyses, tables, and figures in this study is available at
482 <https://github.com/thebrisklab/DeconfoundedfMRI>.

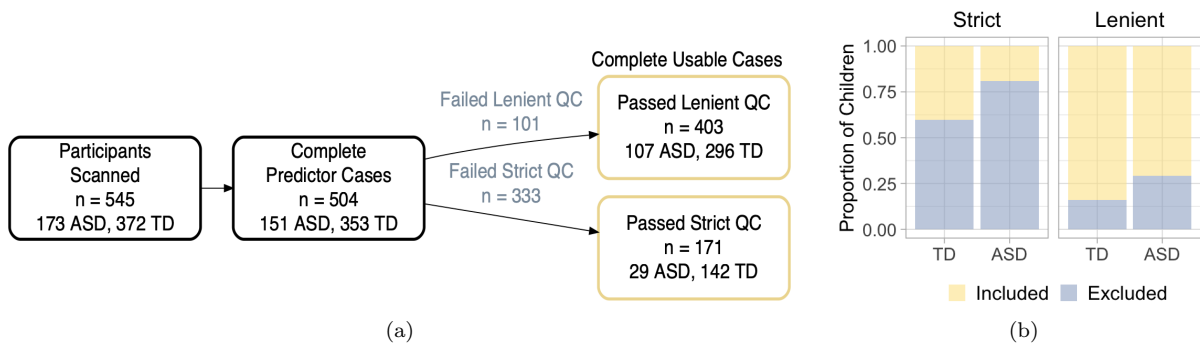


Figure 3: **Motion quality control leads to dramatic reductions in sample size.** a) Flow chart of inclusion criteria for this study showing the number of participants remaining after each exclusion step. Lenient motion quality control (QC) excluded 20% of complete predictor cases, while strict motion QC excluded 66% of complete predictor cases. b) The proportion of children in each diagnosis group whose scans were included (yellow) and excluded (lavender) using the strict (left) and lenient (right panel) gross motion QC. A larger proportion of children in the autism spectrum disorder (ASD) group were excluded compared to typically developing (TD) children using lenient motion QC ($\chi^2=10.3$, $df = 1$, $p=0.001$) and strict ($p<0.001$).

483 3. Results

484 3.1. Impact of motion QC on the study sample and sample bias

485 3.1.1. The impact of motion QC on sample size can be dramatic and differs by diagnosis 486 group

487 Figure 3 illustrates the inclusion criteria used for our analyses and the number of partic-
488 ipants remaining after each exclusion step. Missing covariate data excluded 41 participants,
489 or 7.5% of the total number of participants scanned. Lenient motion QC excluded 20.0% of
490 complete predictor cases, while strict motion QC excluded 66.1% of complete predictor cases.
491 In addition, we found the proportion of excluded children differed by diagnosis group using
492 both levels of motion QC (Figure 3b). Using lenient motion QC, 16.1% of typically develop-
493 ing children were excluded, compared to 29.1% of children in the ASD group ($\chi^2=10.3$, $df =$
494 1, $p=0.001$). Using strict motion QC, 59.8% of typically developing children were excluded,
495 compared to 80.8% of children in the ASD group ($p<0.001$). Thus, commonly used motion
496 QC procedures resulted in large data losses that more severely impacted the size of the ASD
497 group.

498 *3.1.2. rs-fMRI exclusion probability changes with phenotype and age*

499 We observed that children with higher ADOS scores, SRS scores, inattentive symptoms,
500 hyperactive/impulsive symptoms, or poorer motor control were more likely to be excluded,
501 while older children and children with higher GAI were less likely to be excluded when the
502 lenient motion QC was used (all FDR-adjusted $p < 0.01$) as well as the strict motion QC
503 (all FDR-adjusted $p < 0.03$) (Figure 4). In particular, there is a sharp increase in exclusion
504 probability using the lenient motion QC for children with higher ADOS scores (lavender
505 line, left-most panel). The bottom panel of Figure 4 illustrates the covariate distribution
506 for each diagnosis group (pooling included and excluded participants). Interestingly, using
507 the lenient motion QC, the relationship between SRS and exclusion appears flatter over the
508 range of values in the typically developing group and steeper over the range of values in the
509 ASD group (lavender line). In contrast, the relationship between hyperactivity/impulsivity
510 and exclusion appears linear over the range of values present in the typically developing
511 group but fairly flat over the range of values in the ASD group.

512 *3.1.3. Phenotype representations differ between included and excluded children*

513 Figure 5 illustrates distributions of the covariates for included and excluded participants
514 stratified by diagnosis group and motion QC level. For the lenient motion QC, median
515 values for included and excluded participants, U statistics, and FDR-adjusted p values for
516 each measure and diagnosis group are summarized in Web Supplement Table S.1. Using the
517 lenient motion QC, we observed biases in both the ASD and typically developing groups
518 toward the selection of older children (FDR-adjusted $p = 0.04$, 0.04 for the ASD and typically
519 developing groups, respectively) with higher GAI (FDR-adjusted $p = 0.03$, 0.04 for ASD and
520 typically developing groups, respectively). In the ASD group, we also observed biases toward
521 the selection of children who had lower total ADOS, SRS, or motor overflow scores (FDR-
522 adjusted $p = 0.03$, 0.03 , and 0.01 , respectively), but we did not observe differences in terms of
523 inattentive or hyperactive/impulsive symptoms between included and excluded participants
524 (FDR-adjusted $p > 0.6$ for both covariates). In the typically developing group, we did not

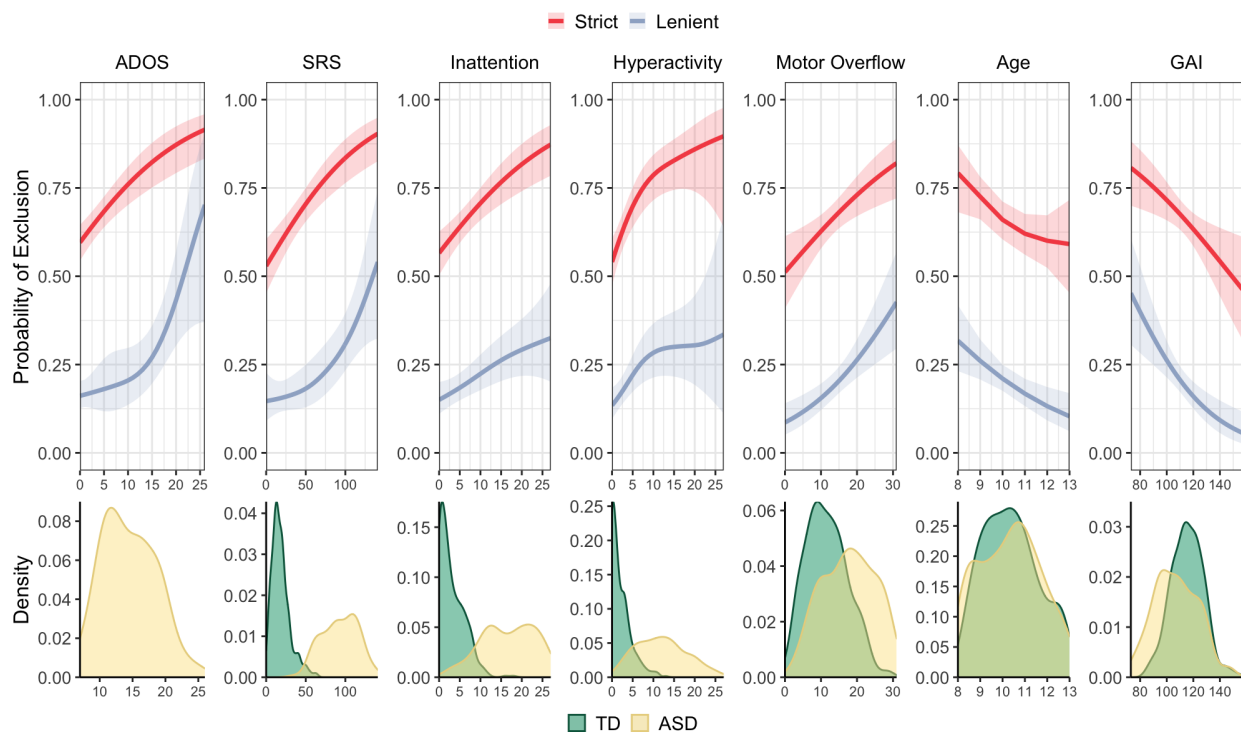


Figure 4: **rs-fMRI exclusion probability changes with phenotype and age.** Univariate analysis of rs-fMRI exclusion probability as a function of participant characteristics. From left to right: Autism Diagnostic Observation Schedule (ADOS) total scores, social responsiveness scale (SRS) scores, inattentive symptoms, hyperactive/impulsive symptoms, total motor overflow, age, and general ability index (GAI) using the lenient (lavender lines, all FDR-adjusted $p < 0.01$), and strict (red lines) motion quality control (all FDR-adjusted $p < 0.03$). Variable distributions for each diagnosis group (included and excluded scans) are displayed across the bottom panel (TD=typically developing, green; ASD=autism spectrum disorder, yellow).

525 observe a bias in terms of SRS or inattention (FDR-adjusted $p=0.5, 0.4$), while there was
526 some evidence of bias for motor overflow ($p=0.08$). We did observe a bias towards the
527 selection of typically developing children with lower hyperactive/impulsive scores (FDR-
528 adjusted $p=0.04$).

529 Differences between included and excluded children also tended to occur using the strict
530 criteria, although in general significance was reduced, owing in part to the reduced sample
531 size in the included group. Typically developing children who were included were less hy-
532 peractive/impulsive than typically developing children who were excluded (FDR-adjusted
533 $p=0.02$). Median values for included and excluded participants, U statistics, and FDR-
534 adjusted p values for each measure and diagnosis group are summarized in Web Supplement
535 Table [S.2](#).

536 *3.1.4. Phenotypes are also related to functional connectivity*

537 The relationships we observed between rs-fMRI data usability and the covariates exam-
538 ined in the preceding analyses may impact our parameter of interest if those measures are also
539 related to functional connectivity. Figure [6](#) illustrates histograms of p values for GAMs of
540 the relationship between edgewise functional connectivity (adjusted for sex, SES, race, and
541 motion, see Section [2.3.2](#)) and ADOS, SRS, inattentive symptoms, hyperactive/impulsive
542 symptoms, total motor overflow, age, and GAI across participants with usable rs-fMRI data
543 using the lenient motion QC (lavender bins) and the strict motion QC (red bins). This
544 analysis is related to the outcome model used in the deconfounded group difference, as it
545 provides insight into whether the sampling bias will lead to confounding. Here, we focus on
546 a single phenotype in each GAM for interpretability. For a given phenotype, a clustering of
547 p values near zero suggests that a covariate is associated with functional connectivity for a
548 greater number of edges. If there is no association between the covariate and functional con-
549 nectivity, we expect the p values to be more uniformly distributed. We see strong clustering
550 of p values near zero for total ADOS across participants with usable rs-fMRI data under
551 lenient and strict motion QC. For SRS, we see clustering using participants who pass strict

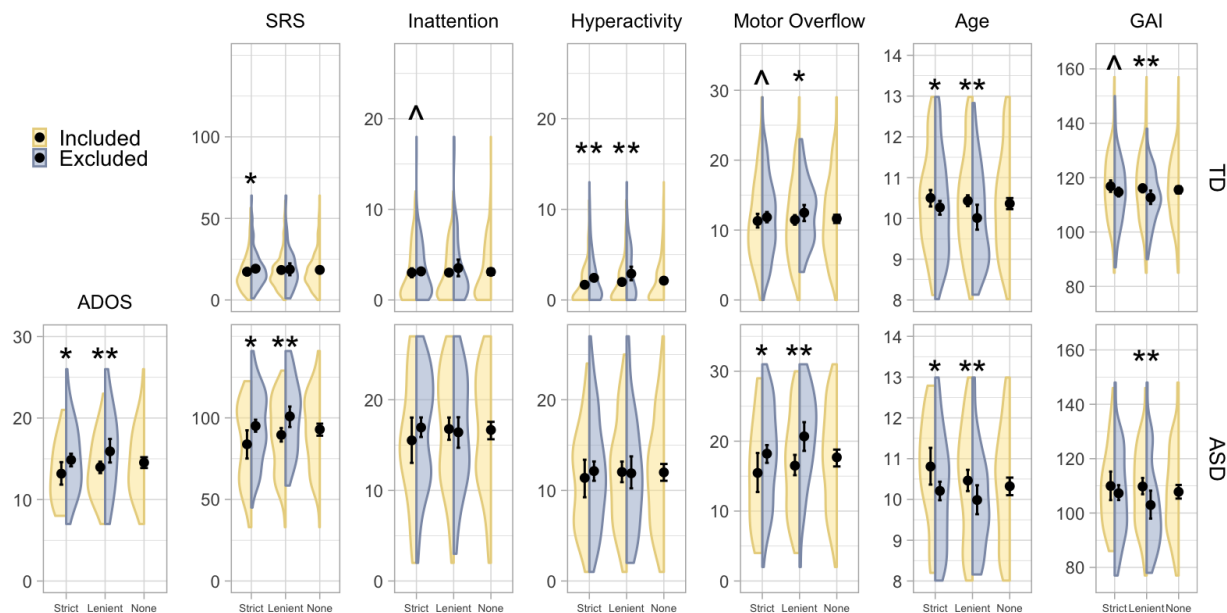


Figure 5: Participants with usable rs-fMRI data differed from participants with unusable rs-fMRI data. Comparison of Autism Diagnostic Observation Schedule (ADOS) scores, social responsiveness scale (SRS) scores, inattentive symptoms, hyperactive/impulsive symptoms, motor overflow, age, and general ability index (GAI) for included (yellow) and excluded (lavender) participants stratified by diagnosis group and motion exclusion level. The deconfounded mean integrates across the diagnosis-specific distribution of usable and unusable covariates for the variables described in Section 2.3.2, which here is labeled as “None.” We controlled for 13 comparisons performed for the lenient and strict motion QC cases using the false discovery rate (FDR). ** indicate differences between included and excluded participants with an FDR-adjusted p value <0.05; * indicate FDR-adjusted p values <0.1; ^ indicate FDR-adjusted p values <0.2. A larger number of significant differences are observed using the lenient motion QC than the strict motion QC, but very few participants pass strict motion QC. autism spectrum disorder (ASD), typically developing (TD). The R code to produce these split violin plots was adapted from DeBruine (2018).

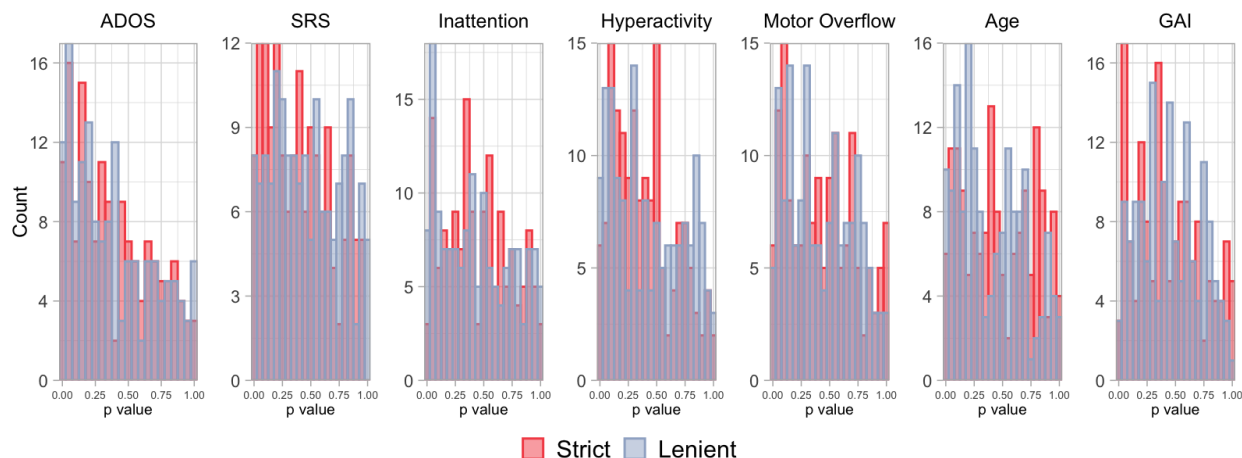


Figure 6: **Some covariates related to rs-fMRI exclusion probability are also related to functional connectivity.** Histograms of p values for generalized additive models of the relationship between edgewise functional connectivity in participants with usable rs-fMRI data and (from left to right) ADOS, social responsiveness scale (SRS) scores, inattentive symptoms, hyperactive/impulsive symptoms, total motor overflow as assessed during the Physical and Neurological Exam for Subtle Signs, age, and general ability index (GAI). For a given covariate, a clustering of p values near zero suggests that covariate is associated with functional connectivity for a greater number of edges. Several covariates appear to be related to functional connectivity using both the lenient motion quality control (lavender bins) and the strict motion quality control (red bins).

552 motion QC, but this pattern is less apparent for participants who pass lenient motion QC.
553 For inattentive symptoms, we see a clustering of p values near zero using participants who
554 pass lenient motion QC. For hyperactive/impulsive symptoms, we see a clustering of p values
555 near zero following both levels of motion QC. For motor overflow, we see some clustering
556 of p values near zero using participants who pass strict motion QC, but this pattern is less
557 clear in participants passing lenient motion QC. For age, we see a clustering of p values near
558 zero using participants who pass the lenient motion QC but not strict. For GAI, we see a
559 clustering of p values near zero using participants who pass the strict motion QC but not
560 lenient.

561 3.2. Application: Deconfounded group difference in the KKI Dataset

562 The deconfounded group difference estimated using DRTMLE revealed more extensive
563 differences between the ASD and typically developing groups than the naïve approach (Fig. 7,
564 Web Supplement Table S.3). At FDR=0.20, the naïve approach indicated three edges show-
565 ing a negative difference in functional connectivity between the ASD and typically developing

566 groups (blue lines) and three edges showing a positive difference (red lines). The DRTMLE
567 approach also indicated these six edges, with five having smaller p values relative to those for
568 the naïve approach. The DRTMLE approach also indicated an additional two edges showing
569 negative group differences and three additional edges showing a positive group difference.
570 Network nodes that gained edges from the DRTMLE versus the naïve method (FDR=0.2)
571 included the default mode network (DMN) (+3), executive control (+3), somatomotor (+1),
572 visual (+1), ventral attention (+1), and dorsal attention (+1). At FDR=0.05 (Web Supple-
573 ment Figure S.2), the naïve approach only indicated one edge showing a negative difference
574 in functional connectivity between the ASD and typically developing groups (primary visual
575 IC-02 to bilateral control IC-27). The DRTMLE approach indicated this edge, while also
576 indicating one other edge showing a negative difference in functional connectivity (postero-
577 lateral cerebellum IC-14 to dorsal attention IC-19).

578 Functional connectivity scores further from zero reflect stronger functional connectivity
579 regardless of sign; positive scores reflect stronger positive partial correlations, or more inte-
580 grated intrinsic activity between nodes. Negative scores reflect negative partial correlations,
581 or more segregated intrinsic activity between nodes. The sign of average group effects re-
582 mained consistent, as did the direction of group differences (Web Supplement Table S.3).
583 Edges showing positive group differences in functional connectivity included edges for which
584 positive correlations were strengthened in the ASD group compared to the typically devel-
585 oping group, as well as connections in which negative correlations were weaker in the ASD
586 group compared to the typically developing group. Similarly, the edges showing negative
587 group differences included connections for which negative correlations were strengthened in
588 the ASD group compared to the typically developing group, as well as connections for which
589 positive correlations were weaker in the ASD group compared to the typically developing
590 group.

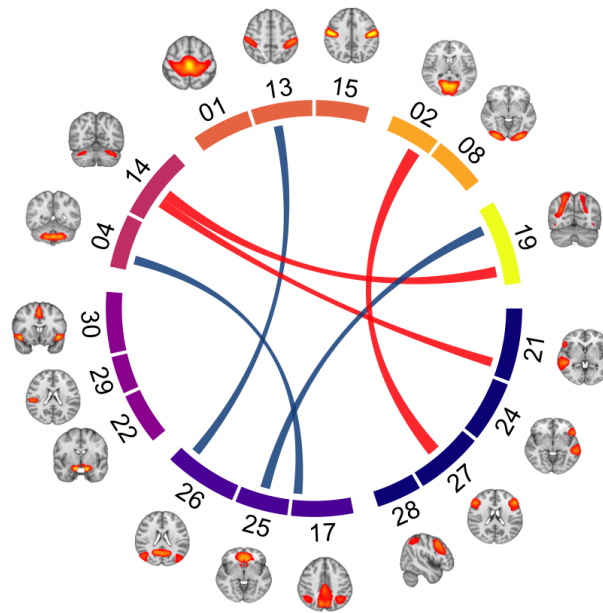
591 In this application, the deconfounded means were very similar to the naive means (Web
592 Supplement Figure S.3). Additionally, partial correlations were highly variable, with the

593 range of partial correlations in the ASD and typically developing groups broadly overlapping.
594 This indicates that the more extensive differences indicated by DRTMLE Fig. 7 are largely
595 driven by smaller standard errors, rather than by extensive confounding following lenient
596 motion QC, which is discussed in Section 4.2.

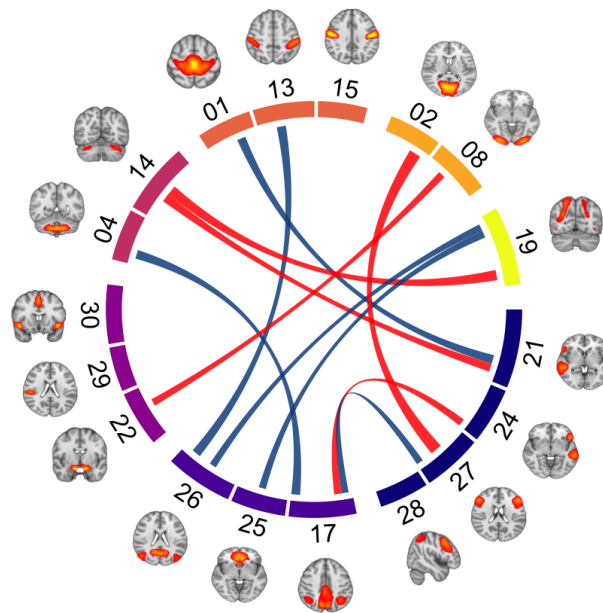
597 4. Discussion

598 We set out to understand what part of the autism spectrum we are characterizing in rs-
599 fMRI analyses: Does excluding high-motion participants allow us to draw conclusions about
600 average brain function/connectivity that are representative of 8-to-13-year-old children across
601 the entire autism spectrum, or does it introduce bias? The primary message that emerges
602 from our findings is that ignoring confounding due to motion exclusion can be problematic
603 scientifically. Using data from a large sample of autistic children without an intellectual
604 disability and typically developing children, we demonstrated that motion exclusion changes
605 the distribution of behavioral and sociodemographic traits in the study sample that are
606 related to functional connectivity. This finding suggests that the generalizability of previous
607 studies reporting naïve analyses may be limited by the selection of older children with less
608 severe clinical profiles because these children are better able to remain still during an rs-fMRI
609 scan. We further propose a statistical approach for addressing the data loss and possible
610 confounding following motion QC using DRTMLE; our findings indicate more extensive
611 differences between autistic and typically developing children using DRTMLE as compared
612 with conventional approaches.

613 In our study, the impact of motion QC on sample size was dramatic and differed by
614 diagnosis group. Additionally, rs-fMRI exclusion probability changed with symptom sever-
615 ity and age. Detailed reporting of the number of participants excluded for excessive head
616 motion is far from standard practice, but we found that motion QC removed a larger pro-
617 portion of autistic children compared to typically developing children, which is consistent
618 with the patterns reported in [Redcay et al. \(2013\)](#) and [Jones et al. \(2010\)](#). Across diagnosis



(a) Naïve Z-Statistic



(b) DRTMLE Z-Statistic

Figure 7: **The DRTMLE deconfounded group difference revealed more extensive differences than the naïve approach.** Z-statistics for autism spectrum disorder (ASD) versus typically developing (TD) using a) the naïve test and b) using DRTMLE. Connections are thresholded using a false discovery rate (FDR) of 0.20. Blue lines indicate ASD > TD (3 in naïve, 6 in DRTMLE). Red lines indicate ASD < TD (3 in naïve, 5 in DRTMLE). Brain regions contributing to each independent component are illustrated and components are grouped by functional assignment. Navy nodes: control. Blue violet: default mode. Purple: salience/ventral attention. Magenta: pontomedullary/cerebellar. Coral: somatomotor. Orange: visual. Yellow: dorsal attention. FDR=0.05 is plotted in Web Supplement Figure S.2. These plots were generated using the circlize package in R (Gu et al., 2014) and the tutorial provided by Mowinckel (2018).

619 groups, children with more severe social deficits, more inattentive symptoms, more hyperac-
620 tive/impulsive symptoms, or poorer motor control were more likely to have unusable rs-fMRI
621 data and be excluded, while older children or children with higher intellectual ability were
622 less likely to be excluded for both levels of motion QC. Similarly, [Simhal et al. \(2021\)](#) found
623 that children with ASD and children with ADHD who failed a mock MRI training protocol
624 were younger, had lower verbal and non-verbal intelligence scores, and more severe ADOS
625 scores than children with ASD and children with ADHD who passed the training protocol.
626 These findings suggest that the mechanisms driving missingness in rs-fMRI studies may be
627 related to scientifically relevant participant characteristics.

628 The estimate of mean functional connectivity should be representative of all children en-
629 rolled in the study, but we observed that participants with usable rs-fMRI data differed from
630 participants with unusable rs-fMRI data that would have been excluded using conventional
631 approaches. Autistic children who were excluded following lenient motion QC tended to
632 be younger, displayed more severe social deficits (both observed by the experimenter using
633 the ADOS and reported by parents/teachers using the SRS), more motor overflow, or lower
634 intellectual ability than autistic children who were included. We observed similar differences
635 between included and excluded autistic children following strict motion QC, although in
636 general power was reduced due to the reduced sample size. Moreover, these characteristics
637 are exactly those that showed relationships with functional connectivity among children with
638 usable data following one or both levels of motion QC. The strength of these relationships
639 between clinically relevant measures and functional connectivity among children with usable
640 data appeared to depend on the level of motion QC used. For instance, evidence of a relation-
641 ship between SRS and functional connectivity was stronger among participants who passed
642 strict motion QC than among participants who passed lenient motion QC (Figure Fig. 6).
643 Given that the criteria used to define usability varies widely among rs-fMRI studies, our
644 findings suggest that differences in the representation of symptom severity among children
645 with usable data following motion QC may have partially contributed to discrepancies in the

646 literature regarding ASD-associated functional connectivity findings. To improve our abil-
647 ity to compare findings across studies, it is critical for rs-fMRI researchers to transparently
648 assess the amount of information lost following motion QC, to consider whether participant
649 characteristics related to usability are also related to the effect of interest, and to try to
650 address the loss of power and potential confounding if they are.

651 Here, we have made progress on this issue using techniques from the missing data and
652 causal inference literature combined with an ensemble of machine learning algorithms. Our
653 approach results in more extensive differences between autistic and typically developing chil-
654 dren than indicated using the naïve approach (Fig. 7). Our framework explicitly treats
655 missingness due to motion QC as a source of confounding, and we define a target param-
656 eter called the deconfounded group difference. The general concept of this framework is
657 to recognize that children with usable data are not representative of all enrolled children
658 within each diagnosis group. DRTMLE combines the results of inverse propensity weighting
659 and G-computation, which improves robustness relative to either approach alone. Inverse
660 propensity weighting gives more weight to children with more severe symptoms who have
661 usable functional connectivity data because a) they are more likely to be missing and b)
662 functional connectivity is related to symptom severity so we need them to stand in for all
663 children with more severe symptoms who are excluded due to data quality concerns. The
664 outcome model estimates functional connectivity for all children, including those with greater
665 symptom severity, and in this sense accounts for children with unusable data. We use an
666 ensemble of machine learning methods to flexibly model possible non-linear relationships
667 between phenotypic traits and data usability (the propensity model) and between pheno-
668 typic traits and functional connectivity (the outcome model). For both the propensity and
669 outcome models, we include a rich collection of variables that we expect to be associated
670 with rs-fMRI usability, functional connectivity, or both. Including variables that contribute
671 to both rs-fMRI usability and functional connectivity represents an opportunity to decrease
672 bias. Including variables that contribute to functional connectivity but not necessarily to

673 rs-fMRI usability represents an opportunity to decrease the variance of our estimate without
674 increasing bias. The propensity and outcome models are then combined using DRTMLE,
675 which results in statistically consistent estimation of the deconfounded group difference and
676 its variances under the assumptions in Section 2.3.1 and discussed in Section 4.3.

677 *4.1. Possible scientific insights gained from DRTMLE*

678 Collectively, the findings suggest that group differences in functional connectivity are
679 more robust using the DRTMLE approach as compared to the naïve approach. What evi-
680 dence do we have to support the validity of these findings? First, the sign of average group
681 effects remained consistent across methods, as did the direction of group differences (Web
682 Supplement Table S.3). In this study, DRTMLE appears to have a larger effect on the vari-
683 ance than on the group means, and we discuss the implications of this and effect size in
684 Section 4.2. Second, the pattern of group differences observed using DRTMLE is consistent
685 with knowledge of DMN-DAN interactions, such that the DMN shows task-induced deactiva-
686 tion, whereas the DAN shows task-induced activation (Padmanabhan et al., 2017). Findings
687 from task-based fMRI studies suggest that individuals with ASD show lower deactivation of
688 the DMN during self-referential processing tasks as compared to typically developing controls
689 (Kennedy et al., 2006; Padmanabhan et al., 2017). Recent findings also suggest a crucial
690 role of the posterolateral cerebellum, a region functionally connected to the DMN (Buckner
691 et al., 2011), in both social mentalizing (Van Overwalle et al., 2020) and behaviors central
692 to a diagnosis of ASD (Lidstone et al., 2021; Stoodley et al., 2017). The cerebellum is also
693 believed to form and update internal models of the world for predictive control in both social
694 and nonsocial contexts (Blakemore et al., 2001). Functional connectivity between the DMN,
695 DAN, and cerebellum networks should be a focus of future research to better understand
696 the neural mechanisms contributing to autism diagnosis.

697 *4.2. Sample size limitations, inference, and effect size*

698 In this study, the differences between the edges selected using DRTMLE versus the naïve
699 approach appear to be largely driven by decreases in the standard error of the estimates
700 rather than by changes in the mean difference (Web Supplement Figure S.3). DRTMLE can
701 be used to address data loss by improving efficiency, which can result in smaller standard
702 errors relative to the naïve approach. The TMLE framework leverages all available covariate
703 data, and when the covariate data are predictive of the outcome, this can improve statistical
704 power (Moore and van der Laan, 2009). One potential limitation is that DRTMLE underes-
705 timates the variance of group estimates for small sample sizes, resulting in anti-conservative
706 p-values (Benkeser et al., 2017). We cannot disentangle the possible gains in efficiency from
707 the possibly anti-conservative p-values (due to a finite sample). This limitation would be
708 more of a concern following strict motion QC; in that case, only 29 autistic children were
709 labeled as having usable scans. However, using the lenient motion QC, more than 100 partic-
710 ipants in each diagnosis group had usable scans. In addition, the FDR corrected p-values we
711 use are conservative in the sense that they do not leverage the positive correlations between
712 some edges. An important avenue for future research is to use permutation tests for inference
713 (Winkler et al., 2014) with DRTMLE. Permutation tests can result in finite sample inference
714 while improving power using max statistics, but they create computational challenges.

715 As in many other rs-fMRI studies, we observed extensive variability among participants
716 in the modified partial correlations used as input to DRTMLE (Web Supplement Figure
717 S.3). This variability resulted in generally small effect sizes from the naïve approach. The
718 maximum Cohen’s D across 153 edges was 0.47 at IC02-IC27, which is a medium effect size,
719 and the average naïve effect size among the eleven edges indicated by DRTMLE as significant
720 was 0.32. Unfortunately, calculating effect sizes in DRTMLE is an open problem.

721 Another limitation of the current study is that machine learning algorithms typically re-
722 quire a relatively large sample size compared to classic approaches. We use cross-validation
723 to guard against overfitting, which has been shown to be effective even without having an in-

724 dependent test dataset (Benkeser et al., 2019). One drawback of cross-validation approaches
725 is that they can be sensitive to the random seed. We addressed this limitation by repeating
726 the cross-validation hundreds of times. Each estimation routine takes approximately 6 hours
727 on a single core (2.60 GHz), which includes fitting the propensity model and the outcomes
728 models at the 153 edges. We used a high performance cluster and 100 cores, and conducted
729 two sets of 200 seeds, such that the full estimation routine took approximately 24 hours.
730 The average z-statistic from the two sets were nearly equivalent.

731 *4.3. Model assumptions and possible violations*

732 Estimating the difference in functional connectivity between autistic and typically devel-
733 oping children in the counterfactual world in which all data are usable from the observable
734 data involves three assumptions: mean exchangeability, positivity, and consistency of the
735 counterfactual and the observed outcome (causal consistency) (Section 2.3.1).

736 With respect to mean exchangeability or the assumption of no unmeasured confounders,
737 we assume that functional connectivity is independent of the missingness mechanism given
738 our variables $\{W, A\}$. As noted, the missingness mechanism is deterministic based on head
739 motion, but we are replacing it with a stochastic model that estimates missingness from
740 $\{W, A\}$. In our application, it is important that summary measures of head motion were *not*
741 included in the propensity and outcome models. To understand the reason for this, consider
742 that children who nearly fail motion QC may have some motion impacts in their functional
743 connectivity signal. The deconfounded group difference assumes that Y reflects the signal of
744 interest, i.e., neural sources of variation that are not corrupted by motion. We took several
745 steps to account for potential motion impacts on functional connectivity in children who
746 nearly fail; we used partial correlations from an ICA that includes some motion artifact
747 components (which removes these sources of variance) and residuals from a linear model
748 including motion, as described in Section 2.1.6 and Section 2.3.2, which results in a Y that
749 more closely captures neural sources of variation. However, if we then included summary
750 motion measures in our propensity and outcome models, the propensity model would up-

751 weight these children who nearly failed, and the outcome model, integrating over the full
752 range of head motion, would potentially reintroduce the motion impacts we tried to carefully
753 remove. Additionally, our statistical estimator has the double robustness property: if at
754 least one of the propensity or outcome models is correctly specified, we obtain a statistically
755 consistent estimator of the deconfounded group difference. We include a rich set of predictors
756 and an ensemble of machine learning algorithms, which helps to address the assumption of
757 no unmeasured confounding.

758 Positivity assumes that there are no values of $\{W, A\}$ such that the data will always be
759 unusable. Violations of positivity assumptions lead to out-of-sample prediction of functional
760 connectivity in the outcome model and instabilities in the propensity model, which can
761 lead to greater variance and bias (Petersen et al., 2010). In Fig. 5, we see that for the
762 lenient criteria, the range of the behavioral traits generally overlap between included and
763 excluded participants, although the most severe ADOS score does not appear among the
764 included children. The highest ADOS score among included children was 23; among all
765 children, 26 (A change in the range also occurs for SRS, but SRS was not included in the
766 propensity and outcome models due to a large proportion of missing values.). As reported
767 in Section 2.3.2, all propensities were greater than 0.30 for the first five random seeds. The
768 lack of propensities close to zero for children with usable or unusable data indicates that the
769 assumption of positivity is reasonable in our application. Regarding the last assumption,
770 causal consistency is a technical assumption that assumes that $Y(1)$ is the same as Y when
771 a child has usable data, which in general cannot be tested but seems reasonable.

772 *4.4. Accounting for variables that should be balanced between diagnosis groups*

773 A possible limitation of the current approach is that we account for covariate imbalance
774 between the ASD and typically developing groups using linear regression prior to attempting
775 to account for bias due to data usability, and it may be desirable to pursue a statistical
776 method that integrates covariate balancing into the deconfounded group difference. The
777 deconfounded group difference estimates the marginal mean of each diagnosis group, where

778 integration is across the distribution of the behavioral variables given diagnosis (defined in
779 Section 2.3.1). However, our typically developing sample was aggregated from multiple rs-
780 fMRI studies conducted at KKI, not all of which involved a comparison sample of autistic
781 children. As a result, sex, race, and socioeconomic status significantly differed between
782 diagnosis groups (Table 1) for this secondary analysis. In an ideal prospective experiment
783 using a random sampling design, these socio-demographic variables would not differ. The
784 naïve approach estimated the difference between autistic and typically developing children
785 while controlling for mean FD, max FD, number of frames with $FD < 0.25$ mm, sex, race, and
786 socioeconomic status in a linear model, which is similar to the approach in Di Martino et al.
787 (2014). Controlling for variables in a linear model corresponds to estimating the conditional
788 mean of functional connectivity given these variables. The residuals of the linear model plus
789 the effect of diagnosis are used as input to estimate the deconfounded group mean for each
790 diagnosis group. If there is no sampling bias due to motion exclusion, then the deconfounded
791 group difference is approximately equivalent to the naïve approach, which is a nice aspect of
792 the present study in that it presents a method for evaluating whether confounding is likely to
793 occur when group differences are estimated using the conventional approach. Our approach
794 accounts for possible confounding due to the demographic and remaining motion imbalances
795 in the ASD and typically developing samples, although it does so using the traditional linear
796 model.

797 We can define two sets of variables: 1) variables that we would like to be balanced in
798 autistic and typically developing children in an ideal sample, and 2) variables whose distribu-
799 tion is specific to diagnosis. The target parameter used in estimating an average treatment
800 effect in causal inference marginalizes with respect to the distribution of variables pooled
801 across treatments, which would address biases introduced by the first set of variables. Our
802 deconfounded group difference addresses the second set of variables. Future work could de-
803 fine a target parameter that marginalizes with respect to the desired distribution of variables
804 that should be balanced and the desired distribution of variables whose distribution depends

805 on diagnosis.

806 *4.5. Other methods to account for missingness*

807 An experimental approach to improve the likelihood of collecting usable data from par-
808 ticipants with more severe symptoms is to perform more extensive training in the mock
809 scanner environment with the hope that this will allow children with more severe ASD to
810 complete a scan session. Regarding statistical approaches, we use DRTMLE to estimate the
811 deconfounded group difference. However, a host of other statistical methods could be ap-
812 plied to the same end including covariate matching, propensity score matching (Stuart, 2010;
813 Bridgeford et al., 2021), inverse propensity weighting (Lewinn et al., 2017), G-computation
814 Robins (1986); Snowden et al. (2011), augmented inverse propensity weighting (Robins et al.,
815 2012), and targeted maximum likelihood estimation. Comparing the performance of these
816 approaches in the context of rs-fMRI studies is an important area for future work but is
817 beyond the scope of this paper.

818 *4.6. Significance to other neurological disorders and developmental studies*

819 We used an rs-fMRI study of ASD to illustrate the unintended cost of motion QC on
820 study generalizability, but the issue of data loss and selection bias due to motion QC is
821 neither specific to ASD or to rs-fMRI. Head motion-induced artifacts are a notorious problem
822 for all magnetic resonance-based neuroimaging modalities, and the relationship between
823 motion and participant characteristics is problematic in studies of developmental and aging
824 trajectories, as well as other neurological disorders. For instance, we found that younger
825 children were more likely to be excluded. Recent studies investigating associations between
826 functional brain organization and measures of maturity during the transition from childhood
827 to adolescence have removed large proportions of data (Marek et al., 2019; Dong et al., 2021).
828 Confounding could occur in analyses of rs-fMRI data collected from such developmental
829 samples if the sample of included children that are able to lay motionless tend to be more
830 mature than the full sample. Diffusion MRI and quantitative susceptibility mapping are

831 also susceptible to motion artifacts (Roalf et al., 2016; He et al., 2015), and as a result,
832 studies using these modalities often exclude participants with gross motion. If quality control
833 procedures in studies using these imaging methods result in a reduced sample in which a
834 variable's distribution differs from the original sample, and there is evidence that this variable
835 is related to the outcome of interest, then we recommend adjusting means using DRTMLE.

836 5. Acknowledgments

837 The data analyzed in this study were provided in part by grants awarded to SHM from
838 Autism Speaks, the National Institute of Mental Health (R01 MH078160, R01 MH106564-
839 03), and the National Institute of Neurological Disorders and Stroke (R01 NS048527); to
840 Keri Rosch from the NIMH (K23 MH101322-05); to Karen Seymour from the NIMH (K23
841 MH107734-05); to Bradley Schlagger from the Eunice Kennedy Shriver National Institute
842 of Child Health & Human Development (U54 HD079123); and to Peter van Zijl from the
843 National Institute of Biomedical Imaging and Bioengineering (P54 EB15909). Analysis,
844 interpretation, and writing of the report were supported by a grant from the NIMH (K01
845 MH109766 to MBN).

846 6. Citation diversity statement

847 Recent work in neuroscience (Dworkin et al., 2020) and other fields has identified a
848 citation bias such that papers from women and other minority scholars are under-cited
849 relative to the number of such papers in the field (Mitchell et al., 2013; Dion et al., 2018;
850 Caplar et al., 2017; Maliniak et al., 2013; Bertolero et al., 2020; Wang et al., 2021; Chatterjee
851 and Werner, 2021; Fulvio et al., 2021). We proactively attempted to choose references that
852 reflect the diversity of the neuroscience and statistics fields in the form of contribution,
853 gender, race, and ethnicity. First, we obtained predicted gender of the first and last authors of
854 each reference using databases that store the probability of a name being carried by a woman
855 or a man (Dworkin et al., 2020; Zhou et al., 2020), with possible combinations including

856 male/male, male/female, female/male, and female/female. Our references contain 15.2%
857 woman(first)/woman(last), 13.0% man/woman, 16.7% woman/man, and 55.2% man/man.
858 Relative to the expected proportions in the field of neuroscience, we over- or under-cited these
859 categories by the following ratios: 8.1%, 3.6%, -9.8%, and -3.8%, respectively. Second, we
860 obtained the predicted racial/ethnic category of the first and last author of each reference by
861 databases that store the probability of a first and last name being carried by an author of color
862 (Ambekar et al., 2009; Sood and Laohaprapanon, 2018). Our references contain 8.3% author
863 of color (first)/author of color(last), 12.8% white author/author of color, 16.9% author of
864 color/white author, and 62.0% white author/white author. Self citations for the first and last
865 author of the current paper, as well as references for this diversity statement were excluded
866 from these proportion calculations. These methods are limited by the databases they use
867 for prediction, but we look forward to future work that could help us to better understand
868 how to support equitable practices in science.

869 References

- 870 Allen, E. A., Erhardt, E. B., Damaraju, E., Gruner, W., Segall, J. M., Silva, R. F., Havlicek,
871 M., Rachakonda, S., Fries, J., Kalyanam, R., Michael, A. M., Caprihan, A., Turner, J. A.,
872 Eichele, T., Adelsheim, S., Bryan, A. D., Bustillo, J., Clark, V. P., Ewing, S. W., Filbey,
873 F., Ford, C. C., Hutchison, K., Jung, R. E., Kiehl, K. A., Kodituwakku, P., Komesu,
874 Y. M., Mayer, A. R., Pearlson, G. D., Phillips, J. P., Sadek, J. R., Stevens, M., Teuscher,
875 U., Thoma, R. J., and Calhoun, V. D. (2011). A baseline for the multivariate comparison
876 of resting-state networks. *Frontiers in Systems Neuroscience*, 5.
- 877 Ambekar, A., Ward, C., Mohammed, J., Male, S., and Skiena, S. (2009). Name-ethnicity
878 classification from open sources. In *Proceedings of the 15th ACM SIGKDD international*
879 *conference on Knowledge Discovery and Data Mining*, pages 49–58.
- 880 American Psychiatric Association (2013). *Diagnostic and statistical manual of mental dis-*
881 *orders, 5th ed. (DSM-5®)*. American Psychiatric Association, 5th edition edition.

- 882 Barber, R. F. and Candès, E. J. (2015). Controlling the false discovery rate via knockoffs.
883 *Annals of Statistics*, 43(5):2055–2085.
- 884 Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and
885 powerful approach to multiple testing. *Journal of the Royal Statistical Society, Series B*
886 *(Methodological)*, 57(1):289–300.
- 887 Benkeser, D., Carone, M., van der Laan, M. J., and Gilbert, P. B. (2017). Doubly robust
888 nonparametric inference on the average treatment effect. *Biometrika*, 104(4):863–880.
- 889 Benkeser, D., Petersen, M., and van der Laan, M. J. (2019). Improved Small-Sample Estima-
890 tion of Nonlinear Cross-Validated Prediction Metrics. *Journal of the American Statistical*
891 *Association*, 115(532):1917–1932.
- 892 Bertolero, M. A., Dworkin, J. D., David, S. U., Lloreda, C. L., Srivastava, P., Stiso, J., Zhou,
893 D., Dzirasa, K., Fair, D. A., Kaczkurkin, A. N., Marlin, B. J., Shohamy, D., Uddin, L. Q.,
894 Zurn, P., and Bassett, D. S. (2020). Racial and ethnic imbalance in neuroscience reference
895 lists and intersections with gender. *bioRxiv*.
- 896 Bijsterbosch, J., Harrison, S. J., Jbabdi, S., Woolrich, M., Beckmann, C., Smith, S., and
897 Duff, E. P. (2020). Challenges and future directions for representations of functional brain
898 organization. *Nature Neuroscience*, 23(12):1484–1495.
- 899 Biswal, B., Zerrin Yetkin, F., Haughton, V. M., and Hyde, J. S. (1995). Functional con-
900 nectivity in the motor cortex of resting human brain using echo-planar MRI. *Magnetic*
901 *Resonance in Medicine*, 34(4):537–541.
- 902 Blakemore, S. J., Frith, C. D., and Wolpert, D. M. (2001). The cerebellum is involved in
903 predicting the sensory consequences of action. *Neuroreport*, 12(9):1879–1884.
- 904 Bridgeford, E. W., Powell, M., Kiar, G., Lawrence, R., Caffo, B., Milham, M., and Vogelstein,

- 905 J. T. (2021). Batch Effects are Causal Effects: Applications in Human Connectomics.
906 *bioRxiv*, page 2021.09.03.458920.
- 907 Buckner, R. L., Krienen, F. M., Castellanos, A., Diaz, J. C., and Thomas Yeo, B. T. (2011).
908 The organization of the human cerebellum estimated by intrinsic functional connectivity.
909 *Journal of Neurophysiology*, 106(5):2322–2345.
- 910 Calhoun, V. D., Adali, T., Pearlson, G. D., and Pekar, J. J. (2001). A method for making
911 group inferences from functional MRI data using independent component analysis. *Human*
912 *Brain Mapping*, 14(3):140–151.
- 913 Calhoun, V. D., Wager, T. D., Krishnan, A., Rosch, K. S., Seymour, K. E., Nebel, M. B.,
914 Mostofsky, S. H., Nyalakanai, P., and Kiehl, K. (2017). The impact of T1 versus EPI spatial
915 normalization templates for fMRI data analyses. *Human Brain Mapping*, 38(11):5331–
916 5342.
- 917 Caplar, N., Tacchella, S., and Birrer, S. (2017). Quantitative evaluation of gender bias in
918 astronomical publications from citation counts. *Nature Astronomy*, 1(6):0141.
- 919 Casey, B. J., Cannonier, T., Conley, M. I., Cohen, A. O., Barch, D. M., Heitzeg, M. M.,
920 Soules, M. E., Teslovich, T., Dellarco, D. V., Garavan, H., Orr, C. A., Wager, T. D.,
921 Banich, M. T., Speer, N. K., Sutherland, M. T., Riedel, M. C., Dick, A. S., Bjork, J. M.,
922 Thomas, K. M., Charani, B., Mejia, M. H., Hagler, D. J., Daniela Cornejo, M., Sicut,
923 C. S., Harms, M. P., Dosenbach, N. U. F., Rosenberg, M., Earl, E., Bartsch, H., Watts,
924 R., Polimeni, J. R., Kuperman, J. M., Fair, D. A., and Dale, A. M. (2018). The Ado-
925 lescent Brain Cognitive Development (ABCD) study: Imaging acquisition across 21 sites.
926 *Developmental Cognitive Neuroscience*, 32:43–54.
- 927 Chatterjee, P. and Werner, R. M. (2021). Gender disparity in citations in high-impact journal
928 articles. *JAMA Netw Open*, 4(7):e2114509.

- 929 Ciric, R., Wolf, D. H., Power, J. D., Roalf, D. R., Baum, G. L., Ruparel, K., Shinohara,
930 R. T., Elliott, M. A., Eickhoff, S. B., Davatzikos, C., Gur, R. C., Gur, R. E., Bassett,
931 D. S., and Satterthwaite, T. D. (2017). Benchmarking of participant-level confound re-
932 gression strategies for the control of motion artifact in studies of functional connectivity.
933 *NeuroImage*, 154:174–187.
- 934 Conners, C. K. (1999). Conners Rating Scales-Revised. In *The use of psychological testing for*
935 *treatment planning and outcomes assessment, 2nd ed*, pages 467–495. Lawrence Erlbaum
936 Associates Publishers, Mahwah, NJ, US.
- 937 Conners, C. K. (2008). *Conners 3*. Multi-Health Systems, Inc, Toronto.
- 938 Constantino, J. N. and Gruber, C. P. (2012). *Social Responsiveness Scale Second Edition*
939 *(SRS-2): Manual*. Western Psychological Services (WPS).
- 940 Constantino, J. N. and Todd, R. D. (2003). Autistic traits in the general population: A twin
941 study. *Archives of General Psychiatry*, 60(5):524–530.
- 942 Dajani, D. R. and Uddin, L. Q. (2016). Local brain connectivity across development in
943 autism spectrum disorder: A cross-sectional investigation. *Autism Research*, 9(1):43–54.
- 944 DeBruine, L. (2018). Plot Comparison [blog post]. *retrieved from*
945 <https://debruine.github.io/post/plot-comparison> on 31/10/2021.
- 946 Deen, B. and Pelphrey, K. (2012). Perspective: Brain scans need a rethink. *Nature*, 491(7422
947 SUPPL.):S20–S20.
- 948 Denckla, M. B. (1985). Revised Neurological Examination for Subtle Signs. *Psychopharma-*
949 *cology Bulletin*, 21(4):773–800.
- 950 Di Martino, A., Kelly, C., Grzadzinski, R., Zuo, X. N., Mennes, M., Mairena, M. A., Lord,
951 C., Castellanos, F. X., and Milham, M. P. (2011). Aberrant striatal functional connectivity
952 in children with autism. *Biological Psychiatry*, 69(9):847–856.

- 953 Di Martino, A., O'Connor, D., Chen, B., Alaerts, K., Anderson, J., Assaf, M., Balsters,
954 J., Baxter, L., Beggiato, A., Bernaerts, S., Blanken, L., Bookheimer, S., Braden, B.,
955 Byrge, L., Castellanos, F., Dapretto, M., Delorme, R., Fair, D., Fishman, I., Fitzgerald,
956 J., Gallagher, L., Keehn, R., Kennedy, D., Lainhart, J., Luna, B., Mostofsky, S., Müller,
957 R.-A., Nebel, M., Nigg, J., O'Hearn, K., Solomon, M., Toro, R., Vaidya, C., Wenderoth,
958 N., White, T., Craddock, R., Lord, C., Leventhal, B., and Milham, M. (2017). Enhancing
959 studies of the connectome in autism using the autism brain imaging data exchange II.
960 *Scientific Data*, 4(1):1–15.
- 961 Di Martino, A., Yan, C.-G., Li, Q., Denio, E., Castellanos, F. X., Alaerts, K., Anderson,
962 J. S., Assaf, M., Bookheimer, S., Dapretto, M., Deen, B., Delmonte, S., Dinstein, I., Ertl-
963 Wagner, B., Fair, D., Gallagher, L., Kennedy, D., Keown, C. L., Keyzers, C., Lainhart,
964 J. E., Lord, C., Luna, B., Menon, V., Minshew, N. J., Monk, C., Mueller, S., Müller,
965 R.-A., Nebel, M. B., Nigg, J. T., O'Hearn, K., Pelphrey, K. A., Peltier, S. J., Rudie, J. D.,
966 Sunaert, S., Thioux, M., Tyszka, J. M., Uddin, L. Q., Verhoeven, J. S., Wenderoth, N.,
967 Wiggins, J. L., Mostofsky, S. H., and Milham, M. P. (2014). The autism brain imaging data
968 exchange: Towards a large-scale evaluation of the intrinsic brain architecture in autism.
969 *Molecular Psychiatry*, 19(6):659–667.
- 970 Dion, M. L., Sumner, J. L., and Mitchell, S. M. (2018). Gendered citation patterns across
971 political science and social science methodology fields. *Political Analysis*, 26(3):312–327.
- 972 Dong, H.-M., Margulies, D. S., Zuo, X.-N., and Holmes, A. J. (2021). Shifting gradients
973 of macroscale cortical organization mark the transition from childhood to adolescence.
974 *Proceedings of the National Academy of Sciences*, 118(28).
- 975 Dosenbach, N. U., Koller, J. M., Earl, E. A., Miranda-Dominguez, O., Klein, R. L., Van,
976 A. N., Snyder, A. Z., Nagel, B. J., Nigg, J. T., Nguyen, A. L., Wesevich, V., Greene,
977 D. J., and Fair, D. A. (2017). Real-time motion analytics during brain MRI improve data
978 quality and reduce costs. *NeuroImage*, 161:80–93.

- 979 D’Souza, N. S., Nebel, M. B., Crocetti, D., Robinson, J., Wymbs, N., Mostofsky, S. H.,
980 and Venkataraman, A. (2021). Deep sr-DDL: Deep structurally regularized dynamic dic-
981 tionary learning to integrate multimodal and dynamic functional connectomics data for
982 multidimensional clinical characterizations. *NeuroImage*, 241:118388.
- 983 DuPaul, G. J., Power, T. J., Anastopoulos, A. D., and Reid, R. (1998). *ADHD Rating*
984 *Scale—IV: Checklists, norms, and clinical interpretation*. Guilford Press.
- 985 Dworkin, J. D., Linn, K. A., Teich, E. G., Zurn, P., Shinohara, R. T., and Bassett, D. S.
986 (2020). The extent and drivers of gender imbalance in neuroscience reference lists. *Nature*
987 *neuroscience*, 23(8):918–926.
- 988 Erhardt, E. B., Rachakonda, S., Bedrick, E. J., Allen, E. A., Adali, T., and Calhoun, V. D.
989 (2011). Comparison of multi-subject ICA methods for analysis of fMRI data. *Human*
990 *Brain Mapping*, 32(12):2075–2095.
- 991 Fassbender, C., Mukherjee, P., and Schweitzer, J. B. (2017). Reprint of: Minimizing noise
992 in pediatric task-based functional MRI; Adolescents with developmental disabilities and
993 typical development. *NeuroImage*, 154:230–239.
- 994 Fulvio, J. M., Akinnola, I., and Postle, B. R. (2021). Gender (im)balance in citation practices
995 in cognitive neuroscience. *J Cogn Neurosci*, 33(1):3–7.
- 996 Greene, D. J., Koller, J. M., Hampton, J. M., Wesevich, V., Van, A. N., Nguyen, A. L., Hoyt,
997 C. R., McIntyre, L., Earl, E. A., Klein, R. L., Shimony, J. S., Petersen, S. E., Schlaggar,
998 B. L., Fair, D. A., and Dosenbach, N. U. (2018). Behavioral interventions for reducing
999 head motion during MRI scans in children. *NeuroImage*, 171:234–245.
- 1000 Greenland, S., Robins, J. M., and Pearl, J. (1999). Confounding and collapsibility in causal
1001 inference. *Statistical Science*, 14(1):29–46.

- 1002 Gu, Z., Gu, L., Eils, R., Schlesner, M., and Brors, B. (2014). circlize implements and
1003 enhances circular visualization in R. *Bioinformatics*, 30(19):2811–2812.
- 1004 He, N., Ling, H., Ding, B., Huang, J., Zhang, Y., Zhang, Z., Liu, C., Chen, K., and Yan, F.
1005 (2015). Region-specific disturbed iron distribution in early idiopathic Parkinson’s disease
1006 measured by quantitative susceptibility mapping. *Human Brain Mapping*, 36(11):4407–
1007 4420.
- 1008 Hernan, M. A. and Robins, J. M. (2020). *Causal Inference: What If*. Chapman Hall/CRC,
1009 Boca Raton.
- 1010 Hollingshead, A. B. (1975). Four factor index of social status. *Yale Journal of Sociology*, 8.
- 1011 Hus, V., Gotham, K., and Lord, C. (2014). Standardizing ADOS domain scores: Separating
1012 severity of social affect and restricted and repetitive behaviors. *Journal of Autism and*
1013 *Developmental Disorders*, 44(10):2400–2412.
- 1014 Jenkinson, M., Bannister, P., Brady, M., and Smith, S. (2002). Improved Optimization for
1015 the Robust and Accurate Linear Registration and Motion Correction of Brain Images.
1016 *NeuroImage*, 17:825–841.
- 1017 Johnstone, T., Ores Walsh, K. S., Greischar, L. L., Alexander, A. L., Fox, A. S., Davidson,
1018 R. J., and Oakes, T. R. (2006). Motion correction and the use of motion covariates in
1019 multiple-subject fMRI analysis. *Human Brain Mapping*, 27(10):779–788.
- 1020 Jones, T. B., Bandettini, P. A., Kenworthy, L., Case, L. K., Milleville, S. C., Martin, A.,
1021 and Birn, R. M. (2010). Sources of group differences in functional connectivity: An
1022 investigation applied to autism spectrum disorder. *NeuroImage*, 49(1):401–414.
- 1023 Kaufman, J., Birmaher, B., Axelson, D., Perepletchikova, F., Brent, D., and Ryan, N.
1024 (2013). Schedule for affective disorders and schizophrenia for school-age children-present

- 1025 and lifetime version (K-SADS-PL 2013, DSM-5). *Pittsburgh, PA: Western Psychiatric*
1026 *Institute and Clinic and Yale University.*
- 1027 Kennedy, D. P., Redcay, E., and Courchesne, E. (2006). Failing to deactivate: Resting
1028 functional abnormalities in autism. *Proceedings of the National Academy of Sciences of*
1029 *the United States of America*, 103(21):8275–8280.
- 1030 Keown, C. L., Shih, P., Nair, A., Peterson, N., Mulvey, M. E., and Müller, R. A. (2013).
1031 Local functional overconnectivity in posterior brain regions is associated with symptom
1032 severity in autism spectrum disorders. *Cell Reports*, 5(3):567–572.
- 1033 Kong, X. Z., Zhen, Z., Li, X., Lu, H. H., Wang, R., Liu, L., He, Y., Zang, Y., and Liu, J.
1034 (2014). Individual differences in impulsivity predict head motion during magnetic reso-
1035 nance imaging. *PLoS ONE*, 9(8):e104989.
- 1036 Lake, E. M., Finn, E. S., Noble, S. M., Vanderwal, T., Shen, X., Rosenberg, M. D., Spann,
1037 M. N., Chun, M. M., Scheinost, D., and Constable, R. T. (2019). The functional brain
1038 organization of an individual allows prediction of measures of social abilities transdiag-
1039 nostically in autism and attention-deficit/hyperactivity disorder. *Biological Psychiatry*,
1040 86(4):315–326.
- 1041 Lewinn, K. Z., Sheridan, M. A., Keyes, K. M., Hamilton, A., and McLaughlin, K. A. (2017).
1042 Sample composition alters associations between age and brain structure. *Nature Commu-*
1043 *nications*, 8(1):1–14.
- 1044 Lidstone, D. E., Rochowiak, R., Mostofsky, S. H., and Nebel, M. B. (2021). A Data Driven
1045 Approach Reveals That Anomalous Motor System Connectivity is Associated With the
1046 Severity of Core Autism Symptoms. *Autism Research*, Epub ahead of print:1–18.
- 1047 Lombardo, M. V., Eyler, L., Moore, A., Datko, M., Barnes, C. C., Cha, D., Courchesne,
1048 E., and Pierce, K. (2019). Default mode-visual network hypoconnectivity in an autism
1049 subtype with pronounced social visual engagement difficulties. *eLife*, 8:e47427.

- 1050 Lord, C. and Jones, R. M. (2012). Annual research review: Re-thinking the classification
1051 of autism spectrum disorders. *Journal of Child Psychology and Psychiatry and Allied*
1052 *Disciplines*, 53(5):490–509.
- 1053 Lord, C., Risi, S., Lambrecht, L., Cook, E. H., Leventhal, B. L., Dilavore, P. C., Pickles,
1054 A., and Rutter, M. (2000). The Autism Diagnostic Observation Schedule-Generic: A
1055 standard measure of social and communication deficits associated with the spectrum of
1056 autism. *Journal of Autism and Developmental Disorders*, 30(3):205–223.
- 1057 Maenner, M. J., Shaw, K. A., Bakian, A. V., Bilder, D. A., Durkin, M. S., Esler, A.,
1058 Furnier, S. M., Hallas, L., Hall-Lande, J., Hudson, A., Hughes, M. M., Patrick, M.,
1059 Pierce, K., Poynter, J. N., Salinas, A., Shenouda, J., Vehorn, A., Warren, Z., Constantino,
1060 J. N., DiRienzo, M., Fitzgerald, R. T., Grzybowski, A., Spivey, M. H., Pettygrove, S.,
1061 Zahorodny, W., Ali, A., Andrews, J. G., Baroud, T., Gutierrez, J., Hewitt, A., Lee, L.-C.,
1062 Lopez, M., Mancilla, K. C., McArthur, D., Schwenk, Y. D., Washington, A., Williams, S.,
1063 and Cogswell, M. E. (2021). Prevalence and Characteristics of Autism Spectrum Disorder
1064 Among Children Aged 8 Years — Autism and Developmental Disabilities Monitoring
1065 Network, 11 Sites, United States, 2018. *MMWR. Surveillance Summaries*, 70(11):1–16.
- 1066 Mahadevan, A. S., Tooley, U. A., Bertolero, M. A., Mackey, A. P., and Bassett, D. S.
1067 (2021). Evaluating the sensitivity of functional connectivity measures to motion artifact
1068 in resting-state fMRI data. *NeuroImage*, 241:118408.
- 1069 Maliniak, D., Powers, R., and Walter, B. F. (2013). The gender citation gap in international
1070 relations. *International Organization*, 67(4):889–922.
- 1071 Marek, S., Tervo-Clemmens, B., Nielsen, A. N., Wheelock, M. D., Miller, R. L., Laumann,
1072 T. O., Earl, E., Foran, W. W., Cordova, M., Doyle, O., Perrone, A., Miranda-Dominguez,
1073 O., Feczko, E., Sturgeon, D., Graham, A., Hermsillo, R., Snider, K., Galassi, A., Nagel,
1074 B. J., Ewing, S. W., Eggebrecht, A. T., Garavan, H., Dale, A. M., Greene, D. J., Barch,

- 1075 D. M., Fair, D. A., Luna, B., and Dosenbach, N. U. (2019). Identifying reproducible indi-
1076 vidual differences in childhood functional brain networks: An ABCD study. *Developmental*
1077 *Cognitive Neuroscience*, 40:100706.
- 1078 Mayes, S. D. and Calhoun, S. L. (2008). WISC-IV and WIAT-II profiles in children with
1079 high-functioning autism. *Journal of Autism and Developmental Disorders*, 38(3):428–439.
- 1080 Mejia, A., Nebel, M., Barber, A., Choe, A., Pekar, J., Caffo, B., and Lindquist, M. (2018).
1081 Improved estimation of subject-level functional connectivity using full and partial corre-
1082 lation with empirical Bayes shrinkage. *NeuroImage*, 172:478–491.
- 1083 Mejia, A. F., Nebel, M. B., Eloyan, A., Caffo, B., and Lindquist, M. A. (2017). PCA
1084 leverage: outlier detection for high-dimensional functional magnetic resonance imaging
1085 data. *Biostatistics*, 18(3):521–536.
- 1086 Mitchell, S. M., Lange, S., and Brus, H. (2013). Gendered citation patterns in international
1087 relations journals. *International Studies Perspectives*, 14(4):485–492.
- 1088 Moore, K. L. and van der Laan, M. J. (2009). Covariate adjustment in randomized trials
1089 with binary outcomes: Targeted maximum likelihood estimation. *Statistics in Medicine*,
1090 28(1):39–64.
- 1091 Mostofsky, S. H., Newschaffer, C. J., and Denckla, M. B. (2003). Overflow movements predict
1092 impaired response inhibition in children with ADHD. *Perceptual and Motor Skills*, 97(3
1093 II):1315–1331.
- 1094 Mowinckel, A. (2018). Circular Plots in R and Adding Images [blog post]. *retrieved*
1095 *from* <https://drmwinkels.io/blog/2018-05-25-circular-plots-in-r-and-adding-images> on
1096 *31/10/2021*.
- 1097 Muschelli, J., Nebel, M., Caffo, B., Barber, A., Pekar, J., and Mostofsky, S. (2014). Reduction
1098 of motion-related artifacts in resting state fMRI using aCompCor. *NeuroImage*, 96:22–35.

- 1099 Nielsen, A. N., Greene, D. J., Gratton, C., Dosenbach, N. U., Petersen, S. E., and Schlaggar,
1100 B. L. (2019). Evaluating the prediction of brain maturity from functional connectivity
1101 after motion artifact denoising. *Cerebral Cortex*, 29(6):2455–2469.
- 1102 Padmanabhan, A., Lynch, C. J., Schaer, M., and Menon, V. (2017). The Default Mode
1103 Network in Autism. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*,
1104 2(6):476–486.
- 1105 Parkes, L., Fulcher, B., Yücel, M., and Fornito, A. (2018). An evaluation of the efficacy,
1106 reliability, and sensitivity of motion correction strategies for resting-state functional MRI.
1107 *NeuroImage*, 171:415–436.
- 1108 Petersen, M. L., Porter, K. E., Gruber, S., Wang, Y., and Laan, M. J. v. d. (2010). Diagnosing
1109 and responding to violations in the positivity assumption:. *Statistical Methods in Medical
1110 Research*, 21(1):31–54.
- 1111 Polley, E., LeDell, E., Kennedy, C., and van der Laan, M. (2019). SuperLearner: Super
1112 Learner Prediction. *R package v. 2.0-26*.
- 1113 Power, J. D. (2017). A simple but useful way to assess fMRI scan qualities. *NeuroImage*,
1114 154:150–158.
- 1115 Power, J. D., Barnes, K. A., Snyder, A. Z., Schlaggar, B. L., and Petersen, S. E. (2012).
1116 Spurious but systematic correlations in functional connectivity MRI networks arise from
1117 subject motion. *NeuroImage*, 59(3):2142–2154.
- 1118 Power, J. D., Lynch, C. J., Adeyemo, B., and Petersen, S. E. (2020). A critical, event-related
1119 appraisal of denoising in resting-state fMRI studies. *Cerebral Cortex*, 30(10):5544–5559.
- 1120 Power, J. D., Mitra, A., Laumann, T. O., Snyder, A. Z., Schlaggar, B. L., and Petersen,
1121 S. E. (2014). Methods to detect, characterize, and remove motion artifact in resting state
1122 fMRI. *NeuroImage*, 84:320–341.

- 1123 Pruijn, R. H. R., Mennes, M., Buitelaar, J. K., and Beckmann, C. F. (2015). Evaluation of
1124 ICA-AROMA and alternative strategies for motion artifact removal in resting state fMRI.
1125 *NeuroImage*, 112:278–287.
- 1126 Redcay, E., Moran, J. M., Mavros, P. L., Tager-Flusberg, H., Gabrieli, J. D., and Whitfield-
1127 Gabrieli, S. (2013). Intrinsic functional network organization in high-functioning adoles-
1128 cents with autism spectrum disorder. *Frontiers in Human Neuroscience*, 7:573.
- 1129 Reich, W. (2000). Diagnostic Interview for Children and Adolescents (DICA). *Journal of*
1130 *the American Academy of Child and Adolescent Psychiatry*, 39(1):59–66.
- 1131 Roalf, D. R., Quarmley, M., Elliott, M. A., Satterthwaite, T. D., Vandekar, S. N., Ruparel,
1132 K., Gennatas, E. D., Calkins, M. E., Moore, T. M., Hopson, R., Prabhakaran, K., Jackson,
1133 C. T., Verma, R., Hakonarson, H., Gur, R. C., and Gur, R. E. (2016). The impact
1134 of quality assurance assessment on diffusion tensor imaging outcomes in a large-scale
1135 population-based cohort. *NeuroImage*, 125:903–919.
- 1136 Robins, J. (1986). A new approach to causal inference in mortality studies with a sustained
1137 exposure period-application to control of the healthy worker survivor effect. *Mathematical*
1138 *Modelling*, 7(9-12):1393–1512.
- 1139 Robins, J. M., Rotnitzky, A., and Zhao, L. P. (2012). Estimation of Re-
1140 gression Coefficients When Some Regressors are not Always Observed.
1141 <https://doi.org/10.1080/01621459.1994.10476818>, 89(427):846–866.
- 1142 Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3):581–592.
- 1143 Satterthwaite, T. D., Elliott, M. A., Gerraty, R. T., Ruparel, K., Loughhead, J., Calkins,
1144 M. E., Eickhoff, S. B., Hakonarson, H., Gur, R. C., Gur, R. E., and Wolf, D. H. (2013).
1145 An improved framework for confound regression and filtering for control of motion artifact
1146 in the preprocessing of resting-state functional connectivity data. *NeuroImage*, 64:240–256.

- 1147 Satterthwaite, T. D., Wolf, D. H., Loughead, J., Ruparel, K., Elliott, M. A., Hakonarson,
1148 H., Gur, R. C., and Gur, R. E. (2012). Impact of in-scanner head motion on multiple
1149 measures of functional connectivity: Relevance for studies of neurodevelopment in youth.
1150 *NeuroImage*, 60(1):623–632.
- 1151 Simhal, A. K., Filho, J. O., Segura, P., Cloud, J., Petkova, E., Gallagher, R., Castellanos,
1152 F. X., Colcombe, S., Milham, M. P., and Di Martino, A. (2021). Predicting multiscan
1153 MRI outcomes in children with neurodevelopmental conditions following MRI simulator
1154 training. *Developmental Cognitive Neuroscience*, 52:101009.
- 1155 Simonoff, E., Pickles, A., Charman, T., Chandler, S., Loucas, T., and Baird, G. (2008).
1156 Psychiatric disorders in children with autism spectrum disorders: Prevalence, comorbidity,
1157 and associated factors in a population-derived sample. *Journal of the American Academy*
1158 *of Child and Adolescent Psychiatry*, 47(8):921–929.
- 1159 Snowden, J. M., Rose, S., and Mortimer, K. M. (2011). Implementation of G-Computation
1160 on a Simulated Data Set: Demonstration of a Causal Inference Technique. *American*
1161 *Journal of Epidemiology*, 173(7):731–738.
- 1162 Sood, G. and Laohaprapanon, S. (2018). Predicting race and ethnicity from the sequence of
1163 characters in a name. *arXiv preprint arXiv:1805.02109*.
- 1164 Stoodley, C. J., D’Mello, A. M., Ellegood, J., Jakkamsetti, V., Liu, P., Nebel, M. B., Gibson,
1165 J. M., Kelly, E., Meng, F., Cano, C. A., Pascual, J. M., Mostofsky, S. H., Lerch, J. P.,
1166 and Tsai, P. T. (2017). Altered cerebellar connectivity in autism and cerebellar-mediated
1167 rescue of autism-related behaviors in mice. *Nature Neuroscience*, 20(12):1744–1751.
- 1168 Stuart, E. A. (2010). Matching methods for causal inference: A review and a look forward.
1169 *Statistical science : a review journal of the Institute of Mathematical Statistics*, 25(1):1–21.
- 1170 Supekar, K., Uddin, L. Q., Khouzam, A., Phillips, J., Gaillard, W. D., Kenworthy, L. E.,

- 1171 Yerys, B. E., Vaidya, C. J., and Menon, V. (2013). Brain Hyperconnectivity in Children
1172 with Autism and its Links to Social Deficits. *Cell Reports*, 5(3):738–747.
- 1173 Tyszka, J. M., Kennedy, D. P., Paul, L. K., and Adolphs, R. (2014). Largely typical patterns
1174 of resting-state functional connectivity in high-functioning adults with autism. *Cerebral*
1175 *Cortex*, 24(7):1894–1905.
- 1176 Uddin, L. Q., Supekar, K., Lynch, C. J., Khouzam, A., Phillips, J., Feinstein, C., Ryali, S.,
1177 and Menon, V. (2013). Salience network-based classification and prediction of symptom
1178 severity in children with autism. *JAMA Psychiatry*, 70(8):869–879.
- 1179 Uddin, L. Q., Yeo, B. T., and Spreng, R. N. (2019). Towards a Universal Taxonomy of
1180 Macro-scale Functional Human Brain Networks. *Brain Topography*, 32(6):926–942.
- 1181 van der Laan, M. J., Polley, E. C., and Hubbard, A. E. (2007). Super learner. *Statistical*
1182 *Applications in Genetics and Molecular Biology*, 6(1).
- 1183 van der Laan, M. J. and Robins, J. M. (2003). *Unified Methods for Censored Longitudinal*
1184 *Data and Causality*. Springer Science & Business Media.
- 1185 van der Laan, M. J. and Rose, S. (2011). *Targeted Learning: Causal Inference for Observa-*
1186 *tional and Experimental Data*. Springer, New York, NY.
- 1187 van Dijk, K. R., Sabuncu, M. R., and Buckner, R. L. (2012). The influence of head motion
1188 on intrinsic functional connectivity MRI. *NeuroImage*, 59(1):431–438.
- 1189 Van Overwalle, F., Van de Steen, F., van Dun, K., and Heleven, E. (2020). Connectiv-
1190 ity between the cerebrum and cerebellum during social and non-social sequencing using
1191 dynamic causal modelling. *NeuroImage*, 206:116326.
- 1192 Vander Weele, T. J. (2009). Concerning the consistency assumption in causal inference.
1193 *Epidemiology*, 20(6):880–883.

- 1194 Vasa, R. A., Mostofsky, S. H., and Ewen, J. B. (2016). The Disrupted Connectivity Hy-
1195 pothesis of Autism Spectrum Disorders: Time for the Next Phase in Research. *Biological*
1196 *Psychiatry: Cognitive Neuroscience and Neuroimaging*, 1(3):245–252.
- 1197 Wang, X., Dworkin, J., Zhou, D., Stiso, J., Falk, E. B., Bassett, D., Zurn, P., and Lydon-
1198 Staley, D. M. (2021). Gendered citation practices in the field of communication. *Annals*
1199 *of the International Communication Association*, 45(2):134–153.
- 1200 Wechsler, D. (2003). *Wechsler Intelligence Scale for Children-WISC-IV*. Psychological Cor-
1201 poration.
- 1202 Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New
1203 York.
- 1204 Winkler, A. M., Ridgway, G. R., Webster, M. A., Smith, S. M., and Nichols, T. E. (2014).
1205 Permutation inference for the general linear model. *NeuroImage*, 92(100):381–397.
- 1206 Wood, S. N. (2017). *Generalized additive models: an introduction with R*. CRC Press.
- 1207 Wymbs, N. F., Nebel, M. B., Ewen, J. B., and Mostofsky, S. H. (2021). Altered inferior
1208 parietal functional connectivity is correlated with praxis and social skill performance in
1209 children with autism spectrum disorder. *Cerebral Cortex*, 31(5):2639–2652.
- 1210 Yan, C. G., Cheung, B., Kelly, C., Colcombe, S., Craddock, R. C., Di Martino, A., Li, Q.,
1211 Zuo, X. N., Castellanos, F. X., and Milham, M. P. (2013). A comprehensive assessment
1212 of regional variation in the impact of head micromovements on functional connectomics.
1213 *NeuroImage*, 76:183–201.
- 1214 Zhou, D., Cornblath, E. J., Stiso, J., Teich, E. G., Dworkin, J. D., Blevins, A. S., and
1215 Bassett, D. S. (2020). Gender diversity statement and code notebook v1.0 [software].
1216 retrieved from <https://github.com/dalejn/cleanBib> on 30/11/2021.